

CARNEGIE MELLON UNIVERSITY

DOCTORAL THESIS

Mathematical and Algorithmic Analysis
of Network and Biological Data

Author:

CHARALAMPOS
E.TSOURAKAKIS

Supervisor:

PROF. ALAN FRIEZE

Thesis Committee

PROF. ALAN FRIEZE (CHAIR)

PROF. TOM BOHMAN

PROF. R. RAVI

PROF. RUSSELL SCHWARTZ

*A thesis submitted in fulfilment of the requirements
for the degree of Doctor of Philosophy*

in the

ALGORITHMS, COMBINATORICS AND OPTIMIZATION
DEPARTMENT OF MATHEMATICAL SCIENCES

April 10, 2013

“I don’t keep checks and balances, I don’t seek to adjust myself. I follow the deep throbbing of my heart.”

Nikos Kazantzakis

Keywords: Data Science, Graphs and Networks, Algorithm design, Random graphs, Computational Cancer Biology, Applied Mathematics

CARNEGIE MELLON UNIVERSITY

DEPARTMENT OF MATHEMATICAL SCIENCES

Abstract

Doctor of Philosophy

Mathematical and Algorithmic Analysis of Network and Biological Data

by CHARALAMPOS E. TSOURAKAKIS

This dissertation contributes to mathematical and algorithmic problems that arise in the analysis of network and biological data. The driving force behind this dissertation is the importance of studying network and biological data. Studies of such data can provide us with insights to important emerging and long-standing problems, such as *why do some mutations cause cancer whereas others do not? What is the structure of the Web graph? How do networks form and how does their structure affect the the spread of ideas and diseases?* This thesis consists of two parts. The first part is devoted to *graphs and networks* and the second part to *computational cancer biology*. Our contributions to *graphs and networks* revolve around the following two axes.

- **Empirical studies:** In order to develop a good model, one has to study structural properties of real-world networks. Given the size of today's networks, such empirical studies and graph-structured computations become challenging tasks. Therefore, in the course of empirically studying real-world networks, both novel and existing well-studied problems but in new computational models need to be solved. We provide both efficient algorithms with strong theoretical guarantees for various graph-structured computations and graph processing systems that allow us to handle big graph data.
- **Models:** The quality of a given model is judged by how well it matches reality. We analyze fundamental graph theoretic properties of random Apollonian networks, a random graph model which mimicks real-world networks. Furthermore, we use random graphs to perform an average case analysis for rainbow connectivity, an intriguing connectivity concept. We provide simple randomized procedures which succeed to solve the problem of interest with high probability on random binomial and regular graphs.

Our contributions to *computational cancer biology* include new models, theoretical insights into existing models, novel algorithmic techniques and detailed experimental analysis of various datasets. Specifically, we contribute the following.

- Novel algorithmic techniques for denoising array comparative genomic hybridization data. Our algorithmic results are of independent interest and provide approximation techniques for speeding up dynamic programming.
- Based on empirical findings which strongly indicate an inherent geometric structure in cancer genomic data, we introduce a geometric model for finding subtypes of cancer which overcomes difficulties of existing methods such as principal/independent component analysis and separation methods for Gaussians. We provide a computational method which solves the optimization problem efficiently and is robust to outliers.
- We find the necessary and sufficient conditions to reconstruct uniquely an oncogenetic tree, a popular tumor phylogenetic method.

Acknowledgements

First, I wish to express my deep gratitude to my advisor Professor Alan M. Frieze. I feel privileged to have worked under his supervision. His insights, advice and support shall remain invaluable to me. His dedication to noble research and open mind are a source of inspiration.

Also, I thank Professor Russell Schwartz and Professor Gary L. Miller. Russell introduced me to computational biology and provided me with insights on how to perform research on this human-centric and interdisciplinary domain. The meetings I have had with Gary have been the best teaching sessions on algorithm design I have had. My interaction with both of them has been a great pleasure.

A key figure during my Ph.D. years was Professor Mihail N. Kolountzakis. I want to thank him for his beneficial impact on my research, our collaboration and for hosting me during the summer of 2010 in the University of Crete. Furthermore, I would like to thank Professor Aristides Gionis for a fruitful collaboration that began in Yahoo! Research and will continue in Aalto University for the next academic semester. I also thank Professor Christos Faloutsos for advising me during the first two years of my Ph.D. in the School of Computer Science.

I would like to thank the members of my thesis committee for their valuable feedback: Professors Tom Bohman, Alan Frieze, Russell Schwartz and R. Ravi.

I would like to thank my collaborators throughout my Ph.D. years in the School of Computer Science and the Department of Mathematical Sciences: Petros Drineas, Aristides Gionis, Christos Faloutsos, Alan Frieze, U Kang, Mihail N. Kolountzakis, Ioannis Koutis, Jure Leskovec, Eirinaios Michelakis, Gary L. Miller, Rasmus Pagh, Richard Peng, Aditya Prakash, Russell Schwartz, Stanley Shackney, Ayshwarya Subramanian, Jimeng Sun, David Tolliver, Maria Tsiarli, Hanghang Tong.

I want to thank all my friends for their support and the good times we have had over the last years. I would like to especially thank Felipe Trevizan, Jamie Flanick, Ioannis Mallios, Evangelos Vlachos, Ioannis Koutis, Don Sheehy, Shay Cohen, Eirinaios Michelakis, Nikos Papadakis, Dimitris Tsaparas and Dimitris Tsagkarakis.

I thank my parents Eftichios and Aliko and my sister Maria for their love and support. Finally, I want to thank my fiancée Maria Tsiarli.

Contents

Acknowledgements	v
Abstract	iii
List of Figures	xi
List of Tables	xviii
1 Introduction	2
1.1 Graphs and Networks	2
1.1.1 Structural properties of real-world networks	4
1.1.1.1 Power laws	4
1.1.1.2 Small-world	5
1.1.1.3 Clustering coefficients	5
1.1.1.4 Communities	6
1.1.1.5 Densification and Shrinking Diameters	7
1.1.1.6 Web graph	8
1.1.2 Network Models	9
1.1.2.1 Erdős-Rényi random graphs	9
1.1.2.2 Configuration model and Random Regular Graphs	10
1.1.2.3 Preferential attachment	11
1.1.2.4 Kronecker graphs	12
1.1.3 Triangles	13
1.1.3.1 Triangle counting methods	15
1.1.4 Dense Subgraphs	18
1.1.5 Graph Partitioning	20
1.1.6 Big Graph Data Analytics	22
1.1.6.1 MAPREDUCE Basics	22
1.2 Computational Cancer Biology	23
1.2.1 Denoising array-based Comparative Genomic Hybridization (aCGH) data	24
1.2.2 Cancer subtypes	26
1.2.3 Oncogenetic trees	28
1.3 Thesis Overview	29
2 Theoretical Preliminaries	38
2.1 Concentration of Measure Inequalities	38
2.2 Random Projections	40
2.3 Extremal Graph Theory	40

2.4	Two useful lemmas	40
2.5	Semidefinite Bounds	41
2.6	Speeding up Dynamic Programming	42
2.7	Reporting Points in a Halfspace	43
2.8	Monge Functions and Dynamic Programming	44
3	Rainbow Connectivity of Sparse Random Graphs	46
3.1	Introduction	46
3.2	Sketch of approach	48
3.3	Proof of Theorem 3.1	48
3.4	Proof of Theorem 3.2	56
3.4.1	Tree building	56
3.4.2	Coloring the edges	59
3.4.3	Case 1: $r \geq 4$:	59
4	Random Apollonian networks	65
4.1	Model & Main Results	65
4.2	Related Work	67
4.3	Proof of Theorem 4.1	68
4.4	Proof of Theorem 4.3	71
4.5	Proof of Theorem 4.4	80
4.6	Proof of Theorem 4.5	83
5	Triangle Counting in Large Graphs	86
5.1	Introduction	86
5.2	Triangle Sparsifiers	87
5.2.1	Proposed Algorithm	87
5.2.2	Experimental Results	91
5.2.2.1	Datasets	91
5.2.2.2	Experimental Setup	91
5.2.2.3	Experimental Results	92
5.2.2.4	The “Doubling” Algorithm	93
5.2.3	Theoretical Ramifications	94
5.3	Efficient Triangle Counting in Large Graphs via Degree-based Vertex Partitioning	95
5.3.1	Proposed Method	95
5.3.1.1	Edge Sparsification	96
5.3.1.2	Triple Sampling	97
5.3.1.3	Hybrid algorithm	98
5.3.1.4	Sampling in the Semi-Streaming Model	99
5.3.2	Experiments	100
5.3.2.1	Data	100
5.3.2.2	Experimental Setup and Implementation Details	100
5.3.2.3	Binomial Sampling in Expected Sublinear time	101
5.3.2.4	Results	101
5.3.2.5	Remarks	102
5.3.3	Theoretical Ramifications	103

5.3.3.1	Random Projections and Triangles	104
5.4	Colorful Triangle Counting	106
5.4.1	Algorithm	106
5.4.2	Analysis	107
5.4.2.1	Second Moment Method	107
5.4.2.2	Concentration via the Hajnal-Szemerédi Theorem	109
5.4.2.3	Complexity	110
5.4.2.4	Discussion	110
5.4.3	A MAPREDUCE Implementation	111
6	Dense Subgraphs	113
6.1	Introduction	113
6.2	A General Framework	114
6.3	Optimal Quasi-cliques	116
6.4	Problem variants	124
6.4.1	Top- k optimal quasi-cliques	124
6.4.2	Constrained optimal quasi-cliques	124
6.5	Experimental evaluation	125
6.5.1	Real-world graphs	126
6.5.2	Synthetic graphs	128
6.6	Applications	131
6.6.1	Thematic groups	132
6.6.2	Correlated genes	133
7	Structure of the Web Graph	135
7.1	Introduction	135
7.2	Related Work	136
7.3	Distinct Elements and the Diameter	137
7.4	HADOOP Algorithms	139
7.4.1	HADI-naive in MAPREDUCE	139
7.4.2	HADI-plain in MAPREDUCE	141
7.4.3	HADI-optimized in MAPREDUCE	143
7.4.4	Analysis	146
7.4.5	HADI in SQL	146
7.5	Wall-clock Times	146
7.5.1	Experimental Setup	147
7.5.2	Running Time and Scale-up	148
7.5.3	Effect of Optimizations	149
7.6	Structure of Real-World Networks	149
7.6.1	Static Patterns	150
7.6.1.1	Diameter	150
7.6.1.2	Shape of Distribution	151
7.6.1.3	Radius plot of GCC	154
7.6.1.4	“Core” and “Whisker” vertices	154
7.6.2	Temporal Patterns	154
8	FENNEL: Streaming Graph Partitioning for Massive Scale Graphs	156

8.1	Introduction	156
8.2	Proposed Framework	157
8.3	One-Pass Streaming Algorithm	167
8.4	Experimental Evaluation	168
8.4.1	Experimental Setup	168
8.4.2	Synthetic Datasets	172
8.4.3	Real-World Datasets	173
8.5	System Evaluation	175
8.6	Discussion	177
9	PEGASUS: A System for Large-Scale Graph Processing	179
9.1	Introduction	179
9.2	Proposed Method	180
9.2.1	GIM-V and PageRank	181
9.2.2	GIM-V and Random Walk with Restart	182
9.2.3	GIM-V and Diameter Estimation	182
9.2.4	GIM-V and Connected Components	183
9.3	HADOOP Implementation	183
9.3.1	GIM-V BASE: Naive Multiplication	184
9.3.2	GIM-V BL: Block Multiplication	185
9.3.3	GIM-V CL: Clustered Edges	186
9.3.4	GIM-V DI: Diagonal Block Iteration	187
9.3.5	GIM-V NR: Node Renumbering	187
9.3.6	Analysis	188
9.4	Scalability	189
9.4.1	Results	189
9.5	Pegasus at Work	191
9.5.1	Connected Components of Real Networks	192
9.5.2	PageRank scores of Real Networks	195
9.5.3	Diameter of Real Networks	195
10	Approximation Algorithms for Speeding up Dynamic Programming and Denoising aCGH data	200
10.1	Problem & Formulation	200
10.2	$O(n^2)$ Dynamic Programming Algorithm	201
10.3	Analysis of The Transition Function	203
10.4	Additive Approximation using Halfspace Queries	204
10.5	Multiscale Monge Decomposition	208
10.6	Validation of Our Model	212
10.6.1	Experimental Setup and Datasets	212
	How to pick C ?	213
10.6.2	Experimental Results and Biological Analysis	214
	10.6.2.1 Synthetic Data	214
	10.6.2.2 Coriell Cell Lines	215
	10.6.2.3 Breast Cancer Cell Lines	217
	Cell Line BT474:	217
	Cell Line HS578T:	219

Cell Line T47D:	220
Cell Line MCF10A	222
11 Robust Unmixing of Tumor States in Array Comparative Genomic Hybridization Data	223
11.1 Introduction	223
11.2 Approach	224
11.2.1 Algorithms and Assumptions	224
11.2.1.1 Hard Geometric Unmixing	227
11.2.1.2 Soft Geometric Unmixing	227
11.2.1.3 Analysis & Efficiency	228
11.3 Experimental Methods	230
11.3.1 Methods: Synthetic Experiments	230
11.3.1.1 Mixture Sampler	230
11.3.1.2 Geometric Sampling of End-members & Noise	231
11.3.1.3 Evaluation	232
11.4 Results	233
11.4.1 Results: Synthetic Data	233
11.4.2 Array Comparative Gene Hybridization (aCGH) Data	233
11.4.2.1 Preprocessing	234
11.4.2.2 Unmixing Analysis and Validation	235
12 Perfect Reconstruction of Oncogenetic Trees	240
12.1 Introduction	240
12.2 Proof of Theorem 12.1	241
13 Conclusion	244
13.1 Open Problems	244
13.1.1 Rainbow Connectivity (Chapter 3)	244
13.1.2 Random Apollonian Graphs (Chapter 4)	244
13.1.3 Triangle Counting (Chapter 5)	245
13.1.4 Densest Subgraphs (Chapter 6)	246
13.1.5 Structure of the Web Graph (Chapter 7)	246
13.1.6 FENNEL: Streaming Graph Partitioning for Massive Scale Graphs (Chapter 8)	246
13.1.7 PEGASUS: A System for Large-Scale Graph Processing (Chapter 9)	247
13.1.8 Approximation Algorithms for Speeding up Dynamic Programming and Denoising aCGH data (Chapter 10)	247
13.1.9 Robust Unmixing of Tumor States in Array Comparative Genomic Hybridization Data (Chapter 11)	248
13.1.10 Perfect Reconstruction of Oncogenetic Trees (Chapter 12)	248
Bibliography	250

List of Figures

1.1	Bow-tie structure of the Web graph (Image source: [90])	8
1.2	Schematic representation of array CGH. Genomic DNA from two cell populations (1) is differentially labeled and hybridized in a microarray (2). Typically the reference DNA comes from a normal subject. For humans this means that the reference DNA comes from a normal diploid genome. The ratios on each spot are measured and normalised so that the median \log_2 ratio is zero. The final result is an ordered tuple containing values of the fluorescent ratios in each genomic position per each chromosome. This is shown in (3) where we see the genomic profile of the cell line GM05296 [365]. The problem of denoising array CGH data is to infer the true DNA copy number T per genomic position from a set of noisy measurements of the quantity $\log_2 \frac{T}{R}$, where $R=2$ for normal diploid humans.	25
3.1	Structure of Lemma 3.6.	50
3.2	Subgraph found in the proof of Lemma 3.6.	51
3.3	Figure shows $\frac{\log n}{101}$ -ary trees T_x, T_y . The two roots are shown respectively at the center of the trees. In our thinking of the random coloring as an evolutionary process, the green edges incident to x survive with probability 1, the red edges incident to y with probability $1 - \frac{1}{q}$ and all the other edges with probability $p_0 = \left(1 - \frac{2k}{q}\right)^2$ where k is the depth of both trees and q the number of available colors. Our analysis in Lemma 3.6 using these probabilities gives a lower bound on the number of alive pairs of leaves after coloring T_x, T_y from the root to the leaves respectively.	53
3.4	Taking care of small vertices.	55
4.1	Snapshots of a Random Apollonian Network (RAN) at: (a) $t = 1$ (b) $t = 2$ (c) $t = 3$ (d) $t = 100$	66
4.2	Coupling used in Lemma 4.11.	73
4.3	An instance of the process for $t = 2$. Each face is labeled with its depth.	84
4.4	RANs as random ternary trees.	85
4.5	The height of the random ternary tree cannot be used to lower bound the diameter. The height of the random ternary tree can be arbitrarily large but the diameter is 2.	85
5.1	Graph with linear number $O(n)$ of triangles.	95
5.2	Conditions of Theorem 5.5 are tight. In order to hope for concentration p has to be greater than (a) $\frac{\Delta}{t}$ and (b) $t^{-1/2}$	109

5.3	Consider the indicator variable X_i corresponding to the i -th triangle. Note that $\Pr[X_i \text{rest are monochromatic}] = p \neq \Pr[X_i] = p^2$. The indicator variables are <i>pairwise</i> but not mutually independent.	111
6.1	Edge density and diameter of the top-10 subgraphs found by our GREEDY-OQC and LOCALSEARCHOQC methods, and Charikar’s algorithm, on the AS-skitter graph (top) and the Wikipedia 2006/11 graph (bottom). . .	129
6.2	Precision and recall for our method and Goldberg’s algorithm vs. the power-law exponent of the host graph.	131
6.3	Authors returned by our LOCALSEARCHOQC algorithm when queried with Papadimitriou and Abiteboul. The set includes well-known database scientists. The induced subgraph has 34 vertices and 457 edges. The edge density is 0.81, the diameter is 3, the triangle density is 0.66.	132
6.4	Authors returned by our LOCALSEARCHOQC algorithm when queried with Papadimitriou and Blum. The set includes well-known theoretical computer scientists. The induced subgraph has 13 vertices and 38 edges. The edge density is 0.49, the diameter is 3, the triangle density is 0.14. . .	132
6.5	Genes returned by our LOCALSEARCHOQC algorithm when queried with p53. The induced subgraph is a clique with 14 vertices.	133
6.6	A tumorigenesis pathway consistent with our findings.	133
7.1	One iteration of HADI-naive. <i>Stage 1</i> . Bitstrings of all vertex are sent to every reducer. <i>Stage 2</i> . Sums up the count of changed nodes.	140
7.2	One iteration of HADI-plain. <i>Stage 1</i> . Edges and bitstrings are matched to create partial bitstrings. <i>Stage 2</i> . Partial bitstrings are merged to create full bitstrings. <i>Stage 3</i> . Sums up the count of changed nodes, and compute $N(h)$, the neighborhood function. Computing $N(h)$ is not drawn in the figure for clarity.	142
7.3	Converting the original edge and bitstring to blocks. The 4-by-4 edge and length-4 bitstring are converted to 2-by-2 super-elements and length-2 super-bitstrings. Notice the lower-left super-element of the edge is not produced since there is no nonzero element inside it.	145
7.4	Running time versus number of edges with HADI-OPT on Kronecker graphs for three iterations. Notice the excellent scalability: linear on the graph size (number of edges).	148
7.5	“Scale-up” (throughput $1/T_M$) versus number of machines M , for the Kronecker graph (2B edges). Notice the near-linear growth in the beginning, close to the ideal(dotted line).	148
7.6	Run time of HADI with/without optimizations for Kronecker and Erdős-Rényi graphs with several billions of edges, on the M45 HADOOP cluster using 90 machines for 3 iterations. HADI-OPT is up to $7.6\times$ faster than HADI-plain.	149
7.7	(a) Radius plot (Count versus Radius) of the YahooWeb graph. Notice the effective diameter is surprisingly small. Also notice the peak (marked ‘S’) at radius 2, due to star-structured disconnected components. (b) Radius plot of GCC (Giant Connected Component) of YahooWeb graph. The <i>only</i> vertex with radius 5 (marked ‘C’) is <code>google.com</code>	150

7.8	Average diameter vs. number of vertices in lin-log scale for the three different Web graphs, where M and B stand for millions and billions respectively. (0.3M): web pages inside nd.edu at 1999, from Albert et al.'s work. (203M): web pages crawled by Altavista at 1999, from Broder et al.'s work (1.4B): web pages crawled by Yahoo at 2002 (YahooWeb in Table 7.1). Notice the relatively small diameters for both the directed and the undirected cases.	151
7.9	Static Radius Plot (Count versus Radius) of U.S. Patent and LinkedIn graphs. Notice the bimodal structure with 'outsiders' (vertices in the DCs), 'core' (central vertices in the GCC), and 'whiskers' (vertices connected to the GCC with long paths).	151
7.10	Radius plot (Count versus radius) for several connected components of the U.S. Patent data in 1985. In blue: the distribution for the giant connected component; rest colors: several disconnected components.	152
7.11	Size distribution of connected components. Notice the size of the disconnected components (DCs) follows a power-law which explains the first peak around radius 0 of the radius plots in Figure 7.9.	153
7.12	Radius-Degree plots of real-world graphs. HD represents the vertex with the highest degree. Notice that HD belongs to core vertices inside the GCC, and whiskers have small degree.	153
7.13	Evolution of the effective diameter of real graphs. The diameter increases until a 'gelling' point, and starts to decrease after the point.	154
7.14	Radius distribution over time. "Expansion": the radius distribution moves to the right until the gelling point. "Contraction": the radius distribution moves to the left after the gelling point.	155
8.1	Fraction of edges cut λ and maximum load normalized ρ as a function of γ , ranging from 1 to 4 with a step of 0.25, over five randomly generated power law graphs with slope 2.5. The straight lines show the performance of METIS.	168
8.2	(Left) CDF of the relative difference $(\frac{\lambda_{\text{Fennel}} - \lambda_c}{\lambda_c}) \times 100\%$ of percentages of edges cut of our method and the best competitor. (Right) Same but for the absolute difference $(\lambda_{\text{Fennel}} - \lambda_c) \times 100\%$	173
8.3	Absolute difference $\delta\lambda$ and Relative gain versus the maximum load imbalance $\delta\rho$	175
8.4	Fennel: time vs. number of edges.	176
9.1	GIM-V BL using 2 x 2 blocks. $B_{i,j}$ represents a matrix block, and v_i represents a vector block. The matrix and vector are joined block-wise, not element-wise.	186
9.2	Clustered vs. non-clustered adjacency matrices for two isomorphic graphs. The edges are grouped into 2 by 2 blocks. The left graph uses only 3 blocks while the right graph uses 9 blocks.	186
9.3	Propagation of component id(=1) when block width is 4. Each element in the adjacency matrix of (a) represents a 4 by 4 block; each column in (b) and (c) represents the vector after each iteration. GIM-V DL finishes in 4 iterations while GIM-V BL requires 8 iterations.	187

9.4	Scalability and Performance of GIM-V. (a) Running time decreases quickly as more machines are added. (b) The performance(=1/running time) of 'BL-CL' wins more than 5x (for n=3 machines) over the 'BASE'. (c) Every version of GIM-V shows linear scalability.	190
9.5	Comparison of GIM-V DI and GIM-V BL-CL for HCC. GIM-V DI finishes in 6 iterations while GIM-V BL-CL finishes in 18 iterations due to long chains.	191
9.6	Degree distribution of LinkedIn. Notice that the original minimum vertex has degree 1, which is highly probable given the power-law behavior of the degree distribution. After the renumbering, the minimum vertex is replaced with a highest-degree node.	192
9.7	Number of iterations vs. the minimum vertex of LinkedIn, for connected components. D_i represents the vertex with i -th largest degree. Notice that the number of iterations decreased by 19 % after renumbering.	192
9.8	Degree distribution of Wikipedia. Notice that the original minimum vertex has degree 1, as in LinkedIn. After the renumbering, the minimum vertex is replaced with a highest-degree node.	193
9.9	Number of iterations vs. the minimum vertex of Wikipedia, for connected components. D_i represents the vertex with i -th largest degree. Notice that the number of iterations decreased by 17 % after renumbering.	193
9.10	The evolution of connected components. (a) The giant connected component grows for each year. However, the second largest connected component do not grow above Dunbar's number(≈ 150) and the slope of the size distribution remains constant after the gelling point at year 2003. As in LinkedIn, notice the growth of giant connected component and the constant slope of the size distribution.	194
9.11	Connected Components of YahooWeb. Notice the two anomalous spikes which are far from the constant-slope line. Most of them are domain selling or porn sites which are replicated from templates.	195
9.12	PageRank distribution of YahooWeb. The distribution follows a power law with an exponent 2.30.	196
9.13	PageRank distribution of WWW-Barabasi. The distribution follows a power law with an exponent 2.25.	196
9.14	The evolution of PageRanks.(a) The distributions of PageRanks follows a power-law. However, the exponent at year 2003, which is around the gelling point, is much different from year 2004, which are after the gelling point. The exponent increases after the gelling point and becomes stable. Also notice the maximum PageRank after the gelling point is about 10 times larger than that before the gelling point due to the emergence of the giant connected component. (b) Again, the distributions of PageRanks follows a power-law. Since the gelling point is before year 2005, the three plots show similar characteristics: the maximum PageRanks and the slopes are similar.	197
9.15	Radius of real graphs. X axis: radius in linear scale. Y axis: number of vertices in log scale. (Row 1) LinkedIn from 2003 to 2006. (Row 2) Wikipedia from 2005 to 2007. (Row 3) DBLP, flickr, Epinion. Notice that all the radius plots have the bimodal structure due to many smaller connected components(first mode) and the giant connected component(second mode).	198

10.1	Answering whether or not $\bar{D}P_i \leq x + C$ reduces to answering whether the point set $\{(j, \bar{D}P_j, 2S_j, S_j^2 + \bar{D}P_j j) \in \mathbb{R}^4, j < i\}$ has a non-empty intersection with the halfspace $\gamma = \{y \in \mathbb{R}^4 : a_i y \leq b_i\}$ where a_i and b_i are a 4-dimensional constant vector and a constant which depend on i respectively. This type of queries can be solved efficiently, see [15].	205
10.2	Solving the k -th Monge problem, $k = 1, \dots, m = O(\log_{1+\epsilon}(i))$, for a fixed index i . The k -th interval is the set of indices $\{j : l_k \leq j \leq r_k, l_k = i - (1+\epsilon)^k, r_k = i - (1+\epsilon)^{k-1}\}$. Ideally, we wish to define a Monge function w'_k whose maximum inside the k -th interval (red color) is smaller than the minimum outside that interval (green color). This ensures that the optimal breakpoint for the k -th Monge subproblem lies inside the red interval.	211
10.3	ROC curve of CGHTRIMMER as a function of C on data from [415]. The red arrow indicates the point (0.91 and 0.98 recall and precision respectively) corresponding to $C=0.2$, the value used in all subsequent results.	214
10.4	Performance of CGHTRIMMER, CBS, and CGHSEG on denoising synthetic aCGH data from [263]. CGHTRIMMER and CGHSEG exhibit excellent precision and recall whereas CBS misses two consecutive genomic positions with DNA copy number equal to 3.	214
10.5	Visualization of the segmentation output of CGHTRIMMER, CBS, and CGHSEG for the cell line BT474 on chromosomes 1 (a,b,c), 5 (d,e,f), and 17 (g,h,i). (a,d,g) CGHTRIMMER output. (b,e,h) CBS output. (c,f,i) CGHSEG output. Segments exceeding the ± 0.3 threshold [61] are highlighted.	218
10.6	Visualization of the segmentation output of CGHTRIMMER, CBS, and CGHSEG for the cell line HS578T on chromosomes 3 (a,b,c), 11 (d,e,f), and 17 (g,h,i). (a,d,g) CGHTRIMMER output. (b,e,h) CBS output. (c,f,i) CGHSEG output. Segments exceeding the ± 0.3 threshold [61] are highlighted.	219
10.7	Visualization of the segmentation output of CGHTRIMMER, CBS, and CGHSEG for the cell line T47D on chromosomes 1 (a,b,c), 11 (d,e,f), and 20 (g,h,i). (a,d,g) CGHTRIMMER output. (b,e,h) CBS output. (c,f,i) CGHSEG output. Segments exceeding the ± 0.3 threshold [61] are highlighted.	220
10.8	Visualization of the segmentation output of CGHTRIMMER, CBS, and CGHSEG for the cell line MCF10A on chromosomes 1 (a,b,c), 8 (d,e,f), and 20 (g,h,i). (a,d,g) CGHTRIMMER output. (b,e,h) CBS output. (c,f,i) CGHSEG output. Segments exceeding the ± 0.3 threshold [61] are highlighted.	221

- 11.1 *Left:* The minimum area fit of a simplex containing the sample points in the plane (shown in black) using the program in § 11.2.1.1. On noiseless data, hard geometric unmixing recovers the locations of the fundamental components at the vertices. *Right:* However, the containment simplex is highly sensitive to noise and outliers in the data. A single outlier, circled above, radically changes the shape of the containment simplex fit (light gray above). In turn, this changes the estimates of basis distributions used to unmix the data. We mitigate this short coming by developing a soft geometric unmixing model (see § 11.2.1.2) that is comparatively robust to noise. The soft fit (shown dark gray) is geometrically very close to the generating sources as seen on the left. 224
- 11.2 An illustration of the reduced coordinates under the unmixing hypothesis: points (shown in gray) sampled from the 3-simplex embedded in \mathbb{R}^3 and then perturbed by log-normal noise, producing points shown in black with sample correspondence given the green arrows. Note that the dominant subspace remains in the planar variation induced by the simplex, and a 2D reduced representation for simplex fitting is thus sufficient. 226
- 11.3 An example sample set generated for §11.3.1.2 shown in the “intrinsic dimensions” of the model. Note that sample points cleave to the lower dimensional substructure (edges) of the simplex. 231
- 11.4 *Left:* mean squared error for the component reconstruction comparing Hard Geometric Unmixing (MVES: [102]) and Soft Geometric Unmixing (SGU) introduced in §11.2.1.2 for the experiment described in §11.3.1.2 with variable γ . The plot demonstrates that robust unmixing more accurately reconstructs the ground truth centers relative to hard unmixing in the presence of noise. *Right:* mean squared error for mixture reconstruction comparing MVES and SGU. 232
- 11.5 Empirical motivation for the $\ell_1 - \ell_1$ -total variation functional for smoothing CGH data. The left plot shows the histogram of values found in the CGH data obtained from the [312] data set. The distribution is well fit by the high kurtosis Laplacian distribution in lieu of a Gaussian. The right plot shows the distribution of differences along the probe array values. As with the values distribution, these frequencies exhibit high kurtosis. . . 234
- 11.6 The simplex fit to the CGH data samples from [312] ductal data set in \mathbb{R}^3 . The gray tetrahedron was returned by the optimization of Program 11.1 and the green tetrahedron was returned by the robust unmixing routine. . 234
- 11.7 Inferred mixture fractions for six-component soft geometric unmixing applied to breast cancer aCGH data. Data is grouped by tumor, with multiple sectors per tumor placed side-by-side. Columns are annotated below by sector or N for normal control and above by cell sorting fraction (D for diploid, H for hypodiploid, A for aneuploid, and A1/A2 for subsets of aneuploid) where cell sorting was used. 237
- 11.8 Copy numbers of inferred components versus genomic position. The average of all input arrays (top) is shown for comparison, with the six components below. Benchmarks loci are indicated by yellow vertical bars. . . 238

11.9	Plot of amplification per probe highlighting regions of shared amplification across components. The lower (blue) dots mark the location of the collected cancer benchmarks set. Bars highlight specific markers of high shared amplification for discussion in the text. <i>Above:</i> A: 1q21 (site of MUC1), B: 9p21 (site of CDKN2B), C: 7q21 (site of HER2), D: 17q12 (site of PGAP3), E: 5q21 (site of APC/MCC).	239
12.1	Illustration of necessity conditions of Theorem 12.1.	241
13.1	By the pigeonhole principle, one of the three initial faces receives $\Theta(t)$ vertices. Using Theorem 4.3 it is not hard to see that the encircled set of vertices S has conductance $\phi(S) \approx \frac{\sqrt{t}}{t} = \frac{1}{\sqrt{t}}$ <i>whp.</i>	245
13.2	Weighted Graph	245
13.3	Lemma 13.1: if $ P_{i_1} - P_{i_2} > 2\sqrt{2C}$ then Segmentation (b) which consists of five segments (two of which are single points) is superior to Segmentation (a) which consists of a single segment.	248

List of Tables

5.1	Dataset sources.	91
5.2	Datasets used in our experiments. Abbreviations are included. Symbol \odot stands for Autonomous Systems graphs, \star for online social networks and \diamond for Web graphs. Notice that the networks with the highest triangle counts are online social networks (Flickr, Livejournal, Orkut), verifying the folklore that online social networks are abundant in triangles.	92
5.3	Results of experiments averaged over 5 trials using $p = 0.1$. All running times are reported in seconds. The first column shows the running time for the exact counting algorithm averaged over 5 runs. The second and third column show the error and running time averaged over 5 runs for each dataset (two decimal digits of accuracy). Standard deviations are also included (three decimal digits of accuracy). The last column shows the running time averaged over 5 runs for the 1-pass algorithm as stated in [96, §2.2] and the corresponding standard deviations. The number of samples for each dataset was set to a value that achieves <i>at most</i> as good accuracy as the ones achieved by our counting method. See Section 5.2.2.3 for all the details.	93
5.4	Doubling procedure for the Wikipedia 2005 graph with 44,667,095 triangles.	94
5.5	Datasets used in our experiments.	100
5.6	Results of experiments averaged over 5 Trials using only triple sampling.	102
5.7	Results of experiments averaged over 5 trials using sparsification and triple sampling.	103
5.8	Values for the variables involved in our formulae for five real-world graphs. Typically, Δ and t_{\max} are significantly less than the obvious upper bounds $n-2$ and $\binom{n-1}{2}$ respectively. Furthermore, $3\Delta t$ is significantly larger than $\sum_{e \in E(G)} \delta_e^2$	112
6.1	Difference between densest subgraph and optimal quasi-clique on some popular graphs. $\delta = e[S]/\binom{ S }{2}$ is the edge density of the extracted subgraph, D is the diameter, and $\tau = t[S]/\binom{ S }{3}$ is the triangle density.	114
6.2	Graphs used in our experiments.	126
6.3	Densest subgraphs extracted with Charikar’s method vs. optimal quasi-cliques extracted with the proposed GREEDYOQC algorithm (GREEDY) and LOCALSEARCHOQC algorithm (LS). $\delta = e[S]/\binom{ S }{2}$ is the edge density of the extracted subgraph S , D is the diameter, and $\tau = t[S]/\binom{ S }{3}$ is the triangle density.	128

6.4	Subgraphs returned by Goldberg’s max-flow algorithm and by our two algorithms (GREEDYOQC, LOCALSEARCHOQC) on Erdős-Rényi graphs with 3 000 vertices and three values of p , and with a planted clique of 30 vertices.	130
7.1	Datasets (B: Billion, M: Million, K: Thousand, G: Gigabytes)	147
8.1	Fraction of edges cut λ and the normalized maximum load ρ for Fennel, the previously best-known heuristic (linear weighted degrees [373]) and hash partitioning of vertices for the Twitter graph with approximately 1.5 billion edges. Fennel and best competitor require around 40 minutes, METIS more than $8\frac{1}{2}$ hours.	156
8.2	Datasets used in our experiments.	169
8.3	Performance of various existing methods on amazon0312, k is set to 32.	170
8.4	Fraction of edges cut λ and normalized maximum load ρ for Fennel and METIS [245] averaged over 5 random graphs generated according to the HP(5000,0.8,0.5) model.	172
8.5	The relative gain $(1 - \frac{\lambda_{\text{Fennel}}}{\lambda_c}) \times 100\%$ and load imbalance, where subindex c stands for the best competitor, averaged over all datasets in Table 9.2 as a function of k	174
8.6	The average duration of a step and the average amount of MB exchanged per node and per step during the execution of PageRank on LiveJournal data set.	176
9.1	GIM-V in terms of SQL.	181
9.2	Order and size of networks.	189
10.1	Summary of proof of Lemma 10.5.	208
10.2	Datasets, papers and the URLs where the datasets can be downloaded. ⊙ and ■ denote which datasets are synthetic and real respectively.	213
10.3	Results from applying CGHTRIMMER, CBS, and CGHSEG to 15 cell lines. Rows with listed chromosome numbers (e.g., GM03563/3) corresponded to known gains or losses and are annotated with a check mark if the expected gain or loss was detected or a “No” if it was not. Additional rows list chromosomes on which segments not annotated in the benchmark were detected; we presume these to be false positives.	216
11.1	Benchmark set of breast cancer markers selected for validation of real data, annotated by gene name, genomic locus, and the set of components exhibiting amplification at the given marker.	237

Dedicated to:
my fiancée Maria Tsiarli
my parents Alik and Eftichios and my sister Maria,
Mr. Andreas Varverakis

Chapter 1

Introduction

The motivating force behind this dissertation is the importance of studying network and biological data. Such studies can provide us with insights or even answers to various significant questions: *How do people establish connections among each other and how does the underlying social graph affect the spread of ideas and diseases? What will the structure of the Web be in some years from now? How can we design better marketing strategies? Why do some mutations cause cancer whereas others do not? Can we use cancer data to improve diagnostics and therapeutics?*

This dissertation contributes to mathematical and algorithmic problems which arise in the course of studying *network* and *cancer* data. This Chapter is organized as follows: in Sections 1.1 and 1.2 we introduce the reader to some basic network and biology concepts respectively. In Section 1.3 we motivate our work and present the contributions of this dissertation.

1.1 Graphs and Networks

Networks appear throughout in nature and society [21]. Furthermore, many applications which use other types of data, such as text or image data, create graphs by data-processing. *Network Science* has emerged over the last years as an interdisciplinary area spanning traditional domains including mathematics, computer science, sociology, biology and economics. Since complexity in social, biological and economical systems, and more generally in complex systems, arises through pairwise interactions there exists a surging interest in understanding networks [144]. Graph theory [76] plays a key role in modeling and studying networks. Given the increasing importance of networks in our lives, Professor Daniel Spielman believes that Graph Theory is the new Calculus [369]. Some important networks and the corresponding graph models follow.

- *Human brain.* Our brain consists of neurons which form a network. The number of neurons in the region of human cortex is estimated to be roughly 10^{10} . Neurons are connected through synapses. The strength of synapses in general varies. The average number of synapses of each neuron is in the range of 24 000-80 000 for humans [401]. The brain can be modeled as a graph, whose vertex set corresponds to neurons and edge set to synapses.
- *World Wide Web (WWW).* The World Wide Web is a particularly important network. Turing Award winner Jim Gray in his 1998 Turing award address speech mentioned “The emergence of ‘cyberspace’ and the World Wide Web is like the discovery of a new continent”. The vertices of the WWW are Web pages and the edges are hyperlinks (URLs) that point from one page to another. In 2011, Google reported that WWW has more than a trillion of edges.
- *Internet.* There exist two different types of Internet graphs. In the first type, vertices are routers and edges correspond to physical connections. The second level is the autonomous systems level, where a single vertex represents a domain, namely multiple routers and computers. An edge is drawn if there is at least one route that connects them. It is worth mentioning that the first type of topology can be studied with the *traceroute* tool, whereas the second with *BGP tables*.
- *Social and online social networks.* According to Aristotle man is by nature a social animal. Social networks precede the individual. These networks are modeled by using one vertex per human. Each edge corresponds to some sort of interaction, for instance friendship.

Since the emergence of the World Wide Web, we live in the information age. Nowadays, online social networks and social media are a part of our daily lives which can affect society immensely. Again, such networks are modeled as graphs where each vertex corresponds to an account. Each edge corresponds to a connection.

- *Collaboration networks.* There exist many types of collaboration networks. Vertices can represent for instance mathematicians or actors. An edge between two vertices is drawn when the corresponding mathematicians have co-authored a paper or when the corresponding actors have played in the same movie respectively. Two famous collaboration networks are the Erdős and the Kevin Bacon collaboration networks.
- *Protein interaction networks.* Protein - protein interactions occur when two or more proteins bind together, often to carry out their biological function. These interactions define protein interaction networks. The corresponding graph has a vertex for each protein. Proteins i and j are connected if they interact.
- *Wireless sensor networks.* Vertices are autonomous sensors. These sensors monitor physical or environmental conditions, such as temperature, brightness and humidity

A directed edge (i, j) suggests that information can be passed from sensor i to sensor j .

- *Power grid.* An electrical grid is a network which delivers electricity from suppliers to consumers. The vertices of the power grid graph correspond to generators, transformers and substations. Edges correspond to high-voltage transmission lines.
- *Financial networks.* One of the globalization effects is the interdependence between financial institutions and markets. Financial network models are represented by graphs whose vertex set corresponds to financial institutions such as banks and countries. Edges correspond to different types of interactions, for instance borrower-lender relations.
- *Blog networks.* A blog is a web log. A blog may contain a blogroll, namely a list of blogs that interest the blogger, or may reference other blogs through posts. This network is represented by a graph whose vertex set corresponds to blogs. An edge (i, j) exists if the i -th blog has blog j in its blogroll or a URL pointing to j .

A remarkable fact which underpins network science is that many real-world networks, despite their different origins, share several common structural characteristics. In the next Section we present some established patterns, shared by numerous real-world networks [144, 180].

1.1.1 Structural properties of real-world networks

Real-world networks possess a variety of remarkable properties. For instance, it is known that they are robust to random but vulnerable to malicious attacks [21, 78]. Here, we review some important empirical properties of real-world networks which are related to our work.

1.1.1.1 Power laws

Real-world networks typically have skewed degree distributions. These heavy tailed empirical distributions are frequently modeled as power laws.

Definition 1.1. The degree sequence of a graph follows a power law distribution if the number of vertices N_d with degree d is given by $N_d \propto d^{-\alpha}$ where $\alpha > 1$ is the power law degree exponent or slope.

One of the early papers that popularized power laws as the modeling choice for empirical degree distributions was the Faloutsos, Faloutsos and Faloutsos paper [158]. The

Faloutsos brothers found that the Internet at the autonomous systems level follows a power law degree distribution with $\alpha \approx 2.4$. In general, there exist three main problems with the initial studies of power laws in networks. First, Internet graphs generated with traceroute sampling [158] may produce power-law distributions due to the bias of the process, even if the true underlying graph is regular [264]. Secondly, there exist methodological flaws in determining the exponent/slope of the power law distribution. Clauset et al. provide a proper methodology for finding the slope of the distribution [118]. Finally, other distributions could potentially fit the data better but were not considered as candidates. Such distribution is the lognormal [369]. A nice review of power law and lognormal distributions appears in [305].

1.1.1.2 Small-world

Six degrees of separation is the theory that anyone in the world is no more than six relationships away from any other person. In the early 20th century Nobel Peace Prize winner Guglielmo Marconi, the father of modern radio, suggested that it would take only six relay stations to cover and connect the earth by radio [298]. It is likely that this idea was the seed for the six degrees of separation theory, which was further supported by Frigyes Karinthy in a short story called Chains. Since then many scientists, including Michael Gurevich, Ithiel De Sola Pool have worked on this theory. In a famous experiment, Stanley Milgram asked people to route a postcard to a fixed recipient by passing them to direct acquaintances [303]. Milgram observed that depending on the sample of people chosen the average number of intermediaries was between 4.4 and 5.7. Milgram's experiment besides its existential aspect has a strong algorithmic aspect as well, which was first studied by Kleinberg [252].

Nowadays, World Wide Web and online social networks provide us with data that reach the planetary scale. Recently, Backstrom, Boldi, Rosa, Ugander and Vigna showed that the world is even smaller than what the six degrees of separation theory predicts [45]. Specifically, they perform the first world-scale social-network graph-distance computation, using the entire Facebook network of active users (at that time 721 million users, 69 billion friendship links) and observe an average distance of 4.74. In Chapter 7 we shall see formal graph theoretical concepts which quantify the small world phenomenon.

1.1.1.3 Clustering coefficients

Watts and Strogatz in their influential paper [412] proposed a simple graph model which reproduces two features of social networks: abundance of triangles and the existence of short paths among any pair of nodes. Their model combines the idea of homophily which

leads to the wealth of triangles in the network and the idea of weak ties which create short paths. In order to quantify the homophily, they introduce the definitions of the clustering coefficient. The definition of the transitivity $T(G)$ of a graph G , introduced by Newman et al. [317], is closely related to the clustering coefficient and measures the probability that two neighbors of a vertex are connected.

Definition 1.2 (Clustering Coefficient). A vertex $v \in V(G)$ with degree $d(v)$ which participates into $t(v)$ triangles has clustering coefficient $C(v)$ equal to the fraction of edges among its neighbors to the maximum number of triangles it could participate:

$$C(v) = \frac{t(v)}{\binom{d(v)}{2}} \quad (1.1)$$

The clustering coefficient $C(G)$ of graph G is the average of $C(v)$ over all $v \in V(G)$.

Definition 1.3 (Transitivity). The transitivity of a graph measures the probability that two neighbors of a vertex are connected:

$$T(G) = \frac{3 \times t}{\text{number of connected triples}} \quad (1.2)$$

1.1.1.4 Communities

Intuitively, communities are sets of vertices which are densely intra-connected and sparsely inter-connected [196, 314]. A large amount of research in network science has focused on finding communities. The goal of community detection methods is to partition the graph vertices into communities so that there are many edges among vertices in the same community and few edges among vertices in different communities. A landmark study of communities by Leskovec et al. [277] placed various folklores revolving around community existence in question. Specifically, Leskovec et al. observed in the majority of the networks they studied that communities exist at small size scales. Specifically, as the size increases up to a value which empirically is close to 100 the quality of the community tends to increase. However, after the critical value of 100, the quality tends to decrease. This indicates that communities blend in with the rest of the network and lose their community-like profile.

There exist various formalizations of the community notion [277]. However, what underpins all these formalizations is the attempt to understand the cut structure of the graph. A popular measure for the quality of a community is the conductance.

Definition 1.4. Given a graph $G(V, E, w)$ and a set $S \subseteq V$ of vertices, the conductance of S is defined as

$$\phi(S) = \frac{\sum_{(i,j) \in E, i \in S, j \in \bar{S}} w_{ij}}{\min(\text{vol}(S), \text{vol}(\bar{S}))},$$

where $\bar{S} = V \setminus S$ and the volume of a given set $A \subseteq V$ of vertices is defined as $\text{vol}(A) = \sum_{i \in A} d(i)$. The conductance of the graph G is defined as

$$\phi = \min_{S \subseteq V} \phi(S).$$

It is not a coincidence that random walks are frequently used to find communities [343] as their use is common in the general setting of graph partitioning [322]. A lot of interest exists into finding dense sets of vertices around a given seed. A popular method for finding such sets was first introduced by Lovász and Simonovits [290] who show that random walks of length $O(\frac{1}{\phi})$ can be used to compute a cut with sparsity at most $\tilde{O}(\sqrt{\phi})$ if the sparsest cut has conductance ϕ . Later, Spielman and Teng [371, 372] provided a local graph partitioning algorithm which implements efficiently the Lovász-Simonovits idea. Furthermore, their algorithm has a bounded work/volume ratio. Another closely related approach which does not explicitly compute sequences of random walk distributions but computes a personalized Pagerank vector was introduced by Andersen, Chung, and Lang [33]. A few other representative approaches for the problem of community detection include methods on minimum cut [170], modularity maximization [196], and spectral methods [245, 318]. In fact the literature on the topic is so extensive that we do not attempt to make a proper review here; a comprehensive survey has been conducted by Fortunato [174].

In the typical setting of finding communities, the vertex set of the graph is partitioned. A relaxation of the latter requirement, allowing overlaps between sets of vertices, yields the notion of overlapping communities [39, 383, 428].

1.1.1.5 Densification and Shrinking Diameters

Leskovec, Kleinberg and Faloutsos [280] studied how numerous real-world networks from a variety of domains evolve over time. In their work two important observations were made. First, networks become denser over time and the densification follows a power law pattern. Secondly, effective diameters shrink over time. The second pattern is particularly surprising and creates a modeling challenge as well, since the vast majority of real-world networks have a diameter that grows as the network grows.

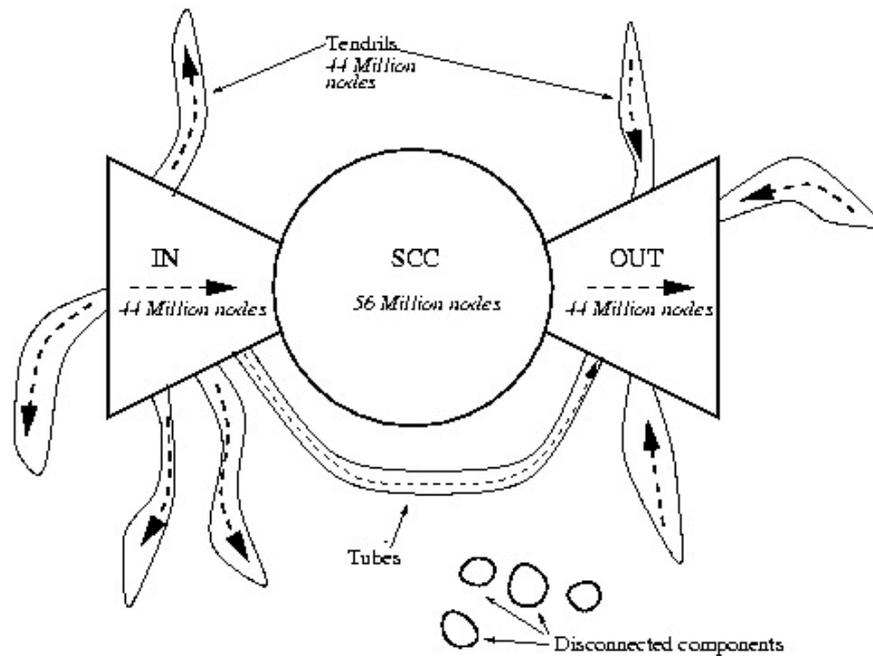


FIGURE 1.1: Bow-tie structure of the Web graph (Image source: [90])

1.1.1.6 Web graph

We focus on the bow *bow-tie structure* of the Web graph. Other important properties of the Web graph include the abundance of bipartite cliques [253, 262] and compressibility [73, 109, 110]. In 1999 Andrei Broder et al. [90] performed an influential study of the Web graph using strongly connected components (SCCs) as their building blocks. Specifically, they proposed the bow-tie model for the structure of the Web graph based on their findings on the index of pages and links of the AltaVista search engine. According to the bow-tie structure of the Web, there exists a single giant SCC. Broder et al. [90] positioned the remaining SCCs with respect to the giant SCC as follows:

- IN: vertices that can reach the giant SCC but cannot be reached from it.
- OUT: vertices that can be reached from the giant SCC but cannot be reach it.
- Tendrils: These are vertices that either are reachable from IN but cannot reach the giant SCC or the vertices that can reach OUT but cannot be reached from the giant SCC.
- Disconnected: vertices that belong to none of the above categories. These are the vertices that even if we ignore the direction of the edges, have no path connecting them to the giant SCC.

A schematic picture of the bow-tie structure of the Web is shown in Figure 1.1. This structure has been verified in other studies as well [69, 137].

1.1.2 Network Models

There exist two main research lines which model network formation. Random graph and strategic models. We focus on random graph models and specifically the models used in this dissertation. The interested reader may consult the cited papers and references therein for more random graph models and [221] for a rich account of results on strategic models.

1.1.2.1 Erdős-Rényi random graphs

Let \mathcal{G} be the family of all labeled graphs with vertex set $V = [n] = \{1, \dots, n\}$. Notice $|\mathcal{G}| = 2^{\binom{n}{2}}$. Random graph models assign to each graph $G \in \mathcal{G}$ a probability. The random binomial graph model $G(n, p)$ has two parameters, n the number of vertices and a probability parameter $0 \leq p \leq 1$. The $G(n, p)$ model assigns to a graph $G \in \mathcal{G}$ the following probability

$$\Pr[G] = p^{|E(G)|} (1-p)^{\binom{n}{2}-|E(G)|}.$$

We will refer to random binomial graphs as Erdős-Rényi graphs interchangeably. Historically, Gilbert [194] introduced originally the $G(n, p)$ model but Erdős and Rényi founded the field of random graphs [152, 153]. They introduced a closely related model known as $G(n, m)$. This model has two parameters, the number of vertices n and the number of edges m , where $0 \leq m \leq \binom{n}{2}$. This model assigns equal probability to all labelled graphs on the vertex set $[n]$ with exactly m edges. In other words,

$$\Pr[G] = \begin{cases} \frac{1}{\binom{\binom{n}{2}}{m}} & \text{if } |E(G)| = m \\ 0 & \text{if } |E(G)| \neq m \end{cases}$$

We shall be interested in understanding various graph theoretic properties.

Definition 1.5. Define a graph property \mathcal{P} as a subset of all possible labelled graphs. Namely $\mathcal{P} \subseteq 2^{\binom{[n]}{2}}$.

For instance \mathcal{P} can be the set of planar graphs or the set of graphs that contain a Hamiltonian cycle. We will call a property \mathcal{P} as monotone increasing if $G \in \mathcal{P}$ implies

$G + e \in \mathcal{P}$. For instance the Hamiltonian property is monotone increasing. Similarly, we will call a property \mathcal{P} as monotone decreasing if $G \in \mathcal{P}$ implies $G - e \in \mathcal{P}$. For instance the planarity property is monotone decreasing. Since there is an underlying probability distribution, we shall be interested in how likely a property is. Our claims relating to random graphs will be probabilistic. We will say that an event A_n holds with high probability (*whp*) if $\lim_{n \rightarrow +\infty} \Pr[A_n] = 1$. The following are background definitions which can be found in any random graph theory textbook [77, 223].

Notice that in the $G(n, p)$ model we toss a coin independently for each possible edge and with probability p we add it to the graph. In expectation there will be $p \binom{n}{2}$ edges. When $p = \frac{m}{\binom{n}{2}}$, then a random binomial graph has in expectation m edges and intuitively $G(n, p)$ and $G(n, m)$ should behave similarly. The following theorem quantifies this intuition.

Theorem 1.6. *Let $0 \leq p_0 \leq 1$, $s(n) = n\sqrt{p(1-p)} \rightarrow +\infty$, and $\omega(n) \rightarrow +\infty$ as $n \rightarrow +\infty$. Then,*

(a) *if \mathcal{P} is any graph property and for all $m \in \mathbb{N}$ such that $|m - \binom{n}{2}p| < \omega(n)s(n)$, the probability $\Pr[G(n, m) \in \mathcal{P}] \rightarrow p_0$, then $\Pr[G(n, p) \in \mathcal{P}] \rightarrow p_0$ as $n \rightarrow +\infty$.*

(b) *if \mathcal{P} is a monotone graph property and $p_- = p_0 - \frac{\omega ns(n)}{n^3}$, $p_+ = p_0 + \frac{\omega ns(n)}{n^3}$ then from the facts that $\Pr[G(n, p_-) \in \mathcal{P}] \rightarrow p_0$, $\Pr[G(n, p_+) \in \mathcal{P}] \rightarrow p_0$, it follows that $\Pr[G(n, p \binom{n}{2}) \in \mathcal{P}] \rightarrow p_0$ as $n \rightarrow +\infty$.*

1.1.2.2 Configuration model and Random Regular Graphs

The configuration model can construct a multigraph in general with a given degree sequence $d = (d_1, \dots, d_n)$. We describe the configuration for random regular graphs [419], as we use it in Chapter 3. We follow the configuration model of Bollobás [77] in our proofs, see [223] for further details. Let $W = [2m = rn]$ be our set of *configuration points* and let $W_i = [(i-1)r + 1, ir]$, $i \in [n]$, partition W . The function $\phi : W \rightarrow [n]$ is defined by $w \in W_{\phi(w)}$. Given a pairing F (i.e. a partition of W into m pairs) we obtain a (multi-)graph G_F with vertex set $[n]$ and an edge $(\phi(u), \phi(v))$ for each $\{u, v\} \in F$. Choosing a pairing F uniformly at random from among all possible pairings Ω_W of the points of W produces a random (multi-)graph G_F . Each r -regular simple graph G on vertex set $[n]$ is equally likely to be generated as G_F . Here, simple means without loops of multiple edges. Furthermore, if $r = O(1)$ then G_F is simple with a probability bounded below by a positive value independent of n . Therefore, any event that occurs *whp* in G_F will also occur *whp* in $G(n, r)$.

1.1.2.3 Preferential attachment

The configuration model can be used to generate graphs with power law degree distributions. Assuming that $d = (d_1, \dots, d_n)$ is a graphical sequence following a power law distribution, each vertex i obtains d_i configuration points. Then a pairing F results in a multigraph with degree sequence d . However, this mechanism does not provide any insights on how a dynamic network that evolves over time exhibits a power law degree distribution.

Albert-László Barabási and Réka Albert in a highly influential paper [50] provide a dynamic mechanism that explains how power law degree sequences emerge in real-world networks. We present a generalized version of their model with two parameters m and $\delta \geq -1$ as presented in §8.1 in [403]. The original version [50] is a subcase of this general model for $\delta = 0$. The model generates a sequence of graphs which we denote by $\{PA_t(m, \delta)\}_{t=1}^{+\infty}$ which for every t yields a graph with t vertices and mt edges. We define the model for $m = 1$. The model $PA_t(m, \delta)$ where $m \geq 2$ is reduced to the case $m = 1$ by running the model $PA_{mt}(1, \delta)$ and collapsing sequences of m consecutive vertices. Let $\{v_1, \dots, v_t\}$ be the set of vertices of $PA_t(m, \delta)$ and let $d_i(t)$ be the degree of vertex v_i at time t . Initially at time 1 the graph $PA_1(1, \delta)$ consists of a single vertex with a loop. The growth follows the following preferential rule. To obtain $PA_{t+1}(1, \delta)$ from $PA_t(1, \delta)$ a new vertex with a single edge is added to the graph. This edge chooses its second endpoint according to the following probability distribution. With probability $\frac{d_i(t)+\delta}{t(2+\delta)+(1+\delta)}$ vertex v_i is chosen, where $i \in [t]$, and with the remaining probability $\frac{1+\delta}{t(2+\delta)+(1+\delta)}$ a self loop is created.

Define $p_k = \left(2 + \frac{\delta}{m}\right) \frac{\Gamma(k+\delta)\Gamma(m+2+\delta+\frac{\delta}{m})}{\Gamma(m+\delta)\Gamma(k+3+\delta+\frac{\delta}{m})}$ for $k \geq m$ where $\Gamma(t) = \int_0^{+\infty} x^{t-1}e^{-x}dx$ is the Γ function and $P_k(t) = \frac{1}{t} \sum_{i=1}^t \mathbf{1}(d_i(t) = k)$. Also, let $N_k(t) = tP_k(t)$. It turns out that $\mathbb{E}[N_k(t)] \approx p_k t$ and that the degree sequence is strongly concentrated around its expectation. Specifically, for any $C > m\sqrt{8}$ as $t \rightarrow +\infty$ the following concentration inequality holds.

$$\Pr \left[\max_k |N_k(t) - \mathbb{E}[N_k(t)]| \geq C\sqrt{t \log t} \right] = o(1).$$

For the special case $\delta = 0$,

$$p_k = \frac{2m(m+1)}{k(k+1)(k+2)},$$

namely the probability distribution follows a power law with slope 3, as $p_k \sim k^{-3}$. It is worth mentioning that the model of preferential attachment was introduced conceptually by Barabási and Albert but it was Bollobás and Riordan with their collaborators who formally defined and studied the model [79, 80]. Power law degree sequences can also emerge with growth models based on optimization [157].

1.1.2.4 Kronecker graphs

Kronecker graphs [278, 279] are inspired by fractal theory [296]. There exist two versions of Kronecker graphs, a deterministic and a randomized one. To define each one, we remind the definition of the Kronecker product.

Definition 1.7. Given two matrices $A^{m \times n} = (a_{ij})$, $B^{m' \times n'} = (b_{ij})$, the Kronecker product matrix $C^{mm' \times nn'}$ is given by

$$C = A \otimes B = \begin{pmatrix} a_{11}B & \dots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \dots & a_{mn}B \end{pmatrix}.$$

Deterministic Kronecker graphs are defined by a small initiator adjacency matrix K_1 and the order k . A deterministic Kronecker graph of order k is defined by

$$K_1^{(k)} = \underbrace{K_1 \otimes \dots \otimes K_1}_{k \text{ times}}.$$

The stochastic Kronecker graphs use an initiator matrix with probabilities. The final adjacency matrix is the outcome of a randomized rounding of the k -th order Kronecker product of the initiator matrix. Kronecker graphs match several empirical properties such as heavy-tailed degree distributions and triangles, low diameters, and also obeys the densification power law. Most properties are analyzed in the deterministic case [282, 387]. Mahdian and Xu in an elegant paper studied stochastic Kronecker graphs. They show a phase transition for the emergence of the giant component and for connectivity, and prove that such graphs have constant diameters beyond the connectivity threshold [294]. Two additional appealing features of Kronecker graphs is the existence of methods to fit the parameters of the 2×2 initiator matrix to a given graph and their generation is embarassingly parallel.

Other popular models are the copying model [253, 262], the Cooper-Frieze model [123], the Aiello-Chung-Lu model [19, 20], protean graphs [345] the Fabrikant-Koutsoupias-Papadimitriou model [157], and the forest fire model [280].

1.1.3 Triangles

Subgraphs play a central role in graph theory. It is not an exaggeration to claim that studying subgraphs has been an active thread of research since the early days of graph theory: Euler paths and cycles, Hamilton paths and cycles, matchings, cliques, neighborhoods of vertices typically sketched via the degree sequence are special types of subgraphs.

Among various subgraphs, triangles play a major role in network analysis. A triangle is a clique of order 3. The number of triangles in a graph is a computationally expensive, crucial graph statistic in complex network analysis, in random graph models and in various important applications. Despite the fact that real-world networks tend to be sparse in edges, they are dense in triangles. This observation implies that when two vertices share a common neighbor, then it is more likely that they are/become connected. For instance, it has been observed in the MSN Messenger social network that if two people have a common contact it is 18 000 times more likely that they are connected [10]. The transitivity of adjacency is striking in social networks and in other types of networks too. There exist two processes that generate triangles in a social network: *homophily* and *transitivity*. According to the former, people tend to choose friends with similar characteristics to themselves (e.g., race, education) [409, 417] and according to the latter friends of friends tend to become friends themselves [409]. We survey a wide range of applications which rely on the number of triangles in a given graph.

Clustering Coefficients and Transitivity of a Graph: Despite the fact that Erdős-Rényi graphs have a short diameter they do not model social networks well. Social networks have many triangles. This was the main motivation of Watts and Strogatz [412] in their influential paper to introduce clustering coefficients and the notion of transitivity which we defined in a previous section.

Uncovering Hidden Thematic Structures: Eckmann and Moses [145] propose the use of the clustering coefficient for detecting subsets of web pages with a common topic. The key idea is that reciprocal links between pages indicate a mutual recognition/respect and then triangles due to their transitivity properties can be used to extend “seeds” to larger subsets of vertices with similar thematic structure in the Web graph. In other words, regions of the World Wide Web with high curvature indicate a common topic, allowing the authors to extract useful meta-information. This idea has found more applications, in natural language processing [138] and in bioinformatics [232, 348].

Exponential Random Graph Model: Frank and Strauss [176] proved under the assumption that two edges are dependent only if they share a common vertex that the

sufficient statistics for Markov graphs are the counts of triangles and stars. Wasserman and Pattison [410] proposed the exponential random graph (ERG) model which generalized the Markov graphs [347]. Triangles are frequently used as one of the sufficient statistics of the ERG model and counting them is necessary for parameter estimation, e.g., using Markov chain Monte Carlo (MCMC) procedures [68].

Spam Detection: Becchetti et al. [55] show that the distribution of triangles among spam hosts and non-spam hosts can be used as a feature for classifying a given host as spam or non-spam. The same result holds also for web pages, i.e., the spam and non-spam triangle distributions differ at a detectable level using standard statistical tests from each other.

Content Quality and Role Behavior Identification: Nowadays, there exist many online forums where acknowledged scientists participate, e.g., MathOverflow, CStheory Stackexchange and discuss problems of their fields. This yields significant information for researchers. Several interesting questions arise such as which participants comment on each other. This question including several others were studied in [414]. The number of triangles that a user participates was shown to play a critical role in answering these questions. For further applications in assessing the role behavior of users see [55].

Structural Balance and Status Theory: Balance theory appeared first in Heider's seminal work [210] and is based on the concept "the friend of my friend is my friend", "the enemy of my friend is my enemy" etc. [409]. To quantify this concept edges become signed, i.e., there is a function $c : E(G) \rightarrow \{+, -\}$. If all triangles are positive, i.e., the product of the signs of the edges is $+$, then the graph is balanced. Status theory is based on interpreting a positive edge (u, v) as u having lower status than v , while the negative edge (u, v) means that u regards v as having a lower status than himself/herself. Recently, Leskovec et al. [283] have performed experiments to quantify which of the two theories better apply to online social networks and predict signs of incoming links. Their algorithms require counts of signed triangles in the graph.

Microscopic Evolution of networks: Leskovec et al. [281] present an extensive experimental study of network evolution using detailed temporal information. One of their findings is that as edges arrive in the network, they tend to close triangles, i.e., connect people with common friends.

Community Detection: Counting triangles is important as subroutine in community detection algorithms. Berry et al. use triangle counting to deduce the edge support measure in their community detection algorithm [65]. Gleich and Seshadhri [197] show that heavy-tailed degree distributions and abundance in triangles imply that there exist vertices which together with their neighbors form a low-conductance set, i.e., community.

Motif Detection: Triangles are abundant not only in social networks but in biological networks [2, 217]. This fact can be used e.g., to correlate the topological and functional properties of protein interaction networks [217].

Triangular Connectivity [54]: Two vertices u, v are triangularly connected if there is a sequence of triangles $(\Delta_1, \dots, \Delta_s)$ such that u is a vertex in the Δ_1 , v in Δ_s and Δ_i shares at least one vertex with Δ_{i-1} .

k -truss: The k -truss of a graph G [120] is the maximum subgraph of G where every edge appears in at least $k - 2$ triangles.

Link recommendation: Triangle listing is used in link recommendation [393, 395].

CAD applications: Fudos and Hoffman [185] introduced a graph-constructive approach to solving systems of geometric constraints, a problem which arises frequently in Computer-Aided design (CAD) applications. One of the steps of their algorithm computes the number of triangles in an appropriately defined graph.

Given the large number of applications, there exists a lot of interest in developing efficient triangle listing and counting algorithms.

1.1.3.1 Triangle counting methods

There exist *exact* and *approximate* triangle counting algorithms. It is worth noting that for most of the applications described in Section 5.1 the exact number of triangles is not crucial. Hence, approximate counting algorithms which are fast and output a high quality estimate are desirable for the applications in which we are interested.

Exact Counting: Naive triangle counting by checking all triples of vertices takes $O(n^3)$ units of time. The state of the art algorithm is due to Alon, Yuster and Zwick [25] and runs in $O(m^{\frac{2\omega}{\omega+1}})$, where currently the fast matrix multiplication exponent ω is 2.3727 [416]. Thus, the Alon, Yuster, Zwick (AYZ) algorithm currently runs in $O(m^{1.407})$ time. It is worth mentioning that from a practical point of view algorithms based on matrix multiplication are not used due to the prohibitive memory requirements. Even for medium sized networks, i.e., networks with hundreds of thousands of edges, matrix-multiplication based algorithms are not applicable. Itai and Rodeh in 1978 showed an algorithm which finds a triangle in any graph in $O(m^{\frac{3}{2}})$ [220]. This algorithm can be extended to list the set of triangles in the graph with the same time complexity. Chiba and Nishizeki showed that triangles can be found in time $O(m\alpha(G))$ where $\alpha(G)$ is the *arboricity* of the graph. Since $\alpha(G)$ is at most $O(\sqrt{m})$ their algorithm runs in $O(m^{3/2})$ in the worst case [108]. For special types of graphs more efficient triangle counting

algorithms exist. For instance in planar graphs, triangles can be found in $O(n)$ time [108, 220, 331].

Even if listing algorithms solve a more general problem than the counting one, they are preferred in practice for large graphs, due to the smaller memory requirements compared to the matrix multiplication based algorithms. Simple representative algorithms are the node- and the edge-iterator algorithms. The former counts for each vertex v the number of triangles t_v it is involved in, i.e., the number of edges among its neighbors, whereas the latter algorithm counts for each edge (i, j) the common neighbors of vertices i, j . Both of these algorithms have the same asymptotic complexity $O(mn)$, which in dense graphs results in $O(n^3)$ time, the complexity of the naive counting algorithm. Practical improvements over this family of algorithms have been achieved using various techniques, such as hashing and sorting by the degree [269, 352].

Approximate Counting: On the approximate counting side, most of the triangle counting algorithms have been developed in the streaming setting. In this scenario, the graph is represented as a stream. Two main representations of a graph as a stream are the edge stream and the incidence stream. In the former, edges arrive one at a time. In the latter scenario all edges incident to the same vertex appear successively in the stream. The ordering of the vertices is assumed to be arbitrary. A streaming algorithm produces a $(1 + \epsilon)$ approximation of the number of triangles *whp* by making only a constant number of passes over the stream. However, sampling algorithms developed in the streaming literature can be applied in the setting where the graph fits in the memory as well. Monte Carlo sampling techniques have been proposed to give a fast estimate of the number of triangles. According to such an approach, a.k.a. naive sampling [351], we choose three nodes at random repetitively and check if they form a triangle or not. If one makes

$$r = \log\left(\frac{1}{\delta}\right) \frac{1}{\epsilon^2} \left(1 + \frac{T_0 + T_1 + T_2}{T_3}\right)$$

independent trials where T_i is the number of triples with i edges and outputs as the estimate of triangles the random variable T'_3 equaling to the fractions of triples picked that form triangles times the total number of triples $\binom{n}{3}$, then

$$(1 - \epsilon)T_3 < T'_3 < (1 + \epsilon)T_3$$

with probability at least $1 - \delta$. This is suitable only when $T_3 = o(n^2)$.

In [47] the authors reduce the problem of triangle counting efficiently to estimating moments for a stream of node triples. Then, they use the Alon-Matias-Szegedy (AMS) algorithms [28] to proceed. The key is that the triangle computation reduces to estimating the zero-th, first and second frequency moments, which can be done efficiently.

Furthermore, as the authors suggest their algorithm is efficient only on graphs with $\Omega(n^2/\log \log n)$ triangles, i.e., triangle dense graphs as in the naive sampling. The AMS algorithms are also used by [230], where simple sampling techniques are used, such as choosing an edge from the stream at random and checking how many common neighbors its two endpoints share considering the subsequent edges in the stream. Along the same lines, Buriol et al. [96] proposed two space-bounded sampling algorithms to estimate the number of triangles. Again, the underlying sampling procedures are simple. For instance, in the case of the edge stream representation, they sample randomly an edge and a node in the stream and check if they form a triangle. The three-pass algorithm presented therein, counts in the first pass the number of edges, in the second pass it samples uniformly at random an edge (i, j) and a node $k \in V - \{i, j\}$ and in the third pass it tests whether the edges $(i, k), (k, j)$ are present in the stream. The number of samples r needed to obtain an $(1 \pm \epsilon)$ -approximation with probability $1 - \delta$ is

$$r = O\left(\log\left(\frac{1}{\delta}\right) \frac{T_1 + 2T_2 + 3T_3}{T_3 \epsilon^2}\right) = O\left(\log\left(\frac{1}{\delta}\right) \frac{mn}{t \epsilon^2}\right).$$

Even if the term T_0 in the nominator is missing¹ compared to the naive sampling, the graph has still to be fairly dense with respect to the number of triangles in order to get a $(1 \pm \epsilon)$ -approximation *whp*. Buriol et al. [96] show how to turn the three-pass algorithm into a single pass algorithm for the edge stream representation and similarly they provide a three- and one-pass algorithm for the incidence stream representation. Kane et al. show how to count other subgraphs in the streaming model [235]. In [55] the semi-streaming model for counting triangles is introduced, which allows $\log n$ passes over the edges. The key observation is that since counting triangles reduces to computing the intersection of two sets, namely the induced neighborhoods of two adjacent nodes, ideas from locality sensitivity hashing [91] are applicable to the problem.

Another line of work is based on linear algebraic arguments. Specifically, in the case of “power-law” networks it was shown in [390] that the spectral counting of triangles can be efficient due to their special spectral properties [113]. This idea was further extended in [388] using the randomized Singular Value Decomposition (SVD) approximation algorithm by [140]. More recently, Avron proposed a new approximate triangle counting method based on a randomized algorithm for trace estimation [44].

Graph Sparsifiers: A sparsifier of a graph $G(V, E, w)$ is a sparse graph H that is similar to G in some useful notion. We discuss in the following the Benczúr-Karger cut sparsifier [63, 64] and the Spielman-Srivastava spectral sparsifier [370].

¹Notice that $m(n-2) = T_1 + 2T_2 + 3T_3$ and $t = T_3$.

Benczúr-Karger Sparsifier: Benczúr and Karger introduced in [64] the notion of cut sparsification to accelerate cut algorithms whose running time depends on the number of edges. Using a non-uniform sampling scheme they show that given a graph $G(V, E, w)$ with $|V| = n, |E| = m$ and a parameter ϵ there exists a graph $H(V, E', w')$ with $O(n \log(n)/\epsilon^2)$ edges such that the weight of every cut in H is within a factor of $(1 \pm \epsilon)$ of its weight in G . Furthermore, they provide a nearly-linear time algorithm which constructs such a sparsifier. The key quantity used in the sampling scheme of Benczúr and Karger is the strong connectivity $c_{(u,v)}$ of an edge $(u, v) \in E$ [63, 64]. The latter quantity is defined to be the maximum value k such that there is an induced subgraph G_0 of G containing both u and v , and every cut in G_0 has weight at least k .

Spielman-Srivastava Sparsifier: In [370] Spielman and Teng introduced the notion of a spectral sparsifier in order to strengthen the notion of a cut sparsifier. A quantity that plays a key role in spectral sparsifiers is the *effective resistance*. The term effective resistance comes from electrical network analysis, see Chapter IX in [76]. In a nutshell, let $G(V, E, w)$ be a weighted graph with vertex set V , edge set E and weight function w . We call the weight $w(e)$ *resistance* of the edge e . We define the *conductance* $r(e)$ of e to be the inverse of the resistance $w(e)$. Let \mathcal{G} be the resistor network constructed from $G(V, E, w)$ by replacing each edge e with an electrical resistor whose electrical resistance is $w(e)$. Typically, in \mathcal{G} vertices are called *terminals*, a convention that emphasizes the electrical network perspective of a graph G . The *effective resistance* $R(i, j)$ between two vertices i, j is the electrical resistance measured across vertices i and j in \mathcal{G} . Equivalently, the effective resistance is the potential difference that appears across terminals i and j when we apply a unit current source between them. Finally, the *effective conductance* $C(i, j)$ between two vertices i, j is defined as $C(i, j) = R^{-1}(i, j)$.

Spielman and Srivastava in their seminal work [370] proposed to include each edge of G in the sparsifier H with probability proportional to its effective resistance. They provide a nearly-linear time algorithm that produces spectral sparsifiers with $O(n \log n)$ edges.

1.1.4 Dense Subgraphs

Finding dense subgraphs is a key problem for many applications and the key primitive for community detection. Here we review some important concepts of dense subgraphs.

Cliques: A clique is a subset of vertices all connected to each other. The problem of finding whether there exists a clique of a given size in a graph is **NP**-complete. A *maximum* clique of a graph is a clique of maximum possible size and its size is called the graph's clique number. Håstad [212] shows that, unless $\mathbf{P} = \mathbf{NP}$, there cannot be any polynomial time algorithm that approximates the maximum clique within a factor

better than $\mathcal{O}(n^{1-\epsilon})$, for any $\epsilon > 0$. Feige [161] proposes a polynomial-time algorithm that finds a clique of size $\mathcal{O}((\frac{\log n}{\log \log n})^2)$ whenever the graph has a clique of size $\mathcal{O}(\frac{n}{\log^b n})$ for any constant b . Based on this, an algorithm that approximates the maximum clique problem within a factor of $\mathcal{O}(n^{\frac{(\log \log n)^2}{\log n^3}})$ is also defined. A *maximal* clique is a clique that is not a subset of any other clique. The Bron-Kerbosch algorithm [92] finds all maximal cliques in a graph in exponential time. A near optimal time algorithm for sparse graphs was introduced in [151].

Densest Subgraph: Let $G(V, E)$ be a graph, $|V| = n$, $|E| = m$. The average degree of a vertex set $S \subseteq V$ is defined as $\frac{2e[S]}{|S|}$, where $e[S]$ is the number of edges in the induced graph $G[S]$. The densest subgraph problem is to find a set S that maximizes the average degree. The densest subgraph can be identified in polynomial time by solving a maximum-flow problem [186, 199]. Charikar [104] shows that the greedy algorithm proposed by Asashiro et al. [40] produces a $\frac{1}{2}$ -approximation of the densest subgraph in linear time. Both algorithms are efficient in terms of running times and scale to large networks. In the case of directed graphs, the densest subgraph problem is solved in polynomial time as well. Charikar [104] provided a linear programming approach which requires the computation of n^2 linear programs and a $\frac{1}{2}$ -approximation algorithm which runs in $\mathcal{O}(n^3 + n^2m)$ time. Khuller and Saha [248] improved significantly the state-of-the-art by providing an exact combinatorial algorithm and a fast $\frac{1}{2}$ -approximation algorithm which runs in $\mathcal{O}(n + m)$ time. Kannan and Vinay [241] gave a spectral $\mathcal{O}(\log n)$ approximation algorithm for a related notion of density.

In the classic definition of densest subgraph there is no size restriction of the output. When restrictions on the size $|S|$ are imposed, the problem becomes **NP**-hard. Specifically, the DkS problem of finding the densest subgraph of k vertices is known to be **NP**-hard [41]. For general k , Feige et al. [163] provide an approximation guarantee of $\mathcal{O}(n^\alpha)$, where $\alpha < \frac{1}{3}$. The greedy algorithm by Asahiro et al. [40] gives instead an approximation factor of $\mathcal{O}(\frac{n}{k})$. Better approximation factors for specific values of k are provided by algorithms based on semidefinite programming [162]. From the perspective of (in)approximability, Khot [247] shows that there cannot exist any PTAS for the DkS problem under a reasonable complexity assumption. Arora et al. [38] propose a PTAS for the special case $k = \Omega(n)$ and $m = \Omega(n^2)$. Finally, two variants of the DkS problem are introduced by Andersen and Chellapilla [31]. The two problems ask for the set S that maximizes the average degree subject to $|S| \leq k$ (DamkS) and $|S| \geq k$ (DalkS), respectively. The authors provide constant factor approximation algorithms for both DamkS and DalkS.

Quasi-cliques: A set of vertices S is an α -quasi-clique (or pseudo-clique) if $e[S] \geq$

$\alpha^{\binom{|S|}{2}}$, i.e., if the edge density of the induced subgraph $G[S]$ exceeds a threshold parameter $\alpha \in (0, 1)$. Similarly to cliques, *maximum* quasi-cliques [333] and *maximal* quasi-cliques [94] are quasi-cliques of maximum size and quasi-cliques not contained into any other quasi-clique, respectively. Abello et al. [8] propose an algorithm for finding a single maximal α -quasi-clique, while Uno [399] introduces an algorithm to enumerate all α -quasi-cliques.

k -core, k -clubs, kd -cliques: A k -core is a maximal connected subgraph in which all vertices have degree at least k . There exists a linear time algorithm for finding k -cores by repeatedly removing the vertex having the smallest degree [53]. A k -club is a subgraph whose diameter is at most k [306]. kd -cliques differ from k -clubs as the shortest paths used to compute the diameter of a kd -clique are allowed to use vertices not belonging to that kd -clique. All these clique variants are clearly conceptually different from the optimal quasi-cliques we study in this paper.

1.1.5 Graph Partitioning

Graph partitioning is a fundamental computer science problem. As we discussed above, the problem of finding communities is reduced to understanding the cut structure of the graph. In distributed computing applications, the following version of the graph partitioning problem plays a key role. The interested reader may consult the cited work and the references therein for more information on the balanced graph partitioning problem.

Balanced graph partitioning: The *balanced graph partitioning* problem is a classic **NP**-hard problem of fundamental importance to parallel and distributed computing [189]. The input of this problem is an undirected graph $G(V, E)$ and an integer $k \in \mathbb{Z}^+$, the output is a partition of the vertex set in k balanced parts such that the number of edges across the clusters is minimized. Formally, the balance constraint is defined by the imbalance parameter ν . Specifically, the (k, ν) -balanced graph partitioning asks to divide the vertices of a graph in k clusters each of size at most $\nu \frac{n}{k}$, where n is the number of vertices in G . The case $k = 2, \nu = 1$ is equivalent to the **NP**-hard minimum bisection problem. Several approximation algorithms, e.g., [160], and heuristics, e.g., [167, 246] exist for this problem. When $\nu = 1 + \epsilon$ for any desired but fixed ϵ there exists a $O(\epsilon^{-2} \log^{1.5} n)$ approximation algorithm [258]. When $\nu = 2$ there exists an $O(\sqrt{\log k \log n})$ approximation algorithm based on semidefinite programming (SDP) [260]. Due to the practical importance of k -partitioning there exist several heuristics,

among which METIS [353] and its parallel version [354] stand out for their good performance. METIS is widely used in many existing systems [242]. There are also heuristics that improve efficiency and partition quality of METIS in a distributed system [349].

Streaming balanced graph partitioning:

Despite the large amount of work on the balanced graph partitioning problem, neither state-of-the-art approximation algorithms nor heuristics such as METIS are well tailored to the computational restrictions that the size of today’s graphs impose. Motivated by this fact, Stanton and Kliot introduced the streaming balanced graph partitioning problem, where the graph arrives in the stream and decisions about the partition need to be taken with on the fly quickly [373]. Specifically, the vertices of the graph arrive in a stream with the set of edges incident to them. When a vertex arrives, a partitioner decides where to place the vertex. A vertex is never moved after it has been assigned to one of the k machines. A realistic assumption that can be used in real-world streaming graph partitioners is the existence of a small-sized buffer. Stanton and Kliot evaluate partitioners with or without buffers. The work of [373] can be adapted to edge streams. Stanton showed that streaming graph partitioning algorithms with a single pass under an adversarial stream order cannot approximate the optimal cut size within $o(n)$. The same bound holds also for random stream orders [374]. Finally, Stanton [374] analyzes two variants of well performing algorithms from [373] on random graphs. Specifically, she proves that if the graph G is sampled according to the planted partition model, then the two algorithms despite their similarity can perform differently and that one of the two can recover the true partition *whp*, assuming that inter-, intra-cluster edge probabilities are constant, and their gap is a large constant.

We conclude our brief exposition by outlining the differences between community detection methods and the balanced partitioning problem. One main difference is the lack of restriction on the number of vertices per subset in the community detection problem. A second difference is that in realistic applications the number of clusters in the balanced partitioning problem is part of the input, as it represents the number of machines/clusters available to distribute the graph. In community detection the number of clusters is not known a priori, or even worse, their existence is not clear. It is worth mentioning at this point that in Chapters 6 and 8 we introduce measures conceptually close to the modularity measure [196, 314, 316]. Despite the popularity of modularity, few rigorous results exist. Specifically, Brandes et al. proved that maximizing modularity is **NP**-hard [88]. Approximation algorithms without theoretical guarantees whose performance is evaluated in practice also exist [13].

1.1.6 Big Graph Data Analytics

Except for the algorithmic ‘dasein’ of computer science, there is an engineering one too. An important engineering law is Moore’s law. Gordon Moore based on observations from 1958 until 1965 extrapolated that the number of components in integrated circuits would keep doubling for at least until 1975 [308]. It is remarkable that Moore’s prediction remains (more or less) valid since then. However, as we are approaching the end of its validity, it is becoming clear that in order to perform demanding computational tasks, we need more than one machine. At the same time, input size increases. Currently, the growth rate is unprecedented. Eron Kelly, the general manager of product marketing for Microsoft SQL Server, predicts that as humankind we will generate more data as humankind than we generated in the previous 5,000 years [7]. The term *big data* describes collections of large and complex datasets which are difficult to manipulate and process using traditional tools. Mainly, for these two reasons, i.e., hardware reaching its limits and big data, parallel and distributed computing are the *de facto* solutions for processing large scale data. For this reason, there exists a lot of interest in developing efficient graph processing systems. Popular graph processing platforms are Pregel [295] and its open-source version Apache Giraph that build on MapReduce , and GraphLab [291]. It is worth mentioning that for dynamic graphs there exist other platforms which are suitable for stream/micro-batch processing, such as Twitter’s Storm [6].

In the following we discuss the details of MAPREDUCE [130], which we use in this dissertation as the underlying distributed system to develop efficient large-scale graph processing algorithms and systems.

1.1.6.1 MapReduce Basics

While the PRAM model [229] and the bulk-synchronous parallel model (BSP) [400] are powerful models, MAPREDUCE has largely “taken over” both industry and academia [4]. In few words, this success is due to two reasons: first, MAPREDUCE is a simple and powerful programming model which makes the programmer’s life easy. Secondly, MAPREDUCE is publicly available via its open source version HADOOP. MAPREDUCE was introduced in [130] by Google, one of the largest users of multiple processor computing in the world, for facilitating the development of scalable and fault tolerant applications. In the MAPREDUCE paradigm, a parallel computation is defined on a set of values and consists of a series of map, shuffle and reduce steps. Let (x_1, \dots, x_n) be the set of values, m denote the mapping function which takes a value x and returns a pair of a key k and a value u and r the reduce function.

1. In the map step a mapping function m is applied to a value x_i and a pair (k_i, u_i) of a key k_i and a value u_i is generated.
2. The shuffle step starts upon having mapped all values x_i for $i = 1$ to n to pairs. In this step, a set of lists is produced using the key-value pairs generated from the map step with an important feature. Each list is characterized by the key k and has the form $L_k = \{k : u_1, \dots, u_{j(k)}\}$ if and only if there exists a pair (k, u_i) for $i = 1$ to j .
3. Finally in the reduce step, the reduce function r is applied to the lists generated from the shuffle step to produce the set of values (w_1, w_2, \dots) .

To illustrate the aforementioned abstract concepts consider the problem of counting how many times each word in a given document appears. The set of values is the “bag-of-words” appearing in the document. For example, if the document is the sentence “The dog runs in the forest”, then $\{x_1, x_2, x_3, x_4, x_5, x_6\} = \{ \text{the, dog, runs, in, the, forest} \}$. One convenient choice for the MAPREDUCE functions is the following and results in the following steps: The map function m will map a value x to a pair of a key and a value. A convenient choice for m is something close to the identity map. Specifically, we choose $m(x) = (x, \$)$, where we assume that the dollar sign $\$$ an especially reserved symbol. The shuffle step for our small example will produce the following set of lists: (the:), (dog:\$), (runs:\$), (in:\$), (runs:\$), (forest:\$) The reduce function r will process each list defined by each different word appearing in the document by counting the number of dollar signs $\$$. This number will also be the count of times that specific word appears in the text.

HADOOP implements MAPREDUCE and was originally created by Doug Cutting. Even if HADOOP is well known for MAPREDUCE it is actually a collection of subprojects that are closely related to distributed computing. For example HDFS (HADOOP filesystem) is a distributed filesystem that provides high throughput access to application data and HBase is a scalable, distributed database that supports structured data storage for large tables (column-oriented database). Another subproject is Pig, which is a high-level data-flow language and execution framework for parallel computation [190]. Pig runs on HDFS and MAPREDUCE. For more details and other subprojects, the interested reader can visit the website that hosts the HADOOP project [4].

1.2 Computational Cancer Biology

Human cancer is caused by the accumulation of genetic alternations in cells [43, 413]. It is a complex phenomenon often characterized by the successive acquisition of combinations

of genetic aberrations that result in malfunction or dysregulation of genes. Finding driver genetic mutations, i.e., mutations which confer growth advantage on the cells carrying them and have been positively selected during the evolution of the cancer and uncovering their temporal sequence have been central goals of cancer research the last decades [376]. In this Section we review three problems arising in computational cancer biology. In Section 1.2.1 we present background on a data denoising problem. In Section 1.2.3 we review oncogenetic trees, a popular model for oncogenesis. Finally, in Section 1.2.2 we discuss the problem of discovering cancer subtypes.

1.2.1 Denoising array-based Comparative Genomic Hybridization (aCGH) data

There are many forms of chromosome aberration that can contribute to cancer development, including polyploidy, aneuploidy, interstitial deletion, reciprocal translocation, non-reciprocal translocation, as well as amplification, again with several different types of the latter (e.g., double minutes, HSR and distributed insertions [23]). Identifying the specific recurring aberrations, or sequences of aberrations, that characterize particular cancers provides important clues about the genetic basis of tumor development and possible targets for diagnostics or therapeutics. Many other genetic diseases are also characterized by gain or loss of genetic regions, such as Down Syndrome (trisomy 21) [275], Cri du Chat (5p deletion) [276], and Prader-Willi syndrome (deletion of 15q11-13) [98] and recent evidence has begun to suggest that inherited copy number variations are far more common and more important to human health than had been suspected just a few years ago [424]. These facts have created a need for methods for assessing DNA copy number variations in individual organisms or tissues.

In Chapter 10, we focus specifically on array-based comparative genomic hybridization (aCGH) [70, 231, 341, 342], a method for copy number assessment using DNA microarrays that remains, for the moment, the leading approach for high-throughput typing of copy number abnormalities. The technique of aCGH is schematically represented in Figure 1.2. A test and a reference DNA sample are differentially labeled and hybridized to a microarray and the ratios of their fluorescence intensities is measured for each spot. A typical output of this process is shown in Figure 1.2 (3), where the genomic profile of the cell line GM05296 [365] is shown for each chromosome. The x -axis corresponds to the genomic position and the y -axis corresponds to a noisy measurement of the ratio $\log_2 \frac{T}{R}$ for each genomic position, typically referred to as “probe” by biologists. For healthy diploid organisms, $R=2$ and T is the DNA copy number we want to infer from the noisy measurements. For more details on the use of aCGH to detect different types of chromosomal aberrations, see [23].

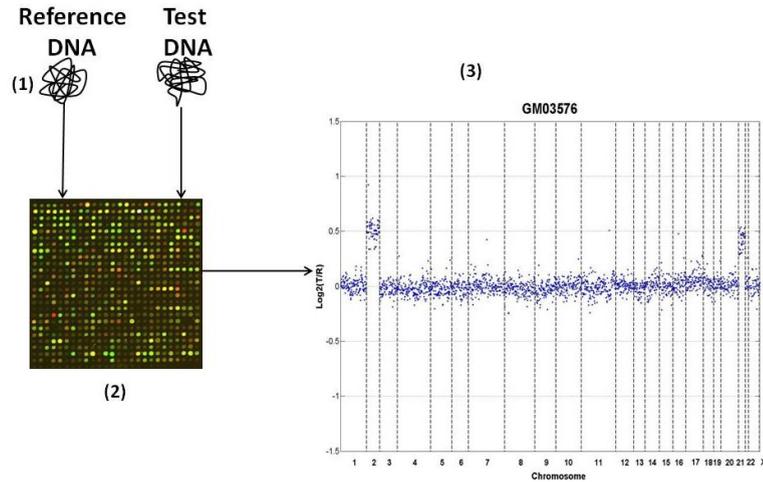


FIGURE 1.2: Schematic representation of array CGH. Genomic DNA from two cell populations (1) is differentially labeled and hybridized in a microarray (2). Typically the reference DNA comes from a normal subject. For humans this means that the reference DNA comes from a normal diploid genome. The ratios on each spot are measured and normalised so that the median \log_2 ratio is zero. The final result is an ordered tuple containing values of the fluorescent ratios in each genomic position per each chromosome. This is shown in (3) where we see the genomic profile of the cell line GM05296 [365]. The problem of denoising array CGH data is to infer the true DNA copy number T per genomic position from a set of noisy measurements of the quantity $\log_2 \frac{T}{R}$, where $R=2$ for normal diploid humans.

Converting raw aCGH log fluorescence ratios into discrete DNA copy numbers is an important but non-trivial problem. Finding DNA regions that consistently exhibit chromosomal losses or gains in cancers provides a crucial means for locating the specific genes involved in development of different cancer types. It is therefore important to distinguish, when a probe shows unusually high or low fluorescence, whether that aberrant signal reflects experimental noise or a probe that is truly found in a segment of DNA that is gained or lost. Furthermore, successful discretization of array CGH data is crucial for understanding the process of cancer evolution, since discrete inputs are required for a large family of successful evolution algorithms, e.g., [132, 134]. It is worth noting that manual annotation of such regions, even if possible [365], is tedious and prone to mistakes due to several sources of noise (impurity of test sample, noise from array CGH method, etc.). A well-established observation that we use in Chapter 10 is that near-by probes tend to have the similar DNA copy number.

Many algorithms and objective functions have been proposed for the problem of discretizing and segmenting aCGH data. Many methods, starting with [178], treat aCGH segmentation as a hidden Markov model (HMM) inference problem. The HMM approach has since been extended in various ways, e.g., through the use of Bayesian HMMs

[206], incorporation of prior knowledge of locations of DNA copy number polymorphisms [362], and the use of Kalman filters [363]. Other approaches include wavelet decompositions [214], quantile regression [147], expectation-maximization in combination with edge-filtering [310], genetic algorithms [228], clustering-based methods [408, 421], variants on Lasso regression [215, 379], and various problem-specific Bayesian [51], likelihood [216], and other statistical models [286]. A dynamic programming approach, in combination with expectation maximization, has been previously used by Picard et al. [340]. In [263] and [415] an extensive experimental analysis of available methods has been conducted. Two methods stand out as the leading approaches in practice. One of these top methods is CGHSEG [339], which assumes that a given CGH profile is a Gaussian process whose distribution parameters are affected by abrupt changes at unknown coordinates/breakpoints. The other method which stands out for its performance is Circular Binary Segmentation [321] (CBS), a modification of binary segmentation, originally proposed by Sen and Srivastava [360], which uses a statistical comparison of mean expressions of adjacent windows of nearby probes to identify possible breakpoints between segments combined with a greedy algorithm to locally optimize breakpoint positions.

1.2.2 Cancer subtypes

Genomic studies have dramatically improved our understanding of the biology of tumor formation and treatment. In part this has been accomplished by harnessing tools that profile the genes and proteins in tumor cells, revealing previously indistinguishable tumor subtypes that are likely to exhibit distinct sensitivities to treatment methods [200, 338]. As these tumor subtypes are uncovered, it becomes possible to develop novel therapeutics more specifically targeted to the particular genetic defects that cause each cancer [42, 71, 336]. While recent advances have had a profound impact on our understanding of tumor biology, the limits of our understanding of the molecular nature of cancer obstruct the burgeoning efforts in “targeted therapeutics” development. These limitations are apparent in the high failure rate of the discovery pipeline for novel cancer therapeutics [233] as well as in the continuing difficulty of predicting which patients will respond to a given therapeutic. A striking example is the fact that trastuzumab, the targeted therapeutic developed to treat HER2-amplified breast cancers, is ineffective in many patients who have HER2-overexpressing tumors and yet effective in some who do not [326]. Furthermore, subtypes typically remain poorly defined — *e.g.*, the “basal-like” breast cancer subtype, for which different studies have inferred very distinct genetic signatures [338, 367] — and yet many patients do not fall into any known subtype. Our belief, then, is that clinical treatment of cancer will reap considerable benefit from the identification of new cancer subtypes and genetic signatures.

One promising approach for better elucidating the common mutational patterns by which tumors develop is to recognize that tumor development is an evolutionary process and apply phylogenetic methods to tumor data to reveal these evolutionary relationships. Much of the work on tumor evolution models flows from the seminal efforts of [131] on inferring *oncogenetic trees* from comparative genomic hybridization (aCGH) profiles of tumor cells. A strength in this model stems from the extraction of ancestral structure from many probe sites per tumor, potentially utilizing measurements of the expression or copy number changes across the entire genome. However, this comes at the cost of overlooking the diversity of cell populations within tumors, which can provide important clues to tumor progression but are conflated with one another in tissue-wide assays like aCGH. The cell-by-cell approaches, such as [337, 361], use this heterogeneity information but at the cost of allowing only a small number of probes per cell. Schwartz and Shackney [358] proposed bridging the gap between these two methodologies by computationally inferring cell populations from tissue-wide gene expression samples. This inference was accomplished through “geometric unmixing,” a mathematical formalism of the problem of separating components of mixed samples in which each observation is presumed to be an unknown convex combination² of several hidden fundamental components. Other approaches to inferring common pathways include mixture models of oncogenetic trees [58], PCA-based methods [218], conjunctive Bayesian networks [191] and clustering [288].

Unmixing falls into the class of methods that seek to recover a set of pure sources from a set of mixed observations. Analogous problems have been coined “the cocktail problem,” “blind source separation,” and “component analysis” and various communities have formalized a collection of models with distinct statistical assumptions. In a broad sense, the classical approach of principal component analysis (PCA) [335] seeks to factor the data under the constraint that, collectively, the fundamental components form an orthonormal system. Independent component analysis (ICA) [121] seeks a set of statistically independent fundamental components. These methods, and their ilk, have been extended to represent non-linear data distributions through the use of kernel methods (see [355, 356] for details), which often confound modeling with black-box data transformations. Both PCA and ICA break down as pure source separators when the sources exhibit a modest degree of correlation. Collectively, these methods place strong independence constraints on the fundamental components that are unlikely to hold for tumor samples, where we expect components to correspond to closely related cell states.

²A point p is a convex combination of basis points v_0, \dots, v_k if and only if the constraints $p = \sum_{i=0}^k \alpha_i v_i$, $\sum_i \alpha_i = 1$ and $\forall i : \alpha_i \geq 0$ obtain. The fractions α_i determine a mixture over the basis points $\{v_i\}$ that produce the location p .

Extracting multiple correlated fundamental components, has motivated the development of new methods for unmixing genetic data. Similar unmixing methods were first developed for tumor samples by Billheimer and colleagues [155] to improve the power of statistical tests on tumor samples in the presence of contaminating stromal cells. Similarly, a hidden Markov model approach to unmixing was developed by Lamy *et al.* [266] to correct for stromal contamination in DNA copy number data. These recent advances demonstrate the feasibility of unmixing-based approaches for separating cell sub-populations in tumor data. Outside the bioinformatics community, geometric unmixing has been successfully applied in the geo-sciences [146] and in hyper-spectral image analysis [102].

The recent work by [358] applied the hard geometric unmixing model to gene expression data with the goal of recovering expression signatures of tumor cell subtypes, with the specific goal of facilitating phylogenetic analysis of tumors. The results showed promise in identifying meaningful sub-populations and improving phylogenetic inferences.

1.2.3 Oncogenetic trees

Among the triumphs of cancer research stands the breakthrough work of Vogelstein and his collaborators [159, 405] which provides significant insight into the evolution of colorectal cancer. Specifically, the so-called “Vogelgram” models colorectal tumorigenesis as a linear accumulation of certain genetic events. Few years later, Desper *et al.* [133] considered more general evolutionary models compared to the “Vogelgram” and presented one of the first theoretical approaches to the problem [43], the so-called *oncogenetic trees*. Before we provide a description of oncogenetic trees which are the focus of our work, we would like to emphasize that since then a lot of research work has followed from several groups of researchers, influenced by the seminal work of Desper *et al.* [133]. Currently there exists a wealth of methods that infer evolutionary models from microarray-based data such as gene expression and array Comparative Genome Hybridization (aCGH) data: distance based oncogenetic trees [135], maximum likelihood oncogenetic trees [406], hidden variable oncogenetic trees [380], conjunctive Bayesian networks [59] and their extensions [56, 192], mixture of trees [57]. The interested reader is urged to read the surveys of Attolini *et al.* [43] and Hainke *et al.* [207] and the references therein on established progression modeling methods. Furthermore, oncogenetic trees have successfully shed light into many types of cancer such as renal cancer [133], hepatic cancer [289] and head and neck squamous cell carcinomas [215].

An oncogenetic tree is a rooted directed tree³. The root represents the healthy state of tissue with no mutations. Any other vertex $v \in V$ represents a mutation. Each edge represents a “cause-and-effect” relationships. Specifically, for a mutation represented by vertex v to occur, all the mutations corresponding to vertices that lie on the directed path from the root to v must be present in the tumor. In other words, if two mutations u, v are connected by an edge (u, v) then v cannot occur if u has not occurred. The edges are labeled with probabilities. Each tumor corresponds to a rooted subtree of the oncogenetic tree and the probability of occurrence is determined as described by [133]. Desper et al. provide an algorithm that finds a likely oncogenetic tree that fits the observed data.

1.3 Thesis Overview

In this Section we motivate our work and provide an overview of this dissertation.

Rainbow Connectivity (Chapter 3)

Connectivity is a fundamental graph theoretic property [83]. The most well-studied connectivity concept asks for the minimum number of vertices or edges which need to be removed in order to disconnect the graph. However, there exist other graph theoretic concepts that strengthen the connectivity concept: imposing bounds on the diameter, existence of edge disjoint spanning trees etc. In 2006 Chartrand et al. [105] defined the concept of *rainbow connectivity*, also referred as *rainbow connection*. We prefer to provide two motivating examples rather than the exact definition which is found in Chapter [Rainbow Connectivity of Sparse Random Graphs](#).

Suppose we wish to route messages in a cellular network G , between any two vertices in a pipeline, and require that each link on the route between the vertices (namely, each edge on the path) is assigned a distinct channel (e.g., a distinct frequency). The minimum number of distinct channels we need to use is the rainbow connectivity of G .

Another motivating example is related to securing communication between government agencies [285]. The Department of Homeland Security of USA was created in 2003 in response to the weaknesses discovered in the secure transfer of classified information after the September 11, 2001 terrorist attacks. Ericksen [154] observed that because of the unexpected aftermath law enforcement and intelligence agencies could not communicate. Given that this situation could not have been easily predicted, the technologies utilized were separate entities and prohibited shared access, meaning that there was no way for

³Typically, the term *tree* is reserved for the undirected case and the term *branching* for the directed case. In the context of oncogenetic trees, we consistently use the term *tree* for a directed tree as in [133].

officers and agents to cross check information between various organizations. While the information needs to be protected since it relates to national security, there must also be procedures that permit access between appropriate parties. This twofold issue can be addressed by assigning information transfer paths between agencies which may have other agencies as intermediaries while requiring a large enough number of passwords and firewalls that is prohibitive to intruders, yet small enough to manage. Equivalently, this number has to be large enough so that one or more paths between every pair of agencies have no password repeated. Rainbow connectivity arises as the natural answer to the following question: what is the minimum number of passwords or firewalls needed that allows at least one path between every two agencies so that the passwords along each path are distinct?

Contributions:

In [182, 184] we prove the following results on the rainbow connectivity of sparse random graphs.

- For an Erdős-Rényi random graph $G = G(n, p)$ at the connectivity threshold, i.e., $p = \frac{\log n + \omega}{n}$, $\omega \rightarrow \infty, \omega = o(\log n)$, we prove Theorem 3.1 which characterizes optimally the rainbow connectivity *whp*. Our proof is constructive in the following sense: a random coloring is *whp* a valid rainbow coloring.
- For random regular graphs [419] we prove Theorem 3.2. The proof of Theorem 3.2 is still constructive, but requires an unexpected use of a Markov Chain Monte Carlo algorithm.

Random Apollonian Graphs (Chapter 4)

In Chapter 4 we analyze Random Apollonian Networks (RANs) [426], a popular random graph model for real-world networks. Compared to other models, RANs generate planar graphs. This makes RANs special for at least two reasons. Firstly, planar graphs form an important family of graphs for various reasons. They model several significant types of spatial real-world networks such as power grids, water distribution networks and road networks. For instance, a street network has edges corresponding to roads and vertices represent roads' intersections and endpoints. Since edges intersect only at vertices, street networks are planar. It is worth mentioning that planarity of street networks is almost always violated in practice because of bridges. However planarity is a good approximation [265]. It has been observed through various experimental studies that real-world planar graphs have distinct features [99, 111, 116, 225, 265] from random planar graphs [300]. One such feature is that the degrees are skewed, obeying a power law degree distribution [265]. A recent paper which surveys properties and models of real-world planar graphs is [52]. Despite the outstanding amount of work on modeling

real-world networks with random graph models, e.g., [19, 50, 84, 85, 143, 157, 172, 173, 270, 284, 294], real-world planar graphs have been overlooked. Secondly, real-world networks tend to have small vertex separators. By the planar separator theorem [287] and the planarity of RANs, this property is satisfied. This should be seen in contrast to the popular preferential attachment model [50] where the generated graph is an expander.

Contributions:

In [179, 183] we perform the first rigorous analysis of RANs.

- We prove in Theorem 4.1 tight results on the degree sequence of RANs. Previous results were weaker or even erroneous, see [420, 425].
- We prove in Theorem 4.3 tight asymptotic expressions for the top- k largest degrees, where k is constant.
- We provide in Theorem 4.4 tight asymptotic expressions for the top- k largest eigenvalues, where k is constant.
- We provide a simple first moment argument that upper-bounds the asymptotic diameter growth. By observing a bijection between RANs and random ternary trees, we are able to prove Theorem 4.5, a refined upper bound on the diameter.

Triangle Counting (Chapter 5)

We motivated the importance of triangles in real-world networks in Section 1.1.3.

Contributions

In Chapter 5 we present results from our work [257, 311, 325, 396].

- In Section 5.2 we present a randomized algorithm for approximately counting the number of triangles in a graph G . The algorithm proceeds as follows: keep each edge independently with probability p , enumerate the triangles in the sparsified graph G' and return the number of triangles found in G' multiplied by p^{-3} . We prove that under mild assumptions on G and p our algorithm returns a good approximation for the number of triangles with high probability. We illustrate the efficiency of our algorithm on various large real-world datasets where we achieve significant speedups. Furthermore, we investigate the performance of existing sparsification procedures namely the Spielman-Srivastava spectral sparsifier [370] and the the Benczúr-Karger cut sparsifier [63, 64] and show that they are not optimal/suitable with respect to triangle counting.

- In Section 5.3 we extend the results from Section 5.2 by introducing a powerful idea of Alon, Yuster and Zwick [25]. As a result, we propose a Monte Carlo algorithm which approximates the true number of triangles within accuracy $(1+\epsilon)$ and runs in $O\left(m + \frac{m^{3/2} \log n}{t\epsilon^2}\right)$ time, where $n, m, t, \epsilon > 0$ are the number of vertices, edges, triangles and a small constant respectively. We extend our method to the semi-streaming model [164] using three passes and a memory overhead of $O\left(m^{1/2} \log n + \frac{m^{3/2} \log n}{t\epsilon^2}\right)$. We propose a random projection based method for triangle counting and provide a sufficient condition to obtain an estimate with low variance.
- In Section 5.4 we present a new sampling approach to approximating the number of triangles in a graph $G(V, E)$, that significantly improves existing sampling approaches. Furthermore, it is easily implemented in parallel. The key idea of our algorithm is to correlate the sampling of edges such that if two edges of a triangle are sampled, the third edge is always sampled. Compared to Section 5.2, this sampling decreases the degree of the multivariate polynomial that expresses the number of sampled triangles. As a result, we are able to obtain more “aggressive” sampling techniques compared to Section 5.2, while strong concentration results remain valid.

Densest Subgraphs (Chapter 6)

Extracting dense subgraphs from large graphs is a key primitive in a variety of application domains [274]. In the Web graph, dense subgraphs may correspond to thematic groups or even spam link farms, as observed by Gibson et al. [193]. In biology, finding dense subgraphs can be used for discovering regulatory motifs in genomic DNA [177], and finding correlated genes [267], and detecting transcriptional modules [156]. In the financial domain, extracting dense subgraphs has been applied to finding price value motifs [141]. Other applications include graph compression [95], graph visualization [30], reachability and distance query indexing [226], and finding stories and events in micro-blogging streams [36].

Contributions In Chapter 6 we present results most of which are included in [385]. Our contributions are summarized as follows.

- We introduce a general framework for finding dense subgraphs, which subsumes popular density functions. We provide theoretical insights into our framework by showing that a large family of objectives are efficiently solvable but there also exist subcases which are **NP**-hard.

- Our framework provides a principled way to derive novel algorithms/heuristics geared towards the requirements of the application of interest. As a special instance of our general framework we introduce the problem of extracting **optimal quasi-clique**, which in general is **NP-hard**.
- We design two efficient algorithms for extracting optimal quasi-cliques. The first one is a greedy algorithm where the smallest-degree vertex is repeatedly removed from the graph, and achieves an additive approximation error. The second algorithm is a heuristic based on the local-search paradigm.
- For a shifted-version of our objective, we show that the problem can be approximated within a constant factor of 0.796 using a semidefinite-programming algorithm.
- We evaluate our efficient algorithms on numerous datasets, both synthetic and real, showing that it produces high quality dense subgraphs. In particular, in the synthetic data experiments, we plant a clique in Erdős-Rényi and in random power-law graphs, and measure precision and recall of the methods in “recovering” the planted clique: our method clearly outperforms the **densest subgraph** in this task. We also develop applications of our method in data mining and bioinformatic tasks, such as forming a successful team of domain experts and finding highly correlated genes from a microarray dataset.
- Finally, motivated by real-world scenarios, we define and evaluate interesting variants of our original problem definition, such as (i) finding the top- k **optimal quasi-cliques**, and (ii) finding **optimal quasi-cliques** that contain a given set of vertices.

Structure of the Web Graph (Chapter 7)

The Web graph describes the directed links between pages of the World Wide Web (WWW) [82, 90]. It is a graph which occupies a special position among real-world networks. The World Wide Web grew in a decentralized way, under the influence and decision of multiple participants. Understanding the structure of WWW and developing realistic models of evolution are two major research problems that have attracted a lot of interest [81].

Contributions

In Chapter 7 we present results from our work [237, 239]. Our contributions can be summarized as follows:

The key contributions of this Chapter are the following:

- We propose HADI, a scalable algorithm to compute the radii and diameter of network.

- We analyze one of the largest public Web graphs, with several *billions* of nodes and edges. We validate the small-world phenomenon and find several new structural patterns.

FENNEL: Streaming Graph Partitioning for Massive Scale Graphs (Chapter 8)

A key computational problem underlying all existing large-scale graph-processing platforms is *balanced graph partitioning*. When the graph is partitioned, the sizes of the partitions have to be balanced to exploit the speedup of parallel computing over different partitions. Furthermore, it is critical that the number of edges between distinct partitions is small in order to minimize the communication cost incurred due to messages exchanged between different partitions. Pregel [295], Apache Giraph [3], PEGASUS [236] and GraphLab [291] use as a default partitioner hash partitioning on vertices, which essentially corresponds to assigning each vertex to one of the k partitions uniformly at random. This heuristic would balance the number of vertices per partition, but as it is entirely oblivious to the graph structure, may well result in grossly suboptimal fraction of edges cut. Balanced graph partitioning becomes even harder in the case of dynamic graphs: whenever a new edge or a new vertex with its neighbors arrives, it has to be assigned to one of partition parts.

Summary of our Contributions

In Chapter 8 we present results from our work [384]. Our contributions can be summarized in the following points:

- We introduce a general framework for graph partitioning that relaxes the hard cardinality constraints on the number of vertices in a cluster [37, 260]. Our formulation provides a unifying framework that subsumes two of the most popular heuristics used for streaming balanced graph partitioning: the folklore heuristic of [344] which places a vertex to the cluster with the fewest non-neighbors, and the degree-based heuristic of [373], which serves as the current state-of-the-art method with respect to performance.
- Our framework allows us to define formally the notion of interpolation between the non-neighbors heuristic [344] and the neighbors heuristic [373]. This provides improved performance for the balanced partitioning problem in the streaming setting.
- We evaluate our proposed streaming graph partitioning method, Fennel, on a wide range of graph datasets, both real-world and synthetic graphs, showing that it produces high quality graph partitions. Table 8.1 shows the performance of Fennel

versus the best previously-known heuristic, which is the linear weighted degrees [373], and the baseline Hash Partition of vertices. We observe that Fennel achieves simultaneously significantly smaller fraction of edges cut and balanced cluster sizes.

- We also demonstrate the performance gains with respect to communication cost and run time while running iterative computations over partitioned input graph data in a distributed cluster of machines. Specifically, we evaluated Fennel and other partitioning methods by computing PageRank in the graph processing platform Apache Giraph. We observe significant gains with respect to the byte count among different clusters and run time in comparison with the baseline Hash Partition of vertices.
- Furthermore, modularity—a popular measure for community detection [196, 314, 316]—is also a special instance of our framework. We establish an approximation algorithm for a shifted objective, achieving a guarantee of $O(\log(k)/k)$ for partitioning into k clusters.

PEGASUS: A System for Large-Scale Graph Processing (Chapter 9)

As we discussed previously, an appealing solution to improve upon scalability is to partition massive graphs into smaller partitions and then use a large distributed system to process them. Designing and implementing efficient graph processing platforms is a major problem.

Contributions

In Chapter 9 we present results from our work [236, 240]. Our main contributions are the following:

- We introduce a generic framework which allows us to perform various important graph mining tasks efficiently by generalizing the standard matrix-vector multiplication (GIM-V).
- We implement PEGASUS, an optimized graph mining library. The source code is available online at <http://www.cs.cmu.edu/~pegasus/>.
- We analyze the performance of our system showing that it scales well to large-scale graphs.
- We apply PEGASUS on several large, real-world networks and we obtain insights into their structure.

Approximation Algorithms for Speeding up Dynamic Programming and Denoising aCGH data (Chapter 10)

As we previously discussed in Section 1.2.1, a major computational problem in cancer biology is assessing DNA copy number variations in individual organisms or tissues.

Contributions

In Chapter 10 we present results from our work [304, 397]. Our contributions can be summarized as follows:

- We propose a new formulation of the array Comparative Genomic Hybridization (aCGH) denoising problem. Specifically, based on the well-established observation that near-by probes tend to have the same DNA copy number, we formulate the problem of denoising aCGH data as the problem of approximating a signal P with another signal F consisting of a few piecewise constant segments. Specifically, let $P = (P_1, P_2, \dots, P_n) \in \mathbb{R}^n$ be the input signal -in our setting the sequence of the noisy aCGH measurements- and let C be a constant. Our goal is to find a function $F : [n] \rightarrow \mathbb{R}$ which optimizes the following objective function:

$$\min_F \sum_{i=1}^n (P_i - F_i)^2 + C \times (|\{i : F_i \neq F_{i+1}\}| + 1). \quad (1.3)$$

- We solve the problem using a dynamic programming algorithm in $O(n^2)$ time.
- We provide a technique which approximates the optimal value of our objective function within additive ϵ error and runs in $\tilde{O}(n^{\frac{4}{3}+\delta} \log(\frac{U}{\epsilon}))$ time, where δ is an arbitrarily small positive constant and $U = \max\{\sqrt{C}, (|P_i|)_{i=1, \dots, n}\}$.
- We provide a technique for approximate dynamic programming which solves the corresponding recurrence within a multiplicative factor of $(1+\epsilon)$ and runs in $O(n \log n/\epsilon)$.
- We validate our proposed model on both synthetic and real data. Specifically, our segmentations result in superior precision and recall compared to leading competitors on benchmarks of synthetic data and real data from the Coriell cell lines. In addition, we are able to find several novel markers not recorded in the benchmarks but supported in the oncology literature.

Robust Unmixing of Tumor States in Array Comparative Genomic Hybridization Data (Chapter 11)

We discussed in Section 1.2 the phenomenon of *inter-tumor heterogeneity*. We propose a geometric approach robust to noise to the problem of detecting cancer subtypes.

Contributions

In Chapter 11 we present results from our work [382]. Our contributions can be summarized as follows:

- We introduce a novel method for finding cancer subtypes using tissue-wide DNA copy number data as assessed by array comparative genomic hybridization (aCGH) data.
- We develop efficient computational tools to solve our optimization problem which is robust to noise.
- We apply our method to an aCGH data set taken from [312] and show that the method identifies state sets corresponding to known subtypes consistent with much of the analysis performed by the authors.

Perfect Reconstruction of Oncogenetic Trees (Chapter 12)

Human cancer is caused by the accumulation of genetic alternations in cells [43, 413]. Finding driver genetic mutations, i.e., mutations which confer growth advantage on the cells carrying them and have been positively selected during the evolution of the cancer and uncovering their temporal sequence have been central goals of cancer research the last decades [376]. Among the triumphs of cancer research stands the breakthrough work of Vogelstein and his collaborators [159, 405] which provides significant insight into the evolution of colorectal cancer. Specifically, the so-called “Vogelgram” models colorectal tumorigenesis as a linear accumulation of certain genetic events. Few years later, Desper et al. [133] considered more general evolutionary models compared to the “Vogelgram” and presented one of the first theoretical approaches to the problem [43], the so-called *oncogenetic trees*.

Oncogenetic trees have been a successful tumorigenesis model for various cancer types. For this reason understanding its properties is an important problem.

Contributions

In Chapter 12 we present results from our work [386]. Our main contribution is the following:

- We provide necessary and sufficient conditions for the unique reconstruction of an oncogenetic tree [133].

It is worth outlining that these conditions may be used to understand better the phenomenon of *intra-tumor heterogeneity* [392].

Conclusion and Future Directions (Chapter 13)

Our work leaves numerous interesting problems open. In Chapter 13 we conclude and provide several new research directions.

Chapter 2

Theoretical Preliminaries

In this Chapter we review theoretical preliminaries that are used in later Chapters.

2.1 Concentration of Measure Inequalities

The use of Chebyshev's inequality is known as *the second moment method*.

Lemma 2.1 (Chebyshev's Inequality [26]). *Let X be a random variable, $\mu = \mathbb{E}[X]$, $\sigma = \sqrt{\text{Var}[X]}$. For any positive $\lambda > 0$*

$$\Pr[|X - \mu| \geq \lambda\sigma] \leq \frac{1}{\lambda^2}.$$

Chernoff bounds allow us to obtain strong concentration results, when applicable. We use the following version in later Chapters.

Lemma 2.2 (Chernoff Inequality [26]). *Let X_1, X_2, \dots, X_k be independently distributed $\{0, 1\}$ variables with $E[X_i] = p$. Then for any $\epsilon > 0$, we have*

$$\Pr\left[\left|\frac{1}{k}\sum_{i=1}^k X_i - p\right| > \epsilon p\right] \leq 2e^{-\epsilon^2 pk/2}$$

The theory of discrete time martingales [26] will be the key to establish concentration inequalities in our proofs for degree sequences.

Lemma 2.3 (Azuma-Hoeffding inequality). *Let $\lambda > 0$. Also, let $(X_t)_{t=0}^n$ be a martingale sequence with $|X_{t+1} - X_t| \leq c$ for $t = 0, \dots, n-1$. Then:*

$$\Pr[|X_n - X_0| \geq \lambda] \leq \exp\left(-\frac{\lambda^2}{2c^2n}\right).$$

The Kim-Vu theorem is an important concentration result since it allows us to obtain strong concentration when the polynomial of interest is not *smooth*. Specifically, for the purposes of our work, let $Y = Y(t_1, \dots, t_m)$ be a positive polynomial of m Boolean variables $[t_i]_{i=1..m}$ which are independent. A common task in combinatorics is to show that Y is concentrated around its expected value. In the following we state the necessary definitions and the main concentration result which we will use in our method. Y is totally positive if all of its coefficients are non-negative variables. Y is homogeneous if all of its monomials have the same degree and we call this value the degree of the polynomial. Given any multi-index $\alpha = (\alpha_1, \dots, \alpha_m) \in \mathbb{Z}_+^m$, define the partial derivative $\partial^\alpha Y = (\frac{\partial}{\partial t_1})^{\alpha_1} \dots (\frac{\partial}{\partial t_m})^{\alpha_m} Y(t_1, \dots, t_m)$ and denote by $|\alpha| = \alpha_1 + \dots + \alpha_m$ the order of α . For any order $d \geq 0$, define $\mathbb{E}_d(Y) = \max_{\alpha: |\alpha|=d} \mathbb{E}(\partial^\alpha Y)$ and $\mathbb{E}_{\geq d}(Y) = \max_{d' \geq d} \mathbb{E}_{d'}(Y)$.

Now, we refer to the main theorem of Kim and Vu of [249, §1.2] as phrased in Theorem 1.1 of [407] or as Theorem 1.36 of [378].

Theorem 2.4. *There is a constant c_k depending on k such that the following holds. Let $Y(t_1, \dots, t_m)$ be a totally positive polynomial of degree k , where t_i can have arbitrary distribution on the interval $[0, 1]$. Assume that:*

$$\mathbb{E}[Y] \geq \mathbb{E}_{\geq 1}(Y) \tag{2.1}$$

Then for any $\lambda \geq 1$:

$$\Pr\left[|Y - \mathbb{E}[Y]| \geq c_k \lambda^k (\mathbb{E}[Y] \mathbb{E}_{\geq 1}(Y))^{1/2}\right] \leq e^{-\lambda + (k-1) \log m}. \tag{2.2}$$

Typically, when a polynomial Y is smooth, it is strongly concentrated. By smoothness one usually means a small Lipschitz coefficient or in other words, when one changes the value of one variable t_j , the value Y changes no more than a constant. However, as stated in [407] this is restrictive in many cases. Thus one can demand “average smoothness” as defined in [407] which is quantified via the expectation of partial derivatives of any order.

2.2 Random Projections

A random projection $x \rightarrow Rx$ from $\mathbb{R}^d \rightarrow \mathbb{R}^k$ approximately preserves all Euclidean distances. One version of the Johnson-Lindenstrauss lemma [227] is the following:

Lemma 2.5 (Johnson Lindenstrauss). *Suppose $x_1, \dots, x_n \in \mathbb{R}^d$ and $\epsilon > 0$ and take $k = C\epsilon^{-2} \log n$. Define the random matrix $R \in \mathbb{R}^{k \times d}$ by taking all $R_{i,j} \sim N(0, 1)$ (standard gaussian) and independent. Then, with probability bounded below by a constant the points $y_j = Rx_j \in \mathbb{R}^k$ satisfy*

$$(1 - \epsilon)|x_i - x_j| \leq |y_i - y_j| \leq (1 + \epsilon)|x_i - x_j|$$

for $i, j = 1, 2, \dots, n$ where $|\cdot|$ represents the Euclidean norm.

2.3 Extremal Graph Theory

Hajnal and Szemerédi [208] proved in 1970 the following conjecture of Paul Erdős:

Theorem 2.6 (Hajnal-Szemerédi Theorem). *Every graph with n vertices and maximum vertex degree at most k is $k + 1$ colorable with all color classes of size $\lfloor \frac{n}{k+1} \rfloor$ or $\lceil \frac{n}{k+1} \rceil$.*

Ahlsvede and Katona consider the following problem: which graph with a given number of vertices n and a given number of edges m maximizes the number of edges in its line graph $L(G)$? The problem is equivalent to maximizing the sum of squares of the degrees of the vertices under the constraint that their sum equals twice the number of the edges. The following theorem was given in [18] and answers this question.

Lemma 2.7 (Ahlsvede-Katona theorem). *The maximum value of the sum of the squares of all vertex degrees $\sum_{v \in V(G)} d(v)^2$ over the set of all graphs with n vertices and m edges occurs at one or both of two special types of graphs, the quasi-star graph or the quasi-complete graph.*

For further progress on other questions related to the above optimization problem such as when does the optimum occur at both graphs, see the work of Abrego, Fernández-Merchant, Neubauer and Watkins [9].

2.4 Two useful lemmas

Two more useful lemmas we will use in later chapters follow.

Lemma 2.8 (Lemma 3.1, [114]). *Suppose that a sequence $\{a_t\}$ satisfies the recurrence*

$$a_{t+1} = \left(1 - \frac{b_t}{t + t_1}\right)a_t + c_t$$

for $t \geq t_0$. Furthermore suppose $\lim_{t \rightarrow +\infty} b_t = b > 0$ and $\lim_{t \rightarrow +\infty} c_t = c$. Then $\lim_{t \rightarrow +\infty} \frac{a_t}{t}$ exists and

$$\lim_{t \rightarrow +\infty} \frac{a_t}{t} = \frac{c}{1 + b}.$$

Graphs can be viewed as electrical networks. Given two vertices $s, t \in V(G)$, we can ensure an electrical current from s to t of value 1. The potential/voltage difference between s and t is defined to be the *effective resistance* $R(s, t)$. For further details the interested reader should read [76]. The following theorem is due to Foster [175].

Theorem 2.9 (Foster's theorem [175]). *Let G be a connected graph of order n . Then*

$$\sum_{(u,v) \in E(G)} R(u, v) = n - 1.$$

2.5 Semidefinite Bounds

Semidefinite programs are generalizations of linear programs which can be solved in polynomial time using interior point methods. Semidefinite programming uses symmetric, positive semidefinite matrices.

Definition 2.10. A matrix $A \in \mathbb{R}^{n \times n}$ is positive semidefinite if and only if for all $x \in \mathbb{R}^n$, $x^T A x \geq 0$.

We define the scalar product $\langle A, B \rangle$ of matrices $A, B \in \mathbb{R}^{n \times n}$ as $\langle A, B \rangle = \text{Trace}(B^T A)$ which is an inner product on the vector space of $n \times n$ matrices. A semidefinite program is defined by the symmetric $n \times n$ matrices C, A_1, \dots, A_m and vector $b \in \mathbb{R}^m$ as follows:

$$\begin{aligned} & \mathbf{max} && \langle C, X \rangle \\ & \text{subject to} && \langle A_i, X \rangle, \text{ for all } i \in \{1, \dots, m\} \\ & && \text{and } X \succeq 0, X \text{ symmetric.} \end{aligned} \tag{2.3}$$

Goemans and Williamson significantly advanced the field of approximation algorithms by introducing a randomized rounding for the MAX-CUT problem[198]. The MAX-CUT problem takes as input a graph $G(V, E, w)$, $w : E \rightarrow \mathbb{R}^+$ and asks for a non-empty set S such that weight $\sum_{(i,j) \in E(G)} w_{ij}$ of the cut $(S, V \setminus S)$ is maximized. This problem is **NP**-hard. The first step of the Goemans-Williamson algorithm is a semidefinite relaxation of the following quadratic integer program which is equivalent to MAX-CUT.

$$p^* = \max \frac{1}{2} \sum_{(i,j) \in E} w_{ij}(1 - x_i x_j) \quad \text{subject to } x \in \{-1, +1\}^n.$$

Notice that $X = xx^T$ is a symmetric, positive semidefinite matrix with rank 1 and $X_{ii} = 1$ for $i \in [n]$. The semidefinite relaxation relaxes the rank 1 condition as follows

$$s^* = \max \frac{1}{2} \sum_{(i,j) \in E} w_{ij}(1 - X_{ij}) \quad \text{subject to } X \succeq 0, X \text{ symmetric.}, X_{ii} = 1 \text{ for all } i.$$

The second step of the Goemans-Williamson algorithm consists of a randomized rounding procedure. Specifically, let X be the optimal solution of the semidefinite relaxation. As X is positive semidefinite $X = U^T U$ where $U \in \mathbb{R}^{d \times n}$, $d \leq n$. Let u_j be the j -th column of U . As $X_{ii} = 1$ for all i , $\|u_i\|_2 = 1$. Goemans and Williamson proposed generating a random unit vector $r \in \mathbb{R}^d$ and letting S be the set of vertices which correspond to columns of U such that $u_j^T r > 0$. They proved that this algorithm provides a β -approximation where $\beta > 0.87856$. Their technique has been extended to the MAX-k-CUT problem [181].

2.6 Speeding up Dynamic Programming

Dynamic programming is a powerful problem solving technique introduced by Bellman [62] with numerous applications in biology, e.g., [211, 339, 411], in control theory, e.g., [66], in operations research and many other fields. Due to its importance, a lot of research has focused on speeding up basic dynamic programming implementations. A successful example of speeding up a naive dynamic programming implementation is the computation of optimal binary search trees. Gilbert and Moore solved the problem efficiently using dynamic programming [195]. Their algorithm runs in $O(n^3)$ time and for several years this running time was considered to be tight. In 1971 Knuth [255] showed that the same computation can be carried out in $O(n^2)$ time. This remarkable result

was generalized by Frances Yao in [422, 423]. Specifically, Yao showed that this dynamic programming speedup technique works for a large class of recurrences. She considered the recurrence $c(i, i) = 0$, $c(i, j) = \min_{i < k \leq j} (c(i, k - 1) + c(k, j)) + w(i, j)$ for $i < j$ where the weight function w satisfies the quadrangle inequality (see Section 2.8) and proved that the solution of this recurrence can be found in $O(n^2)$ time. Eppstein, Galil and Giancarlo have considered similar recurrences where they showed that naive $O(n^2)$ implementations of dynamic programming can run in $O(n \log n)$ time [148]. Larmore and Schieber [268] further improved the running time, giving a linear time algorithm when the weight function is concave. Klawe and Kleitman give in [251] an algorithm which runs in $O(n\alpha(n))$ time when the weight function is convex, where $\alpha(\cdot)$ is the inverse Ackermann function. Furthermore, Eppstein, Galil, Giancarlo and Italiano have also explored the effect of sparsity [149, 150], another key concept in speeding up dynamic programming. Aggarwal, Klawe, Moran, Shor, Wilber developed an algorithm, widely known as the SMAWK algorithm, [16] which can compute in $O(n)$ time the row maxima of a totally monotone $n \times n$ matrix. The connection between the Knuth-Yao technique and the SMAWK algorithm was made clear in [60], by showing that the Knuth-Yao technique is a special case of the use of totally monotone matrices. The basic properties which allow these speedups are the convexity or concavity of the weight function. Such properties date back to Monge [307] and are well studied in the literature, see for example [97].

Close to our work lies the work on histogram construction, an important problem for database applications. Jagadish et al. [222] originally provided a simple dynamic programming algorithm which runs in $O(kn^2)$ time, where k is the number of buckets and n the input size and outputs the best V-optimal histogram. Guha, Koudas and Shim [204] propose a $(1 + \epsilon)$ approximation algorithm which runs in linear time. Their algorithm exploits monotonicity properties of the key quantities involved in the problem. Our $(1 + \epsilon)$ approximation algorithm in Section 2.8 uses a decomposition technique similar to theirs.

2.7 Reporting Points in a Halfspace

Let S be a set of points in \mathbb{R}^d and let k denote the size of the output, i.e., the number of points to be reported. Consider the problem of preprocessing S such that for any halfspace query γ we can report efficiently whether the set $S \cap \gamma$ is empty or not. This problem is a well studied special case of the more general range searching problem. For an extensive survey see the work by Agarwal and Erickson [14]. For $d = 2$, the problem has been solved optimally by Chazelle, Guibas and Lee [107]. For $d = 3$, Chazelle

and Preparata in [106] gave a solution with nearly linear space and $O(\log n + k)$ query time, while Aggarwal, Hansen and Leighton [17] gave a solution with a more expensive preprocessing but $O(n \log n)$ space. When the number of dimensions is greater than 4, i.e., $d \geq 4$, Clarkson and Shor [117] gave an algorithm that requires $O(n^{\lfloor d/2 \rfloor + \epsilon})$ preprocessing time and space, where ϵ is an arbitrarily small positive constant, but can subsequently answer queries in $O(\log n + k)$ time. Matoušek in [299] provides improved results on the problem, which are used by Agarwal, Eppstein, Matoušek [15] in order to create dynamic data structures that trade off insertion and query times. We refer to Theorem 2.1(iii) of their paper [15]:

Theorem 2.11 (Agarwal, Eppstein, Matoušek [15]). *Given a set S of n points in \mathbb{R}^d where $d \geq 3$ and a parameter m between n and $n^{\lfloor \frac{d}{2} \rfloor}$ the halfspace range reporting problem can be solved with the following performance: $O(\frac{n}{m^{\lfloor \frac{d}{2} \rfloor}} \log n)$ query time, $O(m^{1+\epsilon})$ space and preprocessing time, $O(m^{1+\epsilon}/n)$ amortized update time.*

Substituting for $d = 4$, $m = n^{\frac{4}{3}}$ we obtain the following corollary, which will be used as a subroutine in our proposed method:

Corollary 2.12. *Given a set S of n points in \mathbb{R}^4 the halfspace range reporting problem can be solved with $O(n^{\frac{1}{3}} \log n)$ query time, $O(n^{\frac{4}{3}+\delta})$ space and preprocessing time, and $O(n^{\frac{1}{3}+\delta})$ update time, where δ is an arbitrarily small positive constant.*

2.8 Monge Functions and Dynamic Programming

Here, we refer to one of the results in [268] which we use in Section 2.8 as a subroutine for our proposed method. A function w defined on pairs of integer indices is Monge (concave) if for any 4-tuple of indices $i_1 < i_2 < i_3 < i_4$, $w(i_1, i_4) + w(i_2, i_3) \geq w(i_1, i_3) + w(i_2, i_4)$. Furthermore, we assume that f is a function such that the values $f(a_j)$ for all j are easily evaluated. The following results holds:

Theorem 2.13 ([268]). *Consider the one dimensional recurrence $a_i = \min_{j < i} \{f(a_j) + w(j, i)\}$ for $i = 1, \dots, n$, where the basis a_0 is given. There exists an algorithm which solves the recurrence online in $O(n)$ time¹.*

¹Thus, obtaining $O(n)$ speedup compared to the straight-forward dynamic programming algorithm which runs in $O(n^2)$ units of time.

I *Graphs and Networks*

Chapter 3

Rainbow Connectivity of Sparse Random Graphs

3.1 Introduction

Connectivity is a fundamental graph theoretic property. Recently, the concept of *rainbow connectivity* was introduced by Chartrand et al. in [105]. An edge colored graph G is rainbow edge connected if any two vertices are connected by a path whose edges have distinct colors. The rainbow connectivity $rc(G)$ of a connected graph G is the smallest number of colors that are needed in order to make G rainbow edge connected. Notice, that by definition a rainbow edge connected graph is also connected and furthermore any connected graph has a trivial edge coloring that makes it rainbow edge connected, since one may color the edges of a given spanning tree with distinct colors. Other basic facts established in [105] are that $rc(G) = 1$ if and only if G is a clique and $rc(G) = |V(G)| - 1$ if and only if G is a tree. Besides its theoretical interest, rainbow connectivity is also of interest in applied settings, such as securing sensitive information [285], transfer and networking [101].

The concept of rainbow connectivity has attracted the interest of various researchers. Chartrand et al. [105] determine the rainbow connectivity of several special classes of graphs, including multipartite graphs. Caro et al. [100] prove that for a connected graph G with n vertices and minimum degree δ , the rainbow connectivity satisfies $rc(G) \leq \frac{\log \delta}{\delta} n(1 + f(\delta))$, where $f(\delta)$ tends to zero as δ increases. The following simpler bound was also proved in [100], $rc(G) \leq n^{\frac{4 \log n + 3}{\delta}}$. Krivelevich and Yuster [261] removed the logarithmic factor from the Caro et al. [100] upper bound. Specifically they proved that $rc(G) \leq \frac{20n}{\delta}$. Due to a construction of a graph with minimum degree δ and diameter $\frac{3n}{\delta+1} - \frac{\delta+7}{\delta+1}$ by Caro et al. [100], the best upper bound one can hope for is $rc(G) \leq \frac{3n}{\delta}$.

Chandran, Das, Rajendraprasad and Varma [103] have subsequently proved an upper bound of $\frac{3n}{\delta+1} + 3$, which is therefore essentially optimal.

As Caro et al. point out, the random graph setting poses several intriguing questions. Specifically, let $G = G(n, p)$ denote the binomial random graph on n vertices with edge probability p [153]. Caro et al. [100] proved that $p = \sqrt{\log n/n}$ is the sharp threshold for the property $rc(G(n, p)) \leq 2$. He and Liang [209] studied further the rainbow connectivity of random graphs. Specifically, they obtain the sharp threshold for the property $rc(G) \leq d$ where d is constant. For further results and references we refer the interested reader to the recent monograph of Li and Sun [285]. In this work we look at the rainbow connectivity of the binomial graph at the connectivity threshold $p = \frac{\log n + \omega}{n}$ where $\omega = o(\log n)$. This range of values for p poses problems that cannot be tackled with the techniques developed in the aforementioned work. Rainbow connectivity has not been studied in random regular graphs to the best of our knowledge.

Let

$$L = \frac{\log n}{\log \log n} \tag{3.1}$$

and let $A \sim B$ denote $A = (1 + o(1))B$ as $n \rightarrow \infty$.

We establish the following theorems:

Theorem 3.1. *Let $G = G(n, p)$, $p = \frac{\log n + \omega}{n}$, $\omega \rightarrow \infty$, $\omega = o(\log n)$. Also, let Z_1 be the number of vertices of degree 1 in G . Then, with high probability (whp)*

$$rc(G) \sim \max \{Z_1, L\},$$

It is known that whp the diameter of $G(n, p)$ is asymptotic to L for p as in the above range, see for example Theorem 10.17 of Bollobás [76]. Theorem 3.1 gives asymptotically optimal results. Our next theorem is not quite as precise.

Theorem 3.2. *Let $G = G(n, r)$ be a random r -regular graph where $r \geq 3$ is a fixed integer. Then, whp*

$$rc(G) = \begin{cases} O(\log^4 n) & r = 3 \\ O(\log^{2\theta_r} n) & r \geq 4. \end{cases}$$

where $\theta_r = \frac{\log(r-1)}{\log(r-2)}$.

All logarithms whose base is omitted are natural. It will be clear from our proofs that the colorings in the above two theorems can be constructed in a low order polynomial time. The second theorem, while weaker, contains an unexpected use of a Markov Chain Monte-Carlo (MCMC) algorithm for randomly coloring a graph.

The Chapter is organized as follows: After giving a sketch of our approach in Section 3.2, in Sections 3.3, 3.4 we prove Theorems 3.1, 3.2 respectively.

3.2 Sketch of approach

The general idea in the proofs of both theorems is as follows:

1. Randomly color the edges of the graph in question. For Theorem 3.1 we can (in the main) use a uniformly random coloring. The distribution for Theorem 3.2 is a little more complicated.
2. To prove that this works, we have to find, for each pair of vertices x, y , a large collection of edge disjoint paths joining them. It will then be easy to argue that at least one of these paths is rainbow colored.
3. To find these paths we pick a typical vertex x . We grow a regular tree T_x with root x . The depth is chosen carefully. We argue that for a typical pair of vertices x, y , many of the leaves of T_x and T_y can be put into 1-1 correspondence f so that (i) the path P_x from x to leaf v of T_x is rainbow colored, (ii) the path P_y from y to the leaf $f(v)$ of T_y is rainbow colored and (iii) P_x, P_y do not share color.
4. We argue that from most of the leaves of T_x, T_y we can grow a tree of depth approximately equal to half the diameter. These latter trees themselves contain a bit more than $n^{1/2}$ leaves. These can be constructed so that they are vertex disjoint. Now we argue that each pair of trees, one associated with x and one associated with y , are joined by an edge.
5. We now have, by construction, a large set of edge disjoint paths joining leaves v of T_x to leaves $f(v)$ of T_y . A simple estimation shows that *whp* for at least one leaf v of T_x , the path from v to $f(v)$ is rainbow colored and does not use a color already used in the path from x to v in T_x or the path from y to $f(v)$ in T_y .

We now fill in the details of both cases.

3.3 Proof of Theorem 3.1

Observe first that $rc(G) \geq \max\{Z_1, \text{diameter}(G)\}$. First of all, each edge incident to a vertex of degree one must have a distinct color. Just consider a path joining two such vertices. Secondly, if the shortest distance between two vertices is ℓ then we need at

least ℓ colors. Next observe that *whp* the diameter D is asymptotically equal to L , see for example [76]. We break the proof of Theorem 3.1 into several lemmas.

Let a vertex be *large* if $d(x) \geq \log n/100$ and *small* otherwise.

Lemma 3.3. *Whp, there do not exist two small vertices within distance at most $3L/4$.*

Proof.

$$\begin{aligned}
& \Pr \left[\exists x, y \in [n] : d(x), d(y) \leq \log n/100 \text{ and } \text{dist}(x, y) \leq \frac{3L}{4} \right] \\
& \leq \binom{n}{2} \sum_{k=1}^{3L/4} n^{k-1} p^k \left(\sum_{i=0}^{\log n/100} \binom{n-1-k}{i} p^i (1-p)^{n-1-k} \right)^2 \\
& \leq \sum_{k=1}^{3L/4} n (2 \log n)^k \left(2 \binom{n}{\log n/100} p^{\log n/100} (1-p)^{n-1-\log n/100} \right)^2 \\
& \leq \sum_{k=1}^{3L/4} n (2 \log n)^k \left(2(100e^{1+o(1)})^{\log n/100} n^{-1+o(1)} \right)^2 \\
& \leq \sum_{k=1}^{3L/4} n (2 \log n)^k n^{-1.9} \\
& \leq 2n (2 \log n)^{3L/4} n^{-1.9} \\
& \leq n^{-1}.
\end{aligned}$$

□

We use the notation $e[S]$ for the number of edges induced by a given set of vertices S . Notice that if a set S satisfies $e[S] \geq s + t$ where $t \geq 1$, the induced subgraph $G[S]$ has at least $t + 1$ cycles.

Lemma 3.4. *Fix $t \in \mathbb{Z}^+$ and $0 < \alpha < 1$. Then, whp there does not exist a subset $S \subseteq [n]$, such that $|S| \leq \alpha t L$ and $e[S] \geq |S| + t$.*

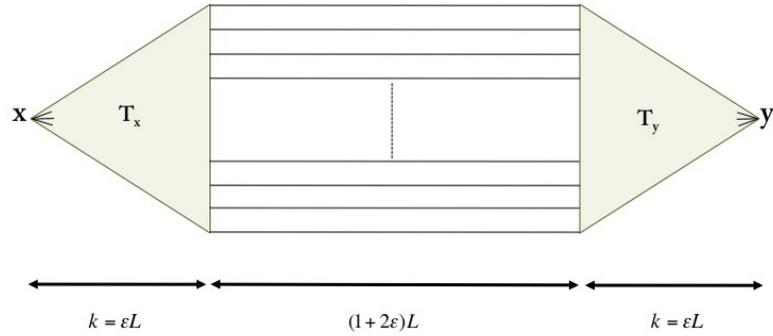


FIGURE 3.1: Structure of Lemma 3.6.

Proof. For convenience, let $s = |S|$ be the cardinality of the set S . Then,

$$\begin{aligned}
\Pr[\exists S : s \leq \alpha t L \text{ and } e[S] \geq s + t] &\leq \sum_{s \leq \alpha t L} \binom{n}{s} \binom{\binom{s}{2}}{s+t} p^{s+t} \\
&\leq \sum_{s \leq \alpha t L} \left(\frac{ne}{s}\right)^s \left(\frac{es^2 p}{2(s+t)}\right)^{s+t} \\
&\leq \sum_{s \leq \alpha t L} (e^{2+o(1)} \log n)^s \left(\frac{es \log n}{n}\right)^t \\
&\leq \alpha t L \left((e^{2+o(1)} \log n)^{\alpha L} \left(\frac{e \alpha t \log^2 n}{n \log \log n}\right)^t \right) \\
&< \frac{1}{n^{(1-\alpha-o(1))t}}.
\end{aligned}$$

□

Remark 3.5. Let T be a rooted tree of depth at most $4L/7$ and let v be a vertex not in T , but with b neighbors in T . Let S consist of v , the neighbors of v in T plus the ancestors of these neighbors. Then $|S| \leq 4bL/7 + 1 \leq 3bL/5$ and $e[S] = |S| + b - 2$. It follows from the proof of Lemma 3.4 with $\alpha = 3/5$ and $t = 8$, that we must have $b \leq 10$ with probability $1 - o(n^{-3})$.

Our next lemma shows the existence of the subgraph $G'_{x,y}$ described next and shown in Figure 3.1 for a given pair of vertices x, y . We first deal with paths between large vertices.

Now let

$$\epsilon = \epsilon(n) = o(1) \text{ be such that } \frac{\epsilon \log \log n}{\log 1/\epsilon} \rightarrow \infty \text{ and let } k = \epsilon L. \quad (3.2)$$

Here L is defined in (3.1) and we could take $\epsilon = 1/(\log \log n)^{1/2}$.

Lemma 3.6. *Whp, for all pairs of large vertices $x, y \in [n]$ there exists a subgraph $G_{x,y}(V_{x,y}, E_{x,y})$ of G as shown in figure 3.1. The subgraph consists of two isomorphic*

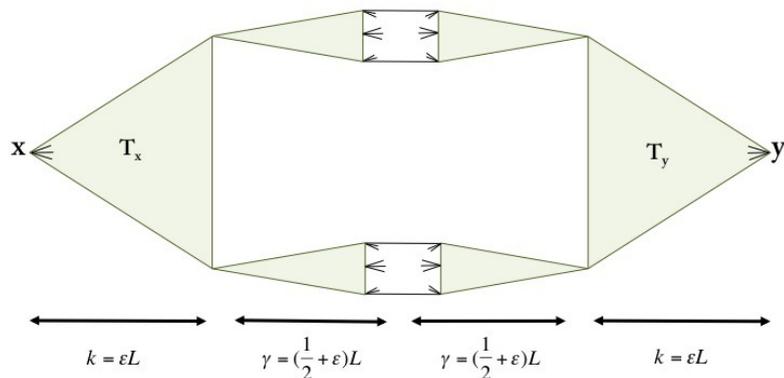


FIGURE 3.2: Subgraph found in the proof of Lemma 3.6.

vertex disjoint trees T_x, T_y rooted at x, y each of depth k . T_x and T_y both have a branching factor of $\log n/101$. I.e. each vertex of T_x, T_y has at least $\log n/101$ neighbors, excluding its parent in the tree. Let the leaves of T_x be x_1, x_2, \dots, x_τ where $\tau \geq n^{4\epsilon/5}$ and those of T_y be y_1, y_2, \dots, y_τ . Then $y_i = f(x_i)$ where f is a natural isomorphism that preserves the parent-child relation. Between each pair of leaves $(x_i, y_i), i = 1, 2, \dots, \tau$ there is a path P_i of length $(1 + 2\epsilon)L$. The paths $P_i, i = 1, 2, \dots, \tau$ are edge disjoint.

Proof. Because we have to do this for all pairs x, y , we note without further comment that likely (resp. unlikely) events will be shown to occur with probability $1 - o(n^{-2})$ (resp. $o(n^{-2})$).

To find the subgraph shown in Figure 3.1 we grow tree structures as shown in Figure 3.2. Specifically, we first grow a tree from x using BFS until it reaches depth k . Then, we grow a tree starting from y again using BFS until it reaches depth k . Finally, we grow trees from the leaves of T_x and T_y using BFS for depth $\gamma = (\frac{1}{2} + \epsilon)L$. Now we analyze these processes. Since the argument is the same we explain it in detail for T_x and we outline the differences for the other trees. We use the notation $D_i^{(\rho)}$ for the number of vertices at depth i of the BFS tree rooted at ρ .

First we grow T_x . As we grow the tree via BFS from a vertex v at depth i to vertices at depth $i + 1$ certain *bad* edges from v may point to vertices already in T_x . Remark 3.5 shows with probability $1 - o(n^{-3})$ there can be at most 10 bad edges emanating from v .

Furthermore, Lemma 3.3 implies that there exists at most one vertex of degree less than $\frac{\log n}{100}$ at each level *whp*. Hence, we obtain the recursion

$$D_{i+1}^{(x)} \geq \left(\frac{\log n}{100} - 10 \right) (D_i^{(x)} - 1) \geq \frac{\log n}{101} D_i^{(x)}. \quad (3.3)$$

Therefore the number of leaves satisfies

$$D_k^{(x)} \geq \left(\frac{\log n}{101} \right)^{\epsilon L} \geq n^{4\epsilon/5}. \quad (3.4)$$

We can make the branching factor exactly $\frac{\log n}{101}$ by pruning. We do this so that the trees T_x are isomorphic to each other.

With a similar argument

$$D_k^{(y)} \geq n^{\frac{4}{5}\epsilon}. \quad (3.5)$$

The only difference is that now we also say an edge is bad if the other endpoint is in T_x . This immediately gives

$$D_{i+1}^{(y)} \geq \left(\frac{\log n}{100} - 20 \right) (D_i^{(y)} - 1) \geq \frac{\log n}{101} D_i^{(y)}$$

and the required conclusion (3.5).

Similarly, from each leaf $x_i \in T_x$ and $y_i \in T_y$ we grow trees $\widehat{T}_{x_i}, \widehat{T}_{y_i}$ of depth $\gamma = (\frac{1}{2} + \epsilon)L$ using the same procedure and arguments as above. Remark 3.5 implies that there are at most 20 edges from the vertex v being explored to vertices in any of the trees already constructed. At most 10 to T_x plus any trees rooted at an x_i and another 10 for y . The numbers of leaves of each \widehat{T}_{x_i} now satisfies

$$\widehat{D}_\gamma^{(x_i)} \geq \frac{\log n}{100} \left(\frac{\log n}{101} \right)^\gamma \geq n^{\frac{1}{2} + \frac{4}{5}\epsilon}.$$

Similarly for $\widehat{D}_\gamma^{(y_i)}$.

Observe next that BFS does not condition the edges between the leaves X_i, Y_i of the trees \widehat{T}_{x_i} and \widehat{T}_{y_i} . I.e., we do not need to look at these edges in order to carry out our construction. On the other hand we have conditioned on the occurrence of certain events to imply a certain growth rate. We handle this technicality as follows. We go through the above construction and halt if ever we find that we cannot expand by the required amount. Let \mathbf{A} be the event that we do not halt the construction i.e. we fail the conditions of Lemmas 3.3 or 3.4. We have $\Pr[\mathbf{A}] = 1 - o(1)$ and so,

$$\Pr[\exists i : e(X_i, Y_i) = 0 \mid \mathbf{A}] \leq \frac{\Pr[\exists i : e(X_i, Y_i) = 0]}{\Pr(\mathbf{A})} \leq 2n^{\frac{4\epsilon}{5}} (1-p)^{n^{1+\frac{8\epsilon}{5}}} \leq n^{-n^\epsilon}.$$

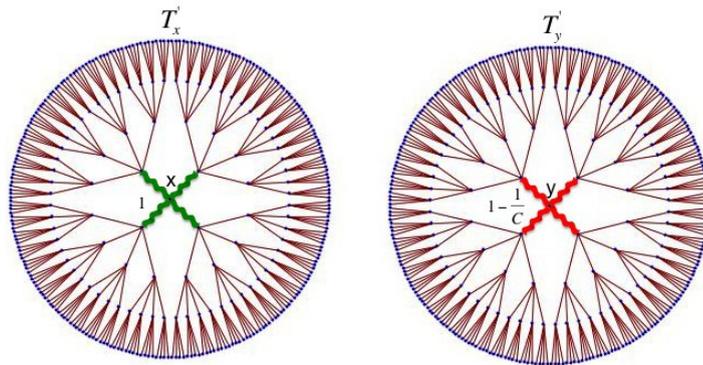


FIGURE 3.3: Figure shows $\frac{\log n}{101}$ -ary trees T_x, T_y . The two roots are shown respectively at the center of the trees. In our thinking of the random coloring as an evolutionary process, the green edges incident to x survive with probability 1, the red edges incident to y with probability $1 - \frac{1}{q}$ and all the other edges with probability $p_0 = \left(1 - \frac{2k}{q}\right)^2$ where k is the depth of both trees and q the number of available colors. Our analysis in Lemma 3.6 using these probabilities gives a lower bound on the number of alive pairs of leaves after coloring T_x, T_y from the root to the leaves respectively.

We conclude that *whp* there is always an edge between each X_i, Y_i and thus a path of length at most $(1 + 2\epsilon)L$ between each x_i, y_i . \square

Let $q = (1 + 5\epsilon)L$ be the number of available colors. We color the edges of G randomly. We show that the probability of having a rainbow path between x, y in the subgraph $G_{x,y}$ of Figure 3.1 is at least $1 - \frac{1}{n^3}$.

Lemma 3.7. *Color each edge of G using one color at random from q available. Then, the probability of having at least one rainbow path between two fixed large vertices $x, y \in [n]$ is at least $1 - \frac{1}{n^3}$.*

Proof. We show that the subgraph $G_{x,y}$ contains such a path. We break our proof into two steps:

Before we proceed, we provide certain necessary definitions. Think of the process of coloring T_x, T_y as an evolutionary process that colors edges by starting from the two roots $x, f(x) = y$ until it reaches the leaves. In the following, we call a vertex u of T_x (T_y) *alive/living* if the path $P(x, u)$ ($P(y, u)$) from x (y) to u is rainbow, i.e., the edges have received distinct colors. We call a pair of vertices $\{u, f(u)\}$ *alive*, $u \in T_x, f(u) \in T_y$ if $u, f(u)$ are both *alive* and the paths $P(x, u), P(y, f(u))$ share no color. Define $A_j = |\{(u, f(u)) : (u, f(u)) \text{ is alive and } \text{depth}(u) = j\}|$ for $j = 1, \dots, k$.

- STEP 1: Existence of at least $n^{\frac{4}{5}\epsilon}$ living pairs of leaves

Assume the pair of vertices $\{u, f(u)\}$ is alive where $u \in T_x, f(u) \in T_y$. It is worth noticing that $u, f(u)$ have the same depth in their trees. We are interested in the number of pairs of children $\{u_i, f(u_i)\}_{i=1, \dots, \log n/101}$ that will be alive after coloring the edges from $\text{depth}(u)$ to $\text{depth}(u) + 1$. A living pair $\{u_i, f(u_i)\}$ by definition has the following properties: edges $(u, u_i) \in E(T_x)$ and $(f(u), f(u_i)) \in E(T_y)$ receive two distinct colors, which are different from the set of colors used in paths $P(x, u)$ and $P(y, f(u))$. Notice the latter set of colors has cardinality $2 \times \text{depth}(u) \leq 2k$.

Let A_j be the number of living pairs at depth j . We first bound the size of A_1 .

$$\Pr \left[A_1 \leq \frac{\log n}{200} \right] \leq 2^{\log n/101} \left(\frac{1}{q} \right)^{\log n/300} = O(n^{-\Omega(\log \log n)}). \quad (3.6)$$

Here $2^{\log n/101}$ bounds the number of choices for A_1 . For a fixed set A_1 there will be at least $\frac{\log n}{101} - \frac{\log n}{200} \geq \frac{\log n}{300}$ edges incident with x that have the same color as their corresponding edges incident with y , under f . The factor $q^{-\log n/300}$ bounds the probability of this event.

For $j > 1$ we see that the random variable equal to the number of living pairs of children of $(u, f(u))$ stochastically dominates the random variable $X \sim \text{Bin} \left(\frac{\log n}{101}, p_0 \right)$, where $p_0 = \left(1 - \frac{2k}{q} \right)^2 = \left(\frac{1+3\epsilon}{1+5\epsilon} \right)^2$. The colorings of the descendants of each live pair are independent and so we have using the Chernoff bounds for $2 \leq j \leq k$,

$$\begin{aligned} \Pr \left[A_j < \left(\frac{\log n}{200} \right)^j p_0^{j-1} \mid A_{j-1} \geq \left(\frac{\log n}{200} \right)^{j-1} p_0^{j-2} \right] \\ \leq \exp \left\{ -\frac{1}{2} \cdot \left(\frac{99}{200} \right)^2 \cdot \frac{\log n}{101} \cdot \left(\frac{\log n}{200} \right)^{j-1} p_0^j \right\} = O(n^{-\Omega(\log \log n)}). \end{aligned} \quad (3.7)$$

(3.6) and (3.7) justify assuming that $A_k \geq \left(\frac{\log n}{200} \right)^k p_0^{k-1} \geq n^{\frac{4}{5}\epsilon}$.

• STEP 2: Existence of rainbow paths between x, y in $G_{x,y}$

Assuming that there are $\geq n^{4\epsilon/5}$ living pairs of leaves (x_i, y_i) for vertices x, y ,

$$\Pr(x, y \text{ are not rainbow connected}) \leq \left(1 - \prod_{i=0}^{2\gamma-1} \left(1 - \frac{2k+i}{q} \right) \right)^{n^{4\epsilon/5}}.$$

But

$$\prod_{i=0}^{2\gamma-1} \left(1 - \frac{2k+i}{q} \right) \geq \left(1 - \frac{2k+2\gamma}{q} \right)^{2\gamma} = \left(\frac{\epsilon}{1+5\epsilon} \right)^{2\gamma}.$$

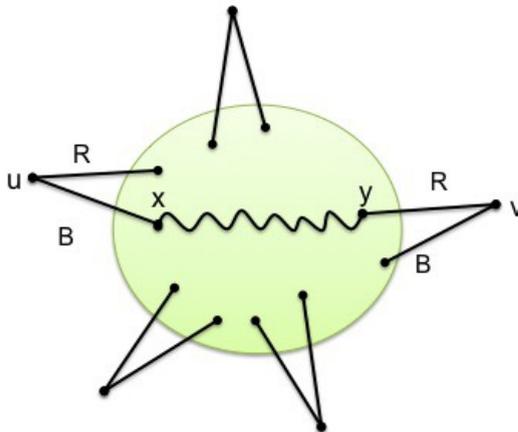


FIGURE 3.4: Taking care of small vertices.

So

$$\begin{aligned} \Pr(x, y \text{ are not rainbow connected}) &\leq \exp \left\{ -n^{4\epsilon/5} \left(\frac{\epsilon}{1+5\epsilon} \right)^{2\gamma} \right\} \\ &= \exp \left\{ -n^{4\epsilon/5 - O(\log(1/\epsilon)/\log \log n)} \right\}. \end{aligned} \quad (3.8)$$

Using (3.2) and the union bound taking (3.8) over all large x, y completes the proof of Lemma 3.7. \square

We now finish the proof of Theorem 3.1 i.e. take care of small vertices.

We showed in Lemma 3.7 that *whp* for any two large vertices, a random coloring results in a rainbow path joining them. We divide the small vertices into two sets: vertices of degree 1, V_1 and the vertices of degree at least 2, V_2 . Suppose that our colors are $1, 2, \dots, q$ and $V_1 = \{v_1, v_2, \dots, v_s\}$. We begin by giving the edge incident with v_i the color i . Then we slightly modify the argument in Lemma 3.7. If x is the neighbor of $v_i \in V_1$ then color i cannot be used in Steps 1 and 2 of that procedure. In terms of analysis this replaces q by $(q-1)$ ($(q-2)$ if y is also a neighbor of V_1) and the argument is essentially unchanged i.e. *whp* there will be a rainbow path between each pair of large vertices. Furthermore, any path starting at v_i can only use color i once and so there will be rainbow paths between V_1 and V_1 and between V_1 and the set of large vertices.

The set V_2 is treated by using only two extra colors. Assume that Red and Blue have not been used in our coloring. Then we use Red and Blue to color two of the edges incident to a vertex $u \in V_2$ (the remaining edges are colored arbitrarily). This is shown in Figure 3.4. Suppose that $V_2 = \{w_1, w_2, \dots, w_t\}$. Then if we want a rainbow path joining w_i, w_j where $i < j$ then we use the red edge to go to its neighbor w'_i . Then we

take the already constructed rainbow path to w_j'' , the neighbor of w_j via a blue edge. Then we can continue to w_j . \square

3.4 Proof of Theorem 3.2

We first observe that simply randomly coloring the edges of $G = G(n, r)$ with $q = n^{o(1)}$ colors will not do. This is because there will *whp* be $\Omega(nq^{1-r^2}) = \Omega(n^{1-o(1)})$ vertices v where all edges at distance at most two from v have the same color.

We follow a similar strategy to the proof in Theorem 3.1. We grow small trees T_x from each vertex x . Then for a pair of vertices x, y we build disjoint trees on the leaves of T_x, T_y so that *whp* we can find edge disjoint paths between any set of leaves S_x of T_x and any set of leaves of S_y of the same size. A bounded number of leaves of T_x, T_y will be excluded from this statement. The main difference will come from our procedure for coloring the edges. Because of the similarities, we will give a little less detail in the common parts of our proofs. We are in effect talking about building a structure like that shown in Figure 3.2. There is one difference, we will have to take care of which leaves of T_x we pair with which leaves of T_y , for a pair of vertices x, y .

Having grown the trees, we have the problem of coloring the edges. Instead of independently and randomly coloring the edges, we use a greedy algorithm that produces a coloring that is guaranteed to color edges differently, if they are close. This will guarantee that the edges of T_x are rainbow, for all vertices x . We then argue that we can find, for each vertex pair x, y , a partial mapping g from the leaves of T_x to the leaves of T_y such that the path from x to leaf v in T_x and the path from y to leaf $g(v)$ in T_y do not share a color. This assumes that v has an image under the partial mapping g . We will have to argue that g is defined on enough vertices in T_x . Given this, we then consider the colors on a set of edge disjoint paths that we can construct from the leaves of T_x to their g -counterpart in the leaves of T_y .

We use the the configuration model of Bollobás [77] in our proofs, see Chapter 1.1.2

3.4.1 Tree building

We will grow a Breadth First Search tree T_x from each vertex. We will grow each tree to depth

$$k = k_r = \begin{cases} \lceil \log_{r-2} \log n \rceil & r \geq 4. \\ \lceil 2 \log_2 \log n - 2 \log_2 \log_2 \log n \rceil & r = 3. \end{cases}$$

Observe that

$$T_x \text{ has at most } r(1 + (r-1) + (r-1)^2 + \dots + (r-1)^{k-1}) = r \frac{(r-1)^k - 1}{r-1} \text{ edges.} \quad (3.9)$$

It is useful to observe that

Lemma 3.8. *Whp, no set of $s \leq \ell_1 = \frac{1}{10} \log_{r-1} n$ vertices contains more than s edges.*

Proof Indeed,

$$\begin{aligned} \Pr(\exists S \subseteq [n], |S| \leq \ell_1, e[S] \geq |S| + 1) &\leq \sum_{s=3}^{\ell_1} \binom{n}{s} \binom{\binom{s}{2}}{s+1} \left(\frac{r^2}{rn-rs} \right)^{s+1} & (3.10) \\ &\leq \frac{r\ell_1}{n} \sum_{s=3}^{\ell_1} \binom{n}{s} \binom{\binom{s}{2}}{s} \left(\frac{r^2}{rn-rs} \right)^s \\ &\leq \frac{r\ell_1}{n} \sum_{s=3}^{\ell_1} \left(\frac{ne}{s} \cdot \frac{se}{2} \cdot \frac{2r}{n} \right)^s \\ &\leq \frac{r\ell_1}{n} \cdot \ell_1 \cdot (e^2 r)^{\ell_1} = o(1). & (3.11) \end{aligned}$$

Explanation of (3.10): The factor $\left(\frac{r^2}{rn-rs} \right)^{s+1}$ can be justified as follows. We can estimate

$$\Pr(e_1, e_2, \dots, e_{s+1} \in E(G_F)) = \prod_{i=0}^s \Pr(e_{i+1} \in E(G_F) \mid e_1, e_2, \dots, e_i \in E(G_F)) \leq \left(\frac{r^2}{rn-rs} \right)^{s+1}$$

if we pair up the lowest index endpoint of each e_i in some arbitrary order. The fraction $\frac{r^2}{rn-rs}$ is an upper bound on the probability that this endpoint is paired with the other endpoint, regardless of previous pairings. \square

Denote the leaves of T_x by L_x .

Corollary 3.9. *Whp, $(r-1)^k \leq |L_x| \leq r(r-1)^{k-1}$ for all $x \in [n]$.*

Proof This follows from the fact that *whp* the vertices spanned by each T_x span at most one cycle. This in turn follows from Lemma 3.8. \square

Consider two vertices $x, y \in V(G)$ where $T_x \cap T_y = \emptyset$. We will show that *whp* we can find a subgraph $G'(V', E'), V' \subseteq V, E' \subseteq E$ with similar structure to that shown in Figure 3.2. Here $k = k_r$ and $\gamma = (\frac{1}{2} + \epsilon) \log_{r-1} n$ for some small positive constant ϵ .

Remark 3.10. In our analysis we expose the pairing F , only as necessary. For example the construction of T_x involves exposing all pairings involving non-leaves of T_x and one

pairing for each leaf. There can be at most one exception to this statement, for the rare case where T_x contains a unique cycle. In particular, if we expose the point q paired with a currently unpaired point p of a leaf of T_x then q is chosen randomly from the remaining unpaired points.

Suppose that we have constructed $i = O(\log n)$ vertex disjoint trees of depth γ rooted at some of the leaves of T_x . We grow the $(i+1)$ st tree \widehat{T}_z via BFS, without using edges that go into y or previously constructed trees. Let a leaf $z \in L_x$ be *bad* if we have to omit a single edge as we construct the first $\ell_1/2$ levels of \widehat{T}_z . The previously constructed trees plus y account for $O(n^{1/2+\epsilon})$ vertices and pairings, so the probability that z is bad, given all the pairings we have exposed so far, is at most $O((r-1)^{\ell_1/2} n^{-1/2+\epsilon}) = O(n^{-1/3})$. Here bad edges can only join two leaves. This probability bound holds regardless of whichever other vertices are bad. This follows from the way we build the pairing F , see the final statement of Remark 3.10. So *whp* there will be at most 3 bad leaves on any T_x . Indeed, $\Pr(\exists x : x \text{ has } \geq 4 \text{ bad leaves}) \leq n^{O(\log n)} n^{-4/3} = o(1)$.

If a leaf is not bad then the first $\ell_1/2$ levels produce $\Theta(n^{1/20})$ leaves. From this, we see that *whp* the next $\gamma - \ell_1$ levels grow at a rate $r - 1 - o(n^{-1/25})$. Indeed, given that a level has L vertices where $n^{1/20} \leq L \leq n^{3/4}$, the number of vertices in the next level dominates $\text{Bin}\left((r-1)L, 1 - O\left(\frac{n^{3/4}}{n}\right)\right)$, after accounting for the configuration points used in building previous trees. Indeed, $(r-1)L$ configuration points associated with good leaves will be unpaired and for each of them, the probability it is paired with a point associated with a vertex in any of the trees constructed so far is $O(n^{1/2+2\epsilon}/n)$. This probability bound holds regardless of the pairings of the other leaf configuration points. We can thus assert that *whp* we will have that all but at most three of the leaves L_x of T_x are roots of vertex disjoint trees $\widehat{T}_1, \widehat{T}_2, \dots$, each with $\Theta(n^{1/2+\epsilon/2})$ leaves. Let L_x^* denote these *good* leaves. The same analysis applies when we build trees $\widehat{T}'_1, \widehat{T}'_2, \dots$, with roots at L_y .

Now the probability that there is no edge joining the leaves of \widehat{T}_i to the leaves of \widehat{T}'_j is at most

$$\left(1 - \frac{(r-1)\Theta(n^{1/2+\epsilon/2})}{rn}\right)^{(r-1)n^{1/2+\epsilon/2}} \leq e^{-\Omega(n^\epsilon)}.$$

To summarise,

Remark 3.11. *Whp* we will succeed in finding in G_F and hence in $G = G(n, r)$, for all $x, y \in V(G_F)$, for all $u \in L_x^*, v \in L_y^*$, a path $P_{u,v}$ from u to v of length $O(\log n)$ such that if $u \neq u'$ and $v \neq v'$ then $P_{u,v}$ and $P_{u',v'}$ are edge disjoint. These paths avoid T_x, T_y except at their start and endpoints.

3.4.2 Coloring the edges

We now consider the problem of coloring the edges of G . Let H denote the line graph of G and let $\Gamma = H^{2k}$ denote the graph with the same vertex set as H and an edge between vertices e, f of Γ if there is a path of length at most k between e and f in H . We will construct a proper coloring of Γ using

$$q = 10(r-1)^{2k} \sim 100 \log^{2\theta_r} n \text{ where } \theta_r = \frac{\log(r-1)}{\log(r-2)}$$

colors. We do this as follows: Let e_1, e_2, \dots, e_m be an arbitrary ordering of the vertices of Γ . For $i = 1, 2, \dots, m$, color e_i with a random color, chosen uniformly from the set of colors not currently appearing on any neighbor in Γ . At this point only e_1, e_2, \dots, e_{i-1} will have been colored.

Suppose then that we color the edges of G using the above method. Fix a pair of vertices x, y of G . We see immediately, that no color appears twice in T_x and no color appears twice in T_y . This is because the distance between edges in T_x is at most $2k$. This also deals with the case where $V(T_x) \cap V(T_y) \neq \emptyset$, for the same reason. So assume now that T_x, T_y are vertex disjoint. We can find lots of paths joining x and y . We know that the first and last k edges of each path will be individually rainbow colored. We will first show that we have many choices of path where these $2k$ edges are rainbow colored when taken together.

3.4.3 Case 1: $r \geq 4$:

We argue now that we can find $\sigma_0 = (r-2)^{k-1}$ leaves $u_1, u_2, \dots, u_\tau \in T_x$ and σ_0 leaves $v_1, v_2, \dots, v_\tau \in T_y$ such for each i the T_x path from x to u_i and the T_y path from y to v_i do not share any colors.

Lemma 3.12. *Let T_1, T_2 be two vertex disjoint copies of an edge colored complete d -ary tree with ℓ levels, where $d \geq 3$. Let T_1, T_2 be rooted at x, y respectively. Suppose that the colorings of T_1, T_2 are both rainbow. Let $\kappa = (d-1)^\ell$. Then there exist leaves $u_1, u_2, \dots, u_\kappa$ of T_1 and leaves $v_1, v_2, \dots, v_\kappa$ of T_2 such that the following is true: If P_i, P'_i are the paths from x to u_i in T_1 and from y to v_i in T_2 respectively, then $P_i \cup P'_i$ is rainbow colored for $i = 1, 2, \dots, \kappa$.*

Proof Let A_ℓ be the minimum number of rainbow path pairs that we can find in any such pair of edge colored trees. We prove that $A_\ell \geq (d-1)^\ell$ by induction on ℓ . This is true trivially for $\ell = 0$. Suppose that x is incident with x_1, x_2, \dots, x_d and that the sub-tree rooted at x_i is $T_{1,i}$ for $i = 1, 2, \dots, d$. Define y_i and $T_{2,i}$, $i = 1, 2, \dots, d$ similarly

with respect to y . Suppose that the color of the edge (x, x_i) is c_i for $i = 1, 2, \dots, d$ and let $Q_x = \{c_1, c_2, \dots, c_d\}$. Similarly, suppose that the color of the edge (y, y_i) is c'_i for $i = 1, 2, \dots, d$ and let $Q_y = \{c'_1, c'_2, \dots, c'_d\}$. Next suppose that Q_j is the set of colors in Q_x that appear on the edges $E(T_{2,j}) \cup \{(y, y_j)\}$. The sets Q_1, Q_2, \dots, Q_d are pair-wise disjoint. Similarly, suppose that Q'_i is the set of colors in Q_y that appear on the edges $E(T_{1,i}) \cup \{(x, x_i)\}$. The sets Q'_1, Q'_2, \dots, Q'_d are pair-wise disjoint.

Now define a bipartite graph H with vertex set $A + B = [d] + [d]$ and an edge (i, j) iff $c_i \notin Q_j$ and $c'_j \notin Q'_i$. We claim that if $S \subseteq A$ then its neighbor set $N_H(S)$ satisfies the inequality

$$d|S| - |N_H(S)| - |S| \leq |S| \cdot |N_H(S)|. \quad (3.12)$$

Here the LHS of (3.12) bounds from below, the size of the set $S : N_H(S)$ of edges between S and $N_H(S)$. This is because there are at most $|S|$ edges missing from $S : N_H(S)$ due to $i \in S$ and $j \in N_H(S)$ and $c_i \in Q_j$. At most $|N_H(S)|$ edges are missing for similar reasons. On the other hand, $d|S|$ is the number there would be without these missing edges. The RHS of (3.12) is a trivial upper bound.

Re-arranging we get that

$$|N_H(S)| - |S| \geq \left\lceil \frac{(d-2-|S|)|S|}{|S|+1} \right\rceil \geq -1.$$

(We get -1 when $|S| = d$).

Thus H contains a matching M of size $d-1$. Suppose without loss of generality that this matching is $(i, i), i = 1, 2, \dots, d-1$. We know by induction that for each i we can find paths $(P_{i,j}, \widehat{P}_{i,j}), j = 1, 2, \dots, (d-1)^{\ell-1}$ where $P_{i,j}$ is a root to leaf path in $T_{1,i}$ and $\widehat{P}_{i,j}$ is a root to leaf path in $T_{2,i}$ and that $P_{i,j} \cup \widehat{P}_{i,j}$ is rainbow for all i, j . Furthermore, (i, i) being an edge of H , means that the edge sets $\{(x, x_i)\} \cup E(P_{i,j}) \cup E(\widehat{P}_{i,j}) \cup \{(y, y_i)\}$ are all rainbow. \square

Let

$$V_1 = \{x : V(T_x) \text{ contains a cycle}\}.$$

When $x, y \notin V_1$ we apply this Lemma to T_x, T_y by deleting one of the r sub-trees attached to each of x, y and applying the lemma directly to the $(r-1)$ -ary trees that remain. This will yield $(r-2)^k$ pairs of paths. If $x \in V_1$, we delete $r-2$ sub-trees attached to x leaving at least two $(r-1)$ -ary trees of depth $k-1$ with roots adjacent to x . We can do the same at y . Let c_1, c_2 be the colors of the two edges from x to the roots of these two trees T_1, T_2 . Similarly, let c'_1, c'_2 be the colors of the two analogous edges from y to the trees T'_1, T'_2 . If color c_1 does not appear in T'_1 then we apply the lemma to T_1 and T'_1 .

Otherwise, we can apply the lemma to T_1 and T_2' . In both cases we obtain $(r-2)^{k-1}$ pairs of paths.

Accounting for bad vertices we put

$$\sigma = \sigma_0 - 6 = (r-2)^{k-1} - 6 \geq \frac{\log n}{r-2} - 6$$

and we see from Remark 3.11 that we can *whp* find σ paths $P_1, P_2, \dots, P_\sigma$ of length $O(\log n)$ from x to y . Path P_i goes from x to a leaf $u_i \in L_x^*$ via T_x and then traverses $Q_i = P(u_i, v_i)$ where $v_i = \phi(u_i) \in L_y^*$ and then goes from v_i to a y via T_y . Here ϕ is some partial map from L_x^* to L_y^* . It is a random variable that depends on the coloring \mathcal{C} of the edges of T_x and T_y . The paths $P_1, P_2, \dots, P_\sigma$ depend on the choice of ϕ and hence \mathcal{C} and so we should write $P_i = P_i(\mathcal{C})$.

We fix the coloring \mathcal{C} and hence $P_1, P_2, \dots, P_\sigma$. Let \mathcal{R} be the event that at least one of the paths $P_1, P_2, \dots, P_\sigma$ is rainbow colored. We show that $\Pr(\neg\mathcal{R} \mid \mathcal{C})$ is small.

We let $c(e)$ denote the color of edge e in a given coloring. We remark next that for a particular coloring c_1, c_2, \dots, c_m of the edges e_1, e_2, \dots, e_m we have

$$\Pr(c(e_i) = c_i, i = 1, 2, \dots, m) = \prod_{i=1}^m \frac{1}{a_i}$$

where $q - \Delta \leq a_i \leq q$ is the number of colors available for the color of the edge e_i given the coloring so far i.e. the number of colors unused by the neighbors of e_i in Γ when it is about to be colored.

Now fix an edge $e = e_i$ and the colors $c_j, j \neq i$. Let C be the set of colors not used by the neighbors of e_i in Γ . The choice by e_i of its color under this conditioning is not quite random, but close. Indeed, we claim that for $c, c' \in C$

$$\frac{\Pr(c(e) = c \mid c(e_j) = c_j, j \neq i)}{\Pr(c(e) = c' \mid c(e_j) = c_j, j \neq i)} \leq \left(\frac{q - \Delta}{q - \Delta - 1} \right)^\Delta.$$

This is because, changing the color of e_i only affects the number of colors available to neighbors of e_i , and only by at most one.

Thus, for $c \in C$, we have

$$\Pr(c(e) = c \mid c(e_j) = c_j, j \neq i) \leq \frac{1}{q - \Delta} \left(\frac{q - \Delta}{q - \Delta - 1} \right)^\Delta.$$

Now $\Delta \leq (r-1)^{2k} = q/10$ and we deduce that

$$\Pr(c(e) = c \mid c(e_j) = c_j, j \neq i) \leq \frac{2}{q}.$$

It follows that for $i \in [\sigma]$,

$$\Pr(P_i \text{ is rainbow colored} \mid \mathcal{C}, \text{ coloring of } \bigcup_{j \neq i} Q_j) \geq \left(1 - \frac{4(k + \gamma)}{q}\right)^{2\gamma}.$$

This is because when we consider the coloring of Q_i there will always be at most $2k + 2\gamma$ colors forbidden by non-neighboring edges, if it is to be rainbow colored.

It then follows that

$$\begin{aligned} \Pr(\neg \mathcal{R} \mid \mathcal{C}) &\leq \left(1 - \left(1 - \frac{4(k + \gamma)}{q}\right)^{2\gamma}\right)^\sigma \\ &\leq \left(\frac{8\gamma(k + \gamma)}{q}\right)^\sigma \\ &\leq \left(\frac{(2 + 10\epsilon) \log_{r-1}^2 n}{10 \log^{\theta_r} n}\right)^\sigma = o(n^{-2}). \end{aligned}$$

This completes the proof of Theorem 3.2 when $r \geq 4$.

Case 2: $r = 3$:

When $r = 3$ we can't use $(r - 2)^k$ to any effect. Also, we need to increase q to $\log^4 n$. This necessary for a variety of reasons. One reason is that we will reduce σ to $2^{k/2}$. We want this to be $\Omega(\log n)$ and this will force k to (roughly) double what it would have been if we had followed the recipe for $r \geq 4$. This makes Δ close to $\log^4 n$ and we need $q \gg \Delta$.

And we need to modify the argument based on Lemma 3.12. Instead of inducting on the trees at depth one from the roots x, y , we now induct on the trees at depth two. Assume first that $x, y \notin V_1$. After ignoring one branch for T_x and T_y we now consider the sub-trees $T_{x,i}, T_{y,i}$, $i = 1, 2, 3, 4$ of T_x, T_y whose roots x_1, \dots, x_4 and y_1, \dots, y_4 are at depth two. We cannot necessarily make this construction when $x \in V_1$. Let P_i be the path from x to x_i in T_x and let \hat{P}_j be the path from y to y_j in T_y . Next suppose that \hat{Q}_j is the set of colors in Q that appear on the edges $E(T_{y,j}) \cup E(\hat{P}_j)$. Similarly, suppose that Q'_i is the set of colors in Q' that appear on the edges $\{E(T_{x,i}) \cup E(P_i)\}$.

Re-define H to be the bipartite graph with vertex set $A + B = [4] + [4]$. The edges of H are as before: (i, j) exists iff $c_i \notin Q_j$ and $c'_j \notin \hat{Q}_i$. This time we can only say that a color is in at most two \hat{Q}_i 's and similarly for the Q'_j 's. The effect of this is to replace (3.12) by

$$4|S| - 2(|N_H(S)| + |S|) \leq |S| \cdot |N_H(S)|$$

from which we can deduce that

$$|S| - |N_H(S)| \leq \frac{|S| \cdot |N_H(S)|}{2} \leq 2|N_H(S)|.$$

It follows that $|N_H(S)| \geq \lceil |S|/3 \rceil \geq |S| - 2$ and so H contains a matching of size two. An inductive argument then shows that we are able to find $2^{\lfloor k/2 \rfloor}$ rainbow pairs of paths. The proof now continues as in the case $r \geq 4$, arguing about the coloring of paths $P_1, P_2, \dots, P_\sigma$ where now $\sigma = 2^{\lfloor k/2 \rfloor}$.

We finally deal with the vertices in V_1 . We classify them according to the size of the cycle C_x that is contained in $V(T_x)$. If T_x contains a cycle C_x then necessarily $|C_x| \leq 2k$ and so there are at most $2k$ types in our classification. It follows from Lemma 3.8 that if $x, y \in V_1$ and $T_x \cap T_y \neq \emptyset$ then $C_x = C_y$ *whp*. Note next that the distance from x to C_x is at most $k - |C_x|/2$. If C is a cycle of length at most $2k$, let $V_C = \{x : C = C_x\}$ and let E_C be the set of edges contained in V_C . We have

$$|V_C| = O(|C|2^{k-|C|/2}) = O(2^k) = O(\log^2 n / \log \log n). \quad (3.13)$$

We introduce $2k$ new sets $\widehat{Q}_i, i = 3, 4, \dots, 2k$ of $O(\log^2 n / \log \log n)$ colors, distinct from Q . Thus we introduce $O(\log^2 n)$ new colors overall. We re-color each E_C with the colors from $\widehat{Q}_{|C|}$. It is important to observe that if $|C| = |C'|$ then the graphs induced by V_C and $V_{C'}$ are isomorphic and so we can color them isomorphically. By the latter we mean that we choose some isomorphism f from V_C to $V_{C'}$ and then if e is an edge of V_C then we color e and $f(e)$ with the same color. After this re-coloring, we see that if T_x and T_y are not vertex disjoint, then they are contained in the same V_C . The edges of V_C are rainbow colored and so now we only need to concern ourselves with $x, y \in V_1$ such that T_x and T_y are vertex disjoint. Assume now that $x, y \in V_1$.

Assume first that x, y are of the same type and that they are at the same distance from C_x, C_y respectively. Our aim now is to define binary trees T'_x, T'_y “contained“ in T_x, T_y that can be used as in Lemma 3.12. If we delete an edge $e = (u, v)$ of C_x then the graph that remains on $V(T_x)$ is a tree with at most two vertices u, v of degree two. Now delete one of the three sub-trees of T_x . If there are vertices of degree two, make sure one of them is in this sub-tree. If necessary, shrink the path of length two with the remaining vertex of degree two in the middle to an edge e_x . It has leaves at depth $k - 1$ and leaves at depth $k - 2$. The resulting binary tree will be our T'_x . The leaves at depth $k - 1$ come in pairs. Delete one vertex from each pair and shrink the paths of length two through the vertex at depth $k - 2$ to an edge.

The edges that are obtained by shrinking paths of length two will have two colors. Because x, y are at the same distance from their cycles, we can delete $f(e)$ from C_y and

do the construction so that T'_x and T'_y will be isomorphically colored.

It is now easy to find 2^{k-2} pairs of paths whose unions are rainbow colored. Each leaf of T_x, T_y can be labelled by a $\{0, 1\}$ string of length $k - 2$. We pair string $\xi_1\xi_2 \cdots \xi_{k-1}\xi_{k-2}$ in T_x with $(1 - \xi_1)\xi_2 \cdots \xi_{k-1}\xi_{k-2}$ in T_y . The associated paths will have a rainbow union. The proof now continues as in the case $r \geq 4$, arguing about the coloring of paths $P_1, P_2, \dots, P_\sigma$ where now $\sigma = 2^{k-2}$.

If x is further from C_x than y is from C_y then let z be the vertex on the path from x to C_x at the same distance from C_x as y is from C_y . We have a rainbow path from z to y and adding the T_x path from x to z gives us a rainbow path from x to y . This relies on the fact that V_{C_x} and V_{C_y} are isomorphically colored.

If x, y are of a different type, then T_x and T_y are re-colored with distinct colors and we can proceed as as in the case $r \geq 4$, arguing about the coloring of paths $P_1, P_2, \dots, P_\sigma$ where now $\sigma = 2^k$, using Corollary 3.9.

If $x \in V_1$ and $y \notin V_1$ then we can proceed as if both are not in V_1 . This is because of the re-coloring of the edges of T_x . We can proceed as as in the case $r \geq 4$, arguing about the coloring of paths $P_1, P_2, \dots, P_\sigma$ where now $\sigma = 2^k$, using Corollary 3.9.

This completes our proof of Theorem 3.2.

We conclude this Chapter with mentioning that if the degree r in Theorem 3.2 is allowed to grow as fast as $\log n$ then one can prove a result closer to that of Theorem 3.1.

Chapter 4

Random Apollonian networks

4.1 Model & Main Results

As we outlined in Chapter 1 planar graphs model several significant types of spatial real-world networks such as power grids and road networks. Despite the outstanding amount of work on modeling real-world networks with random graph models [19, 50, 84, 85, 143, 157, 172, 173, 270, 284, 294], real-world planar graph generators have received considerably less attention. In this Chapter we focus on Random Apollonian Networks (RANs), a popular random graph model for generating planar graphs with power law properties [426]. Before we state our main results we briefly describe the model.

Model: An example of a RAN is shown in Figure 4.1. At time $t = 1$ the RAN is shown in Figure 4.1(a). At each step $t \geq 2$ a face F is chosen uniformly at random among the faces of G_t . Let i, j, k be the vertices of F . We add a new vertex inside F and we connect it to i, j, k . Higher dimensional RANs also exist where instead of triangles we have k -simplexes $k \geq 3$, see [425]. It is easy to see that the number of vertices n_t , edges m_t and faces F_t at time $t \geq 1$ in a RAN G_t satisfy:

$$n_t = t + 3, \quad m_t = 3t + 3, \quad F_t = 2t + 1.$$

Note that a RAN is a maximal planar graph since for any planar graph $m_t \leq 3n_t - 6 \leq 3t + 3$.

Surprisingly, despite the popularity of the model various important properties have been analyzed experimentally and heuristically with lack of rigor. In this Chapter, we prove the following theorems.

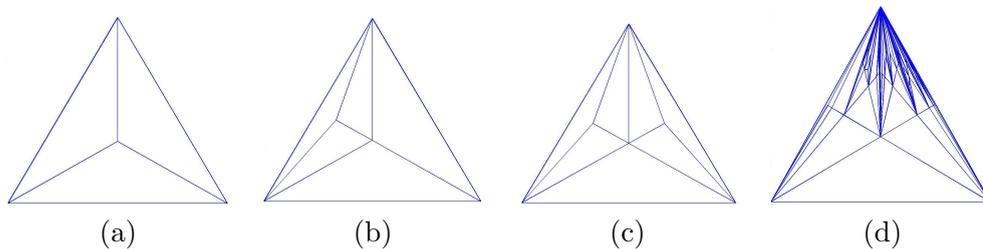


FIGURE 4.1: Snapshots of a Random Apollonian Network (RAN) at: (a) $t = 1$ (b) $t = 2$ (c) $t = 3$ (d) $t = 100$.

Theorem 4.1 (Degree Sequence). *Let $Z_k(t)$ denote the number of vertices of degree k at time t , $k \geq 3$. For any $t \geq 1$ and any $k \geq 3$ there exists a constant b_k depending on k such that*

$$|\mathbb{E}[Z_k(t)] - b_k t| \leq K, \quad \text{where } K = 3.6.$$

Furthermore, for t sufficiently large and any $\lambda > 0$

$$\Pr[|Z_k(t) - \mathbb{E}[Z_k(t)]| \geq \lambda] \leq e^{-\frac{\lambda^2}{72t}}. \quad (4.1)$$

For previous weaker results on the degree sequence see [420, 426]. An immediate corollary which proves strong concentration of $Z_k(t)$ around its expectation is obtained from Theorem 4.1 and a union bound by setting $\lambda = 10\sqrt{t \log t}$. Specifically:

Corollary 4.2. *For all possible degrees k*

$$\Pr\left[|Z_k(t) - \mathbb{E}[Z_k(t)]| \geq 10\sqrt{t \log t}\right] = o(1).$$

The next theorem provides insight into the asymptotic growth of the highest degrees of RANs and is crucial in proving Theorem 4.4.

Theorem 4.3 (Highest Degrees). *Let $\Delta_1 \geq \Delta_2 \geq \dots \geq \Delta_k$ be the k highest degrees of the RAN G_t at time t where k is a fixed positive integer. Also, let $f(t)$ be a function such that $f(t) \rightarrow +\infty$ as $t \rightarrow +\infty$. Then whp¹*

$$\frac{t^{1/2}}{f(t)} \leq \Delta_1 \leq t^{1/2} f(t)$$

and for $i = 2, \dots, k$

¹An event A_t holds with high probability (whp) if $\lim_{t \rightarrow +\infty} \Pr[A_t] = 1$.

$$\Delta_{i-1} - \Delta_i \geq \frac{t^{1/2}}{f(t)}.$$

The growing function $f(t)$ cannot be removed, see [171]. Using Theorem 4.3 and the technique of Mihail and Papadimitriou [302] we show how the top eigenvalues of the adjacency matrix representation of a RAN grow asymptotically as $t \rightarrow +\infty$ whp.

Theorem 4.4 (Largest Eigenvalues). *Let k be a fixed positive integer. Also, let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ be the largest k eigenvalues of the adjacency matrix of G_t . Then whp $\lambda_i = (1 \pm o(1))\sqrt{\Delta_i}$.*

Also, we show the following refined upper bound for the asymptotic growth of the diameter.

Theorem 4.5 (Diameter). *The diameter $d(G_t)$ of G_t satisfies in probability $d(G_t) \leq \rho \log t$ where $\frac{1}{\rho} = \eta$ is the unique solution less than 1 of the equation $\eta - 1 - \log \eta = \log 3$.*

A straight-forward calculation of η , gives us the following corollary.

Corollary 4.6 (Diameter). *The diameter $d(G_t)$ of G_t satisfies asymptotically*

$$\Pr[d(G_t) > 7.1 \log t] \rightarrow 0.$$

The outline of this Chapter is as follows: in Section 4.2 we present briefly related work and technical preliminaries needed for our analysis. We prove Theorems 4.1, 4.3, 4.4 and 4.5 in Sections 4.3, 4.4, 4.5 and 4.6 respectively.

4.2 Related Work

Apollonius of Perga was a Greek geometer and astronomer noted for his writings on conic sections. He introduced the problem of space filling packing of spheres whose classical solution, the so-called Apollonian packing [202], exhibits a power law behavior. Specifically, the circle size distribution follows a power law with exponent around 1.3 [86]. Apollonian Networks (ANs) were introduced in [35] and independently in [139]. Zhou et al. [426] introduced Random Apollonian Networks (RANs). Their degree sequence was analyzed inaccurately in [426] (see comment in [420]) and subsequently using physicist's

methodology in [420]. Eigenvalues of RANs have been studied only experimentally [34]. Concerning the diameter of RNAs it has been shown to grow logarithmically [426] using heuristic arguments (see for instance equation B6, Appendix B in [426]). RANs are planar 3-trees, a special case of random k -trees [254]. Cooper and Uehara [124] and Gao [188] analyzed the degree distribution of random k -trees, a closely related model to RANs. In RANs –in contrast to random k -trees– the random k clique chosen at each step has never previously been selected. For example, in the two dimensional RAN any chosen face is being subdivided into three new faces by connecting the incoming vertex to the vertices of the boundary. Random k -trees due to their power law properties have been proposed as a model for complex networks, see, e.g., [124, 187] and references therein. Recently, a variant of k -trees, namely ordered increasing k -trees has been proposed and analyzed in [329]. Closely related to RANs but not the same are random Apollonian network structures which have been analyzed by Darrasse, Soria et al. [72, 128, 129].

Bollobás, Riordan, Spencer and Tusnády [80] proved rigorously the power law distribution of the Barabási-Albert model [50]. Chung, Lu, Vu [113] Flaxman, Frieze, Fenner [171] and Mihail, Papadimitriou [302] have proved rigorous results for eigenvalue related properties of real-world graphs using various random graph models.

4.3 Proof of Theorem 4.1

We decompose our proof in a sequence of Lemmas. For brevity let $N_k(t) = \mathbb{E}[Z_k(t)]$, $k \geq 3$. Also, let $d_v(t)$ be the degree of vertex v at time t and $\mathbf{1}(d_v(t) = k)$ be an indicator variable which equals 1 if $d_v(t) = k$, otherwise 0. Then, for any $k \geq 3$ we can express the expected number $N_k(t)$ of vertices of degree k as a sum of expectations of indicator variables:

$$N_k(t) = \sum_v \mathbb{E}[\mathbf{1}(d_v(t) = k)]. \quad (4.2)$$

We distinguish two cases in the following.

- CASE 1: $k = 3$:

Observe that a vertex of degree 3 is created only by an insertion of a new vertex. The expectation $N_3(t)$ satisfies the following recurrence²

²The three initial vertices participate in one less face than their degree. However, this leaves our results unchanged.

$$N_3(t+1) = N_3(t) + 1 - \frac{3N_3(t)}{2t+1}. \quad (4.3)$$

The basis for Recurrence (4.3) is $N_3(1) = 4$. We prove the following lemma which shows that $\lim_{t \rightarrow +\infty} \frac{N_3(t)}{t} = \frac{2}{5}$.

Lemma 4.7. $N_3(t)$ satisfies the following inequality:

$$|N_3(t) - \frac{2}{5}t| \leq K, \quad \text{where } K = 3.6 \quad (4.4)$$

Proof. We use induction. Assume that $N_3(t) = \frac{2}{5}t + e_3(t)$, where $e_3(t)$ stands for the error term. We wish to prove that for all t , $|e_3(t)| \leq K$. The result trivially holds for $t = 1$. We also see that for $t = 1$ inequality (4.4) is tight. Assume the result holds for some t . We show it holds for $t + 1$.

$$\begin{aligned} N_3(t+1) &= N_3(t) + 1 - \frac{3N_3(t)}{2t+1} \Rightarrow \\ e_3(t+1) &= e_3(t) + \frac{3}{5} - \frac{6t + 15e_3(t)}{10t+5} = e_3(t) \left(1 - \frac{3}{2t+1}\right) + \frac{3}{5(2t+1)} \Rightarrow \\ |e_3(t+1)| &\leq K \left(1 - \frac{3}{2t+1}\right) + \frac{3}{5(2t+1)} \leq K \end{aligned}$$

Therefore inductively Inequality (4.4) holds for all $t \geq 1$. □

• CASE 2: $k \geq 4$:

For $k \geq 4$ the following holds:

$$\mathbb{E}[\mathbf{1}(d_v(t+1) = k)] = \mathbb{E}[\mathbf{1}(d_v(t) = k)] \left(1 - \frac{k}{2t+1}\right) + \mathbb{E}[\mathbf{1}(d_v(t) = k-1)] \frac{k-1}{2t+1} \quad (4.5)$$

Therefore, we can rewrite Equation (4.2) for $k \geq 4$ as follows:

$$N_k(t+1) = N_k(t) \left(1 - \frac{k}{2t+1}\right) + N_{k-1}(t) \frac{k-1}{2t+1} \quad (4.6)$$

Lemma 4.8. For any $k \geq 3$, the limit $\lim_{t \rightarrow +\infty} \frac{N_k(t)}{t}$ exists. Specifically, let $b_k = \lim_{t \rightarrow +\infty} \frac{N_k(t)}{t}$. Then, $b_3 = \frac{2}{5}, b_4 = \frac{1}{5}, b_5 = \frac{4}{35}$ and for $k \geq 6$ $b_k = \frac{24}{k(k+1)(k+2)}$. Furthermore, for all $k \geq 3$

$$|N_k(t) - b_k t| \leq K, \quad \text{where } K = 3.6. \quad (4.7)$$

Proof. For $k = 3$ the result holds by Lemma 4.7 and specifically $b_3 = \frac{2}{5}$. Assume the result holds for some k . We show that it holds for $k + 1$ too. Rewrite Recursion (4.6) as: $N_k(t + 1) = (1 - \frac{b_t}{t+t_1})N_k(t) + c_t$ where $b_t = k/2, t_1 = 1/2, c_t = N_{k-1}(t) \frac{k-1}{2t+1}$. Clearly $\lim_{t \rightarrow +\infty} b_t = k/2 > 0$ and $\lim_{t \rightarrow +\infty} c_t = \lim_{t \rightarrow +\infty} b_{k-1} t \frac{k-1}{2t+1} = b_{k-1}(k-1)/2$. Hence by Lemma 2.8:

$$\lim_{t \rightarrow +\infty} \frac{N_k(t)}{t} = \frac{(k-1)b_{k-1}/2}{1+k/2} = b_{k-1} \frac{k-1}{k+2}.$$

Since $b_3 = \frac{2}{5}$ we obtain that $b_4 = \frac{1}{5}, b_5 = \frac{4}{35}$ for any $k \geq 6, b_k = \frac{24}{k(k+1)(k+2)}$. This shows that the degree sequence of RANs follows a power law distribution with exponent 3.

Now we prove Inequality (4.7). The case $k = 3$ was proved in Lemma 4.7. Let $e_k(t) = N_k(t) - b_k t$. Assume the result holds for some $k \geq 3$, i.e., $|e_k(t)| \leq K$ where $K = 3.6$. We show it holds for $k + 1$ too. Substituting in Recurrence (4.2) and using the fact that $b_{k-1}(k-1) = b_k(k+2)$ we obtain the following:

$$\begin{aligned} e_k(t+1) &= e_k(t) + \frac{k-1}{2t+1} e_{k-1}(t) - \frac{k}{2t+1} e_k(t) \Rightarrow \\ |e_k(t+1)| &\leq |(1 - \frac{k}{2t+1})e_k(t)| + |\frac{k-1}{2t+1} e_{k-1}(t)| \leq K(1 - \frac{1}{2t+1}) \leq K \end{aligned}$$

Hence by induction, Inequality (4.7) holds for all $k \geq 3$. □

Using integration and a first moment argument, it can be seen that Lemma 4.8 agrees with Theorem 4.3 where it is shown that the maximum degree is $\approx t^{1/2}$. (While $b_k = O(k^{-3})$ suggests a maximum degree of order $t^{1/3}$, summing b_k over $k \geq K$ suggests a maximum degree of order $t^{1/2}$).

Finally, the next Lemma proves the concentration of $Z_k(t)$ around its expected value for $k \geq 3$. This lemma applies Lemma 2.3 and completes the proof of Theorem 4.1.

Lemma 4.9. *Let $\lambda > 0$. For $k \geq 3$*

$$\Pr [|Z_k(t) - \mathbb{E}[Z_k(t)]| \geq \lambda] \leq e^{-\frac{\lambda^2}{72t}}. \quad (4.8)$$

Proof. Let $(\Omega, \mathcal{F}, \mathbb{P})$ be the probability space induced by the construction of a RAN after t insertions. Fix k , where $k \geq 3$, and let $(X_i)_{i \in \{0, 1, \dots, t\}}$ be the martingale sequence defined by $X_i = \mathbb{E}[Z_k(t) | \mathcal{F}_i]$, where $\mathcal{F}_0 = \{\emptyset, \Omega\}$ and \mathcal{F}_i is the σ -algebra generated by the RAN process after i steps. Notice $X_0 = \mathbb{E}[Z_k(t) | \{\emptyset, \Omega\}] = N_k(t)$, $X_t = Z_k(t)$. We show that $|X_{i+1} - X_i| \leq 6$ for $i = 0, \dots, t-1$. Let $P_j = (Y_1, \dots, Y_{j-1}, Y_j)$, $P'_j = (Y_1, \dots, Y_{j-1}, Y'_j)$ be two sequences of face choices differing only at time j . Also, let \bar{P}, \bar{P}' continue from P_j, P'_j until t . We call the faces Y_j, Y'_j special with respect to \bar{P}, \bar{P}' . We define a measure preserving map $\bar{P} \mapsto \bar{P}'$ in the following way: for every choice of a non-special face in process \bar{P} at time l we make the same face choice in \bar{P}' at time l . For every choice of a face inside the special face Y_j in process \bar{P} we make an isomorphic (w.r.t., e.g., clockwise order and depth) choice of a face inside the special face Y'_j in process \bar{P}' . Since the number of vertices of degree k can change by at most 6, i.e., the (at most) 6 vertices involved in the two faces Y_j, Y'_j the following holds:

$$|\mathbb{E}[Z_k(t) | P] - \mathbb{E}[Z_k(t) | P']| \leq 6.$$

Furthermore, this holds for any P_j, P'_j . We deduce that X_{i-1} is a weighted mean of values, whose pairwise differences are all at most 6. Thus, the distance of the mean X_{i-1} is at most 6 from each of these values. Hence, for any one step refinement $|X_{i+1} - X_i| \leq 6 \forall i \in \{0, \dots, t-1\}$. By applying the Azuma-Hoeffding inequality as stated in Lemma 2.3 we obtain

$$\Pr [|Z_k(t) - \mathbb{E}[Z_k(t)]| \geq \lambda] \leq 2e^{-\frac{\lambda^2}{72t}}. \quad (4.9)$$

□

4.4 Proof of Theorem 4.3

We decompose the proof of Theorem 4.3 into several lemmas which we prove in the following. Specifically, the proof follows directly from Lemmas 4.11, 4.12, 4.13, 4.14, 4.15. We partition the vertices into three sets: those added before t_0 , between t_0 and t_1 and after t_1 where $t_0 = \log \log \log(f(t))$ and $t_1 = \log \log(f(t))$. Recall that $f(t)$ is a function

such that $\lim_{t \rightarrow +\infty} f(t) = +\infty$. We define a supernode to be a collection of vertices and the degree of the supernode the sum of the degrees of its vertices.

Lemma 4.10. *Let $d_t(s)$ denote the degree of vertex s at time t . and let $a^{(k)} = a(a+1)\dots(a+k-1)$ denote the rising factorial function. Then, for any positive integer k*

$$\mathbb{E} \left[d_t(s)^{(k)} \right] \leq \frac{(k+2)!}{2} \left(\frac{2t}{s} \right)^{\frac{k}{2}}. \quad (4.10)$$

Proof. As we mentioned in the proof of Theorem 4.1 the three initial vertices 1, 2, 3 have one less face than their degree whereas all other vertices have degree equal to the number of faces surrounding them. In this proof we treat both cases but we omit it in all other proofs.

• CASE 1: $s \geq 4$

Note that $d_s(s) = 3$. By conditioning successively we obtain

$$\begin{aligned} \mathbb{E} \left[d_t(s)^{(k)} \right] &= \mathbb{E} \left[\mathbb{E} \left[d_t(s)^{(k)} | d_{t-1}(s) \right] \right] \\ &= \mathbb{E} \left[(d_{t-1}(s))^{(k)} \left(1 - \frac{d_{t-1}(s)}{2t-1} \right) + (d_{t-1}(s) + 1)^{(k)} \frac{d_{t-1}(s)}{2t-1} \right] \\ &= \mathbb{E} \left[(d_{t-1}(s))^{(k)} \left(1 - \frac{d_{t-1}(s)}{2t-1} \right) + (d_{t-1}(s))^{(k)} \frac{d_{t-1}(s) + k}{d_{t-1}(s)} \frac{d_{t-1}(s)}{2t-1} \right] \\ &= \mathbb{E} \left[(d_{t-1}(s))^{(k)} \right] \left(1 + \frac{k}{2t-1} \right) = \dots = 3^{(k)} \prod_{t'=s+1}^t \left(1 + \frac{k}{2t'-1} \right) \\ &\leq 3^{(k)} \exp \left(\sum_{t'=s+1}^t \frac{k}{2t'-1} \right) \leq 3^{(k)} \exp \left(k \int_s^t \frac{dx}{2x-1} \right) \\ &\leq \frac{(k+2)!}{2} \exp \left(\frac{k}{2} \log \frac{t-1/2}{s-1/2} \right) \leq \frac{(k+2)!}{2} \left(\frac{2t}{s} \right)^{\frac{k}{2}}. \end{aligned}$$

• CASE 2: $s \in \{1, 2, 3\}$

Note that initially the degree of any such vertex is 2. For any $k \geq 0$

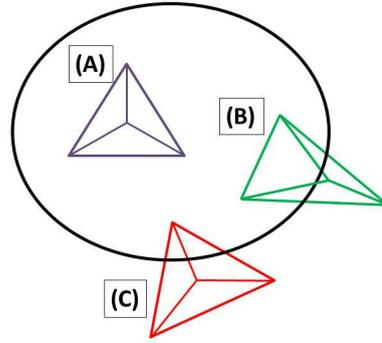


FIGURE 4.2: Coupling used in Lemma 4.11.

$$\begin{aligned}
 \mathbb{E} \left[d_t(s)^{(k)} \right] &= \mathbb{E} \left[\mathbb{E} \left[d_t(s)^{(k)} | d_{t-1}(s) \right] \right] \\
 &= \mathbb{E} \left[(d_{t-1}(s))^{(k)} \left(1 - \frac{d_{t-1}(s) - 1}{2t - 1} \right) + (d_{t-1}(s) + 1)^{(k)} \frac{d_{t-1}(s) - 1}{2t - 1} \right] \\
 &= \mathbb{E} \left[(d_{t-1}(s))^{(k)} \left(1 + \frac{k}{2t - 1} \right) - (d_{t-1}(s))^{(k)} \frac{k}{(2t - 1)d_{t-1}(s)} \right] \\
 &\leq \mathbb{E} \left[(d_{t-1}(s))^{(k)} \right] \left(1 + \frac{k}{2t - 1} \right) \leq \dots \leq \frac{(k + 2)!}{2} \left(\frac{2t}{s} \right)^{\frac{k}{2}}.
 \end{aligned}$$

□

Lemma 4.11. *The degree X_t of the supernode V_{t_0} of vertices added before time t_0 is at least $t_0^{1/4} \sqrt{t}$ whp.*

Proof. We consider a modified process \mathcal{Y} coupled with the RAN process, see also Figure 4.2. Specifically, let Y_t be the modified degree of the supernode in the modified process \mathcal{Y} which is defined as follows: for any type of insertion in the original RAN process –note there exist three types of insertions with respect to how the degree X_t of the supernode (black circle) gets affected, see also Figure 4.2– Y_t increases by 1. We also define $X_{t_0} = Y_{t_0}$. Note that $X_t \geq Y_t$ for all $t \geq t_0$. Let $d_0 = X_{t_0} = Y_{t_0} = 6t_0 + 6$ and $p^* = \Pr[Y_t = d_0 + r | Y_{t_0} = d_0]$.

Claim 1 (1).

$$p^* \leq \binom{d_0 + r - 1}{d_0 - 1} \left(\frac{2t_0 + 3}{2t + 1} \right)^{d_0/2} e^{\frac{3}{2} + t_0 - \frac{d_0}{2} + \frac{2r}{3\sqrt{t}}}$$

Proof. Let $\tau = (t_0 \equiv \tau_0, \underbrace{\tau_1, \dots, \tau_r}_{\text{insertion times}}, \tau_{r+1} \equiv t)$ be a vector denoting that Y_t increases by 1 at τ_i for $i = 1, \dots, r$. We upper bound the probability p_τ of this event in the

following. Note that we consider the case where the vertices have same degree as the number of faces around them. As we mentioned earlier, the other case is analyzed in exactly the same way, modulo a negligible error term.

$$\begin{aligned}
 p_r &= \left[\prod_{k=1}^r \frac{d_0 + k - 1}{2\tau_k + 1} \right] \left[\prod_{k=0}^r \prod_{j=\tau_k+1}^{\tau_{k+1}-1} \left(1 - \frac{d_0 + k}{2j + 1} \right) \right] \\
 &= d_0(d_0 + 1) \dots (d_0 + r - 1) \left[\prod_{k=1}^r \frac{1}{2\tau_k + 1} \right] \exp \left(\sum_{k=0}^r \sum_{j=\tau_k+1}^{\tau_{k+1}-1} \log \left(1 - \frac{d_0 + k}{2j + 1} \right) \right) \\
 &= \frac{(d_0 + r - 1)!}{(d_0 - 1)!} \left[\prod_{k=1}^r \frac{1}{2\tau_k + 1} \right] \exp \left(\sum_{k=0}^r \sum_{j=\tau_k+1}^{\tau_{k+1}-1} \log \left(1 - \frac{d_0 + k}{2j + 1} \right) \right)
 \end{aligned}$$

Consider now the inner sum which we upper bound using an integral:

$$\begin{aligned}
 \sum_{j=\tau_k+1}^{\tau_{k+1}-1} \log \left(1 - \frac{d_0 + k}{2j + 1} \right) &\leq \int_{\tau_k+1}^{\tau_{k+1}} \log \left(1 - \frac{d_0 + k}{2x + 1} \right) dx \\
 &\leq -\left(\tau_{k+1} + \frac{1}{2}\right) \log(2\tau_{k+1} + 1) + \\
 &\quad \frac{2\tau_{k+1} + 1 - (d_0 + k)}{2} \log(2\tau_{k+1} + 1 - (d_0 + k)) + \\
 &\quad \left(\tau_k + \frac{3}{2}\right) \log(2\tau_k + 3) - \frac{2\tau_k + 3 - (d_0 + k)}{2} \log(2\tau_k + 3 - (d_0 + k))
 \end{aligned}$$

since

$$\int \log \left(1 - \frac{d_0 + k}{2x + 1} \right) = -\left(x + \frac{1}{2}\right) \log(2x + 1) + \frac{2x + 1 - (d_0 + k)}{2} \log(2x + 1 - (d_0 + k))$$

Hence we obtain $\sum_{k=0}^r \sum_{j=\tau_k+1}^{\tau_{k+1}-1} \log \left(1 - \frac{d_0+k}{2j+1} \right) \leq A + \sum_{k=1}^r B_k$ where

$$\begin{aligned}
 A &= \left(\tau_0 + \frac{3}{2}\right) \log(2\tau_0 + 3) - \frac{2\tau_0 + 3 - d_0}{2} \log(2\tau_0 + 3 - d_0) \\
 &\quad - \left(\tau_{r+1} + \frac{1}{2}\right) \log(2\tau_{r+1} + 1) + \frac{2\tau_{r+1} + 1 - (d_0 + r)}{2} \log(2\tau_{r+1} + 1 - (d_0 + r))
 \end{aligned}$$

and

$$\begin{aligned}
 B_k &= \left(\tau_k + \frac{3}{2}\right) \log(2\tau_k + 3) - \frac{2\tau_k + 3 - (d_0 + k)}{2} \log(2\tau_k + 3 - (d_0 + k)) \\
 &\quad - \left(\tau_k + \frac{1}{2}\right) \log(2\tau_k + 1) + \frac{2\tau_k + 1 - (d_0 + k - 1)}{2} \log(2\tau_k + 1 - (d_0 + k - 1)).
 \end{aligned}$$

We first upper bound the quantities B_k for $k = 1, \dots, r$. By rearranging terms and using the identity $\log(1 + x) \leq x$ we obtain

$$\begin{aligned}
 B_k &= \left(\tau_k + \frac{1}{2}\right) \log\left(1 + \frac{1}{\tau_k + \frac{1}{2}}\right) + \log(2\tau_k + 3) \\
 &\quad - \frac{1}{2} \log(2\tau_k + 3 - (d_0 + k)) - \frac{2\tau_k + 2 - (d_0 + k)}{2} \log\left(1 + \frac{1}{2\tau_k + 2 - (d_0 + k)}\right). \\
 &\leq \frac{1}{2} + \frac{1}{2} \log(2\tau_k + 3) - \frac{1}{2} \log\left(1 - \frac{d_0 + k}{2\tau_k + 3}\right)
 \end{aligned}$$

First we rearrange terms and then we bound the term e^A by using the inequality $e^{-x-x^2/2} \geq 1 - x$ which is valid for $0 < x < 1$:

$$\begin{aligned}
 A &= -\left(\tau_0 + \frac{3}{2}\right) \log\left(1 - \frac{d_0}{2\tau_0 + 3}\right) + \left(\tau_{r+1} + \frac{1}{2}\right) \log\left(1 - \frac{d_0 + r}{2\tau_{r+1} + 1}\right) + \frac{d_0}{2} \log(2\tau_0 + 3 - d_0) \\
 &\quad - \frac{d_0 + r}{2} \log(2\tau_{r+1} + 1 - (d_0 + r)). \Rightarrow \\
 e^A &= \left(1 - \frac{d_0}{2\tau_0 + 3}\right)^{-(\tau_0 + \frac{3}{2})} \left(1 - \frac{d_0 + r}{2\tau_{r+1} + 1}\right)^{\tau_{r+1} + \frac{1}{2}} (2\tau_0 + 3 - d_0)^{\frac{d_0}{2}} (2\tau_{r+1} + 1 - (d_0 + r))^{-\frac{d_0 + r}{2}} \\
 &= \left(\frac{2\tau_0 + 3}{2\tau_{r+1} + 1}\right)^{d_0/2} (2\tau_{r+1} + 1)^{-r/2} \left(1 - \frac{d_0}{2\tau_0 + 3}\right)^{-(\tau_0 + \frac{3}{2}) + \frac{d_0}{2}} \left(1 - \frac{d_0 + r}{2\tau_{r+1} + 1}\right)^{\tau_{r+1} + \frac{1}{2} - \frac{d_0 + r}{2}} \\
 &\leq \left(\frac{2t_0 + 3}{2t + 1}\right)^{d_0/2} (2t + 1)^{-r/2} \left(1 - \frac{d_0}{2\tau_0 + 3}\right)^{-(\tau_0 + \frac{3}{2}) + \frac{d_0}{2}} e^{\left(-\frac{d_0 + r}{2t + 1} - \left(-\frac{d_0 + r}{2t + 1}\right)^2 / 2\right)(t + 1/2 - \frac{d_0 + r}{2})} \\
 &= \left(\frac{2t_0 + 3}{2t + 1}\right)^{d_0/2} (2t + 1)^{-r/2} \left(1 - \frac{d_0}{2\tau_0 + 3}\right)^{-(\tau_0 + \frac{3}{2}) + \frac{d_0}{2}} e^{-\frac{d_0 + r}{2} + \frac{(d_0 + r)^2}{8t + 4} + \frac{(d_0 + r)^3}{4(2t + 1)^2}}
 \end{aligned}$$

Now we upper bound the term $\exp(A + \sum_{k=1}^r B_k)$ using the above upper bounds:

$$\begin{aligned}
 e^{A+\sum_{k=1}^r B_k} &\leq e^A e^{r/2} \prod_{i=1}^r \sqrt{\frac{2\tau_k + 3}{1 - \frac{d_0+k}{2\tau_k+3}}} \\
 &\leq \left(1 - \frac{d_0}{2\tau_0 + 3}\right)^{-(\tau_0 + \frac{3}{2}) + \frac{d_0}{2}} e^{-\frac{d_0}{2} + \frac{(d_0+r)^2}{8t+4} + \frac{(d_0+r)^3}{4(2t+1)^2}} \left(\frac{2t_0 + 3}{2t + 1}\right)^{d_0/2} \times \\
 &\quad (2t + 1)^{-r/2} \prod_{i=1}^r \sqrt{\frac{2\tau_k + 3}{1 - \frac{d_0+k}{2\tau_k+3}}}
 \end{aligned}$$

Using the above upper bound we get that

$$p_\tau \leq C(r, d_0, t_0, t) \prod_{k=1}^r \left[(2\tau_k + 3 - (d_0 + k))^{-1/2} \left(1 + \frac{1}{\tau_k + 1/2}\right) \right]$$

where

$$C(r, d_0, t_0, t) = \frac{(d_0 + r - 1)!}{(d_0 - 1)!} \left(1 - \frac{d_0}{2\tau_0 + 3}\right)^{-(\tau_0 + \frac{3}{2}) + \frac{d_0}{2}} e^{-\frac{d_0}{2} + \frac{(d_0+r)^2}{8t+4} + \frac{(d_0+r)^3}{4(2t+1)^2}} \left(\frac{2t_0 + 3}{2t + 1}\right)^{d_0/2} (2t+1)^{-r/2}$$

We need to sum over all possible insertion times to bound the probability of interest p^* .

We set $\tau'_k \leftarrow \tau_k - \lceil \frac{d_0+k}{2} \rceil$ for $k = 1, \dots, r$. For $d = o(\sqrt{t})$ and $r = o(t^{2/3})$ we obtain:

$$\begin{aligned}
 p^* &\leq C(r, d_0, t_0, t) \sum_{t_0+1 \leq \tau_1 < \dots < \tau_r \leq t} \prod_{k=1}^r \left[(2\tau_k + 3 - (d_0 + k))^{-1/2} \left(1 + \frac{1}{\tau_k + 1/2} \right) \right] \\
 &\leq C(r, d_0, t_0, t) \sum_{t_0 - \lceil \frac{d_0}{2} \rceil + 1 \leq \tau'_1 \leq \dots \leq \tau'_r \leq t - \lceil \frac{d_0+r}{2} \rceil} \prod_{k=1}^r \left[(2\tau'_k + 3)^{-1/2} \left(1 + \frac{1}{\tau'_k + \frac{d_0+k}{2} + 1/2} \right) \right] \\
 &\leq \frac{C(r, d_0, t_0, t)}{r!} \left(\sum_{t_0 - \lceil \frac{d_0}{2} \rceil}^{t - \lceil \frac{d_0+r}{2} \rceil} (2\tau'_k + 3)^{-1/2} + \frac{1}{\sqrt{2}} (\tau'_k + 3/2)^{-3/2} \right)^r \\
 &\leq \frac{C(r, d_0, t_0, t)}{r!} \left(\int_0^{t - \frac{d_0+r}{2}} \left[(2x + 3)^{-1/2} + \frac{1}{\sqrt{2}} (x + 3/2)^{-3/2} \right] dx \right)^r \\
 &\leq \frac{C(r, d_0, t_0, t)}{r!} \left(\sqrt{2t + 3 - (d_0 + r)} + 2/3 \right)^r \\
 &\leq \frac{C(r, d_0, t_0, t)}{r!} (2t)^{r/2} e^{-\frac{r}{2} \frac{d_0+r-3}{2t}} e^{\frac{2r}{3\sqrt{2t - (d_0+r)+3}}} \\
 &\leq \binom{d_0 + r - 1}{d_0 - 1} \left(\frac{2t_0 + 3}{2t + 1} \right)^{d_0/2} \left[\left(1 - \frac{d_0}{2t_0 + 3} \right)^{-(1 - \frac{d_0}{2t_0 + 3})} \right]^{t_0 + 3/2} \times \\
 &\quad \left(\frac{2t}{2t + 1} \right)^{r/2} \exp \left(-\frac{d_0}{2} + \frac{(d_0 + r)^2}{8t + 4} + \frac{(d_0 + r)^3}{4(2t + 1)^2} - \frac{r(d_0 + r - 3)}{4t} + \frac{2r}{3\sqrt{2t + 3 - (d_0 + r)}} \right)
 \end{aligned}$$

By removing the $o(1)$ terms in the exponential and using the fact that $x^{-x} \leq e$ we obtain the following bound on the probability p^* .

$$p^* \leq \binom{d_0 + r - 1}{d_0 - 1} \left(\frac{2t_0 + 3}{2t + 1} \right)^{d_0/2} e^{\frac{3}{2} + t_0 - \frac{d_0}{2} + \frac{2r}{3\sqrt{t}}}.$$

□

Let \mathcal{A}_1 denote the event that the supernode consisting of the first t_0 vertices has degree Y_t in the modified process \mathcal{Y} less than $t_0^{1/4} \sqrt{t}$. Note that since $\{X_t \leq t_0^{1/4} \sqrt{t}\} \subseteq \{Y_t \leq t_0^{1/4} \sqrt{t}\}$ it suffices to prove that $\Pr[Y_t \leq t_0^{1/4} \sqrt{t}] = o(1)$. Using Claim (1) we obtain

$$\begin{aligned}
 \Pr[\mathcal{A}_1] &\leq \sum_{r=0}^{t_0^{1/4}\sqrt{t}-(6t_0+6)} \binom{r+6t_0+5}{6t_0+5} \left(\frac{2t_0+3}{2t+1}\right)^{3t_0+3} e^{-\frac{3}{2}-2t_0+\frac{2t_0^{1/4}}{3}} \\
 &\leq t_0^{1/4} t^{1/2} \frac{(t_0^{1/4} t^{1/2})^{6t_0+5}}{(6t_0+5)!} \left(\frac{2t_0+3}{2t+1}\right)^{3t_0+3} e^{-\frac{3}{2}-2t_0+\frac{2t_0^{1/4}}{3}} \\
 &\leq \left(\frac{t}{2t+1}\right)^{3t_0+3} \frac{t_0^{3t_0/2+3/2} (2t_0+3)^{3t_0+3}}{(6t_0+5)^{6t_0+5}} e^{4t_0+7/2+2/3t_0^{1/4}} \\
 &\leq 2^{-(3t_0+3)} \frac{e^{4t_0+7/2+2/3t_0^{1/4}}}{(6t_0+5)^{\frac{3}{2}t_0+\frac{1}{2}}} = o(1).
 \end{aligned}$$

□

Lemma 4.12. *No vertex added after t_1 has degree exceeding $t_0^{-2}t^{1/2}$ whp.*

Proof. Let \mathcal{A}_2 denote the event that some vertex added after t_1 has degree exceeding $t_0^{-2}t^{1/2}$. We use a union bound, a third moment argument and Lemma 4.10 to prove that $\Pr[\mathcal{A}_2] = o(1)$. Specifically

$$\begin{aligned}
 \Pr[\mathcal{A}_2] &\leq \sum_{s=t_1}^t \Pr[d_t(s) \geq t_0^{-2}t^{1/2}] = \sum_{s=t_1}^t \Pr[d_t(s)^{(3)} \geq (t_0^{-2}t^{1/2})^{(3)}] \\
 &\leq t_0^6 t^{-3/2} \sum_{s=t_1}^t \mathbb{E}[d_t(s)^{(3)}] \leq 5! \sqrt{2} t_0^6 \sum_{s=t_1}^t s^{-3/2} \leq 5! 2 \sqrt{2} t_0^6 t_1^{-1/2} = o(1).
 \end{aligned}$$

□

Lemma 4.13. *No vertex added before t_1 has degree exceeding $t_0^{1/6}t^{1/2}$ whp.*

Proof. Let \mathcal{A}_3 denote the event that some vertex added before t_1 has degree exceeding $t_0^{1/6}t^{1/2}$. We use again a third moment argument and Lemma 4.10 to prove that $\Pr[\mathcal{A}_3] = o(1)$.

$$\begin{aligned}
 \Pr[\mathcal{A}_3] &\leq \sum_{s=1}^{t_1} \Pr[d_t(s) \geq t_0^{1/6}t^{1/2}] = \sum_{s=1}^{t_1} \Pr[d_t(s)^{(3)} \geq (t_0^{1/6}t^{1/2})^{(3)}] \\
 &\leq t_0^{-1/2} t^{-3/2} \sum_{s=1}^{t_1} \mathbb{E}[d_t(s)^{(3)}] \leq t_0^{-1/2} t^{-3/2} \sum_{s=1}^{t_1} 5! \sqrt{2} \frac{t^{3/2}}{s^{3/2}} \\
 &\leq 5! \sqrt{2} \zeta(3/2) t_0^{-1/2} = o(1)
 \end{aligned}$$

where $\zeta(3/2) = \sum_{s=1}^{+\infty} s^{-3/2} \approx 2.612$.

□

Lemma 4.14. *The k highest degrees are added before t_1 and have degree Δ_i bounded by $t_0^{-1}t^{1/2} \leq \Delta_i \leq t_0^{1/6}t^{1/2}$ whp.*

Proof. For the upper bound it suffices to show that $\Delta_1 \leq t_0^{1/6}t^{1/2}$. This follows immediately by Lemmas 4.12 and 4.13. The lower bound follows directly from Lemmas 4.11, 4.12 and 4.13. Assume that at most $k - 1$ vertices added before t_1 have degree exceeding the lower bound $t_0^{-1}t^{1/2}$. Then the total degree of the supernode formed by the first t_0 vertices is $O(t_0^{1/6}\sqrt{t})$. This contradicts Lemma 4.11. Finally, since each vertex $s \geq t_1$ has degree at most $t_0^{-2}\sqrt{t} \ll t_0^{-1}t^{1/2}$ the k highest degree vertices are added before t_1 whp. □

The proof of Theorem 4.3 is completed with the following lemma.

Lemma 4.15. *The k highest degrees satisfy $\Delta_i \leq \Delta_{i-1} - \frac{\sqrt{t}}{f(t)}$ whp.*

Proof. Let \mathcal{A}_4 denote the event that there are two vertices among the first t_1 with degree $t_0^{-1}t^{1/2}$ and within $\frac{\sqrt{t}}{f(t)}$ of each other. By the definition of conditional probability and Lemma 4.12

$$\Pr[\mathcal{A}_4] = \Pr[\mathcal{A}_4|\bar{\mathcal{A}}_3] \Pr[\bar{\mathcal{A}}_3] + \Pr[\mathcal{A}_4|\mathcal{A}_3] \Pr[\mathcal{A}_3] \leq \Pr[\mathcal{A}_4|\bar{\mathcal{A}}_3] + o(1)$$

it suffices to show that $\Pr[\mathcal{A}_4|\bar{\mathcal{A}}_3] = o(1)$. Note that by a simple union bound

$$\Pr[\mathcal{A}_4] \leq \sum_{1 \leq s_1 < s_2 \leq t_1} \sum_{l=-\frac{\sqrt{t}}{f(t)}}^{\frac{\sqrt{t}}{f(t)}} p_{l,s_1,s_2} = O\left(t_1^2 \frac{\sqrt{t}}{f(t)} \max p_{l,s_1,s_2}\right)$$

where $p_{l,s_1,s_2} = \Pr[d_t(s_1) - d_t(s_2) = l|\bar{\mathcal{A}}_3]$.

We consider two cases and we show that in both cases $\max p_{l,s_1,s_2} = o\left(\frac{f(t)}{t_1^2\sqrt{t}}\right)$.

- CASE 1 $(s_1, s_2) \notin E(G_t)$:

Note that at time t_1 there exist $m_{t_1} = 3t_1 + 3 < 4t_1$ edges in G_{t_1} .

$$\begin{aligned}
 p_{l,s_1,s_2} &\leq \sum_{r=t_0^{-1}t^{1/2}}^{t_0^{1/6}t^{1/2}} \sum_{d_1,d_2=3}^{4t_1} \Pr [d_t(s_1) = r \wedge d_t(s_2) = r - l | d_{t_1}(s_1) = d_1, d_{t_1}(s_2) = d_2] \quad (2) \\
 &\leq t_0^{1/6}t^{1/2} \sum_{d_1,d_2=3}^{4t_1} \binom{2t_0^{1/6}t^{1/2}}{d_1-1} \binom{2t_0^{1/6}t^{1/2}}{d_2-1} \left(\frac{2t_0+3}{2t+1}\right)^{(d_1+d_2)/2} e^{\frac{3}{2}+t_1+\frac{2t_0^{1/6}}{3}} \quad (3) \\
 &\leq t_0^{1/6}t^{1/2} \sum_{d_1,d_2=3}^{4t_1} (2t_0^{1/6}t^{1/2})^{d_1+d_2-2} \left(\frac{2t_0+3}{2t+1}\right)^{(d_1+d_2)/2} e^{2t_1} \\
 &\leq t_0^{1/6}t^{1/2} e^{2t_1} t_1^2 (2t_0^{1/6}t^{1/2})^{8t_1-2} \left(\frac{2t_0+3}{2t+1}\right)^{4t_1} \\
 &= t_0^{4t_1/3+1/6} t^{-1/2} e^{2t_1} t_1^2 2^{8t_1} (2t_0+3)^{4t_1} \left(\frac{t}{2t+1}\right)^{4t_1} \\
 &= o\left(\frac{f(t)}{t_1^2 \sqrt{t}}\right)
 \end{aligned}$$

Note that we omitted the tedious calculation justifying the transition from (2) to (3) since calculating the upper bound of the joint probability distribution is very similar to the calculation of Lemma 4.11.

• CASE 2 $(s_1, s_2) \in E(G_t)$:

Notice that in any case (s_1, s_2) share at most two faces (which may change over time). Note that the two connected vertices s_1, s_2 share a common face only if $s_1, s_2 \in \{1, 2, 3\}$ ³. Consider the following modified process \mathcal{Y}' : whenever an incoming vertex “picks” one of the two common faces we don’t insert it. We choose two other faces which are not common to s_1, s_2 and add one vertex in each of those. Notice that the number of faces increases by 1 for both s_1, s_2 as in the original process and the difference of the degrees remains the same. An algebraic manipulation similar to Case 1 gives the desired result. \square

4.5 Proof of Theorem 4.4

Having computed the highest degrees of a RAN in Section 4.4, eigenvalues are computed by adapting existing techniques [113, 171, 302]. We decompose the proof of Theorem 4.4 in Lemmas 4.16, 4.17, 4.18, 4.19. Specifically, in Lemmas 4.16, 4.17 we bound the degrees and co-degrees respectively. Having these bounds, we decompose the graph into

³We analyze the case where $s_1, s_2 \geq 4$. The other case is treated in the same manner.

a star forest and show in Lemmas 4.18 and 4.19 that its largest eigenvalues, which are $(1 \pm o(1))\sqrt{\Delta_i}$, dominate the eigenvalues of the remaining graph. This technique was pioneered by Mihail and Papadimitriou [302].

We partition the vertices into three set S_1, S_2, S_3 . Specifically, let S_i be the set of vertices added after time t_{i-1} and at or before time t_i where

$$t_0 = 0, t_1 = t^{1/8}, t_2 = t^{9/16}, t_3 = t.$$

In the following we use the recursive variational characterization of eigenvalues [115]. Specifically, let A_G denote the adjacency matrix of a simple, undirected graph G and let $\lambda_i(G)$ denote the i -th largest eigenvalue of A_G . Then

$$\lambda_i(G) = \min_S \max_{x \in S, x \neq 0} \frac{x^T A_G x}{x^T x}$$

where S ranges over all $(n - i + 1)$ dimensional subspaces of \mathbb{R}^n .

Lemma 4.16. *For any $\epsilon > 0$ and any $f(t)$ with $f(t) \rightarrow +\infty$ as $t \rightarrow +\infty$ the following holds whp: for all s with $f(t) \leq s \leq t$, for all vertices $r \leq s$, then $d_s(r) \leq s^{\frac{1}{2} + \epsilon} r^{-\frac{1}{2}}$.*

Proof. Set $q = \lceil \frac{4}{\epsilon} \rceil$. We use Lemma 4.10, a union bound and Markov's inequality to obtain:

$$\begin{aligned} \Pr \left[\bigcup_{s=f(t)}^t \bigcup_{r=1}^s \{d_s(r) \geq s^{1/2+\epsilon} r^{-1/2}\} \right] &\leq \sum_{s=f(t)}^t \sum_{r=1}^s \Pr \left[d_s(r)^{(q)} \geq (s^{1/2+\epsilon} r^{-1/2})^{(q)} \right] \\ &\leq \sum_{s=f(t)}^t \sum_{r=1}^s \Pr \left[d_s(r)^{(q)} \geq (s^{-(q/2+q\epsilon)} r^{q/2}) \right] \\ &\leq \sum_{s=f(t)}^t \sum_{r=1}^s \frac{(q+2)!}{2} \left(\frac{2s}{r}\right)^{q/2} s^{-q/2} s^{-q\epsilon} r^{q/2} \\ &= \frac{(q+2)!}{2} 2^{q/2} \sum_{s=f(t)}^t s^{1-q\epsilon} \\ &\leq \frac{(q+2)!}{2} 2^{q/2} \int_{f(t)-1}^t x^{1-q\epsilon} dx \\ &\leq \frac{(q+2)!}{2(q\epsilon-2)} 2^{q/2} (f(t)-1)^{2-q\epsilon} = o(1). \end{aligned}$$

□

Lemma 4.17. *Let S'_3 be the set of vertices in S_3 which are adjacent to more than one vertex of S_1 . Then $|S'_3| \leq t^{1/6}$ whp.*

Proof. First, observe that when vertex s is inserted it becomes adjacent to more than one vertex of S_1 if the face chosen by s has at least two vertices in S_1 . We call the latter property \mathcal{A} and we write $s \in \mathcal{A}$ when s satisfies it. At time t_1 there exist $2t_1 + 1$ faces total, which consist of faces whose three vertices are all from S_1 . At time $s \geq t_2$ there can be at most $6t_1 + 3$ faces with at least two vertices in S_1 since each of the original $2t_1 + 1$ faces can give rise to at most 3 new faces with at least two vertices in s_1 . Consider a vertex $s \in S_3$, i.e., $s \geq t_2$. By the above argument, $\Pr[|N(s) \cap S_1| \geq 2] \leq \frac{6t_1+3}{2t_1+1}$. Writing $|S'_3|$ as a sum of indicator variables, i.e., $|S'_3| = \sum_{s=t_2}^t I(s \in \mathcal{A})$ and taking the expectation we obtain

$$\begin{aligned} \mathbb{E}[|S'_3|] &\leq \sum_{s=t_2}^t \frac{6t_1+3}{2t+1} \leq (6t_1+3) \int_{t_2}^t (2x+1)^{-1} dx \\ &\leq (3t^{\frac{1}{8}} + \frac{3}{2}) \ln \frac{2t+1}{2t_2+1} = o(t^{1/7}) \end{aligned}$$

By Markov's inequality:

$$\Pr[|S'_3| \geq t^{1/6}] \leq \frac{\mathbb{E}[|S'_3|]}{t^{1/6}} = o(1).$$

Therefore, we conclude that $|S'_3| \leq t^{1/6}$ whp. □

Lemma 4.18. *Let $F \subseteq G$ be the star forest consisting of edges between S_1 and $S_3 - S'_3$. Let $\Delta_1 \geq \Delta_2 \geq \dots \geq \Delta_k$ denote the k highest degrees of G . Then $\lambda_i(F) = (1 - o(1))\sqrt{\Delta_i}$ whp.*

Proof. It suffices to show that $\Delta_i(F) = (1 - o(1))\Delta_i(G)$ for $i = 1, \dots, k$. Note that since the k highest vertices are inserted before t_1 whp, the edges they lose are the edges between S_1 and the ones incident to S'_3 and S_2 and we know how to bound the cardinalities of all these sets. Specifically by Lemma 4.17 $|S'_3| \leq t^{1/6}$ whp and by Theorem 4.3 the maximum degree in G_{t_1}, G_{t_2} is less than $t_1^{1/2+\epsilon_1} = t^{1/8}, t_2^{1/2+\epsilon_2} = t^{5/16}$ for $\epsilon_1 = 1/16, \epsilon_2 = 1/32$ respectively whp. Also by Theorem 4.3, $\Delta_i(G) \geq \frac{\sqrt{t}}{\log t}$. Hence, we obtain

$$\Delta_i(F) \geq \Delta_i(G) - t^{1/8} - t^{5/16} - t^{1/6} = (1 - o(1))\Delta_i(G).$$

□

To complete the proof of Theorem 4.4 it suffices to prove that $\lambda_1(H)$ is $o(\lambda_k(F))$ where $H = G - F$. We prove this in the following lemma. The proof is based on bounding maximum degree of appropriately defined subgraphs using Lemma 4.16 and standard inequalities from spectral graph theory [115].

Lemma 4.19. $\lambda_1(H) = o(t^{1/4})$ whp.

Proof. From Gershgorin's theorem [375] the maximum eigenvalue of any graph is bounded by the maximum degree. We bound the eigenvalues of H by bounding the maximum eigenvalues of six different induced subgraphs. Specifically, let $H_i = H[S_i]$, $H_{ij} = H(S_i, S_j)$ where $H[S]$ is the subgraph induced by the vertex set S and $H(S, T)$ is the subgraph containing only edges with one vertex in S and other in T . We use Lemma 4.18 to bound $\lambda_1(H(S_1, S_3))$ and Lemma 4.17 for the other eigenvalues. We set $\epsilon = 1/64$.

$$\begin{aligned} \lambda_1(H_1) &\leq \Delta_1(H_1) \leq t_1^{1/2+\epsilon} = t^{33/512}. \\ \lambda_1(H_2) &\leq \Delta_1(H_2) \leq t_2^{1/2+\epsilon} t_1^{-1/2} = t^{233/1024}. \\ \lambda_1(H_3) &\leq \Delta_1(H_3) \leq t_3^{1/2+\epsilon} t_2^{-1/2} = t^{15/64}. \\ \lambda_1(H_{12}) &\leq \Delta_1(H_{12}) \leq t_2^{1/2+\epsilon} = t^{297/1024}. \\ \lambda_1(H_{23}) &\leq \Delta_1(H_{23}) \leq t_3^{1/2+\epsilon} t_1^{-1/2} = t^{29/64}. \\ \lambda_1(H_{13}) &\leq \Delta_1(H_{13}) \leq t^{1/6}. \end{aligned}$$

Therefore whp we obtain

$$\lambda_1(H) \leq \sum_{i=1}^3 \lambda_1(H_i) + \sum_{i<j} \lambda_1(H_{i,j}) = o(t^{1/4}).$$

□

4.6 Proof of Theorem 4.5

Before we give the proof of Theorem 4.5, we give a simple proof that the diameter of a RAN is $O(\log t)$ whp.

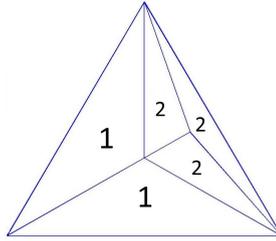


FIGURE 4.3: An instance of the process for $t = 2$. Each face is labeled with its depth.

We begin with a necessary definition for the proof of Claim (2). We define the *depth of a face* recursively. Initially, we have three faces, see Figure 4.1(a), whose depth equals 1. For each new face β created by subdividing a face γ , we have $depth(\beta) = depth(\gamma) + 1$. An example is shown in Figure 4.3, where each face is labeled with its corresponding depth.

Claim 2 (2). The diameter $d(G_t)$ satisfies $d(G_t) = O(\log t)$ *whp*.

Proof. A simple but key observation is that if k^* is the maximum depth of a face then $d(G_t) = O(k^*)$. Hence, we need to upper bound the depth of a given face after t rounds. Let $F_t(k)$ be the number of faces of depth k at time t , then:

$$\mathbb{E}[F_t(k)] = \sum_{1 \leq t_1 < t_2 < \dots < t_k \leq t} \prod_{j=1}^k \frac{1}{2t_j + 1} \leq \frac{1}{k!} \left(\sum_{j=1}^t \frac{1}{2j + 1} \right)^k \leq \frac{1}{k!} \left(\frac{1}{2} \log t \right)^k \leq \left(\frac{e \log t}{2k} \right)^{k+1}$$

By the first moment method we obtain $k^* = O(\log t)$ *whp* and by our observation $d(G_t) = O(\log t)$ *whp*. □

The depth of a face can be formalized via a bijection between random ternary trees and RANs. Using this bijection we prove Theorem 4.5 which gives a refined upper bound on the asymptotic growth of the diameter.

Proof. Consider the random process which starts with a single vertex tree and at every step picks a random leaf and adds three children to it. Let T be the resulting tree after t steps. There exists a natural bijection between the RAN process and this process, see [128] and also Figure 4.4. The depth of T in probability is $\frac{t}{2} \log t$ where $\frac{1}{\rho} = \eta$ is the unique solution less than 1 of the equation $\eta - 1 - \log \eta = \log 3$, see Broutin and Devroye

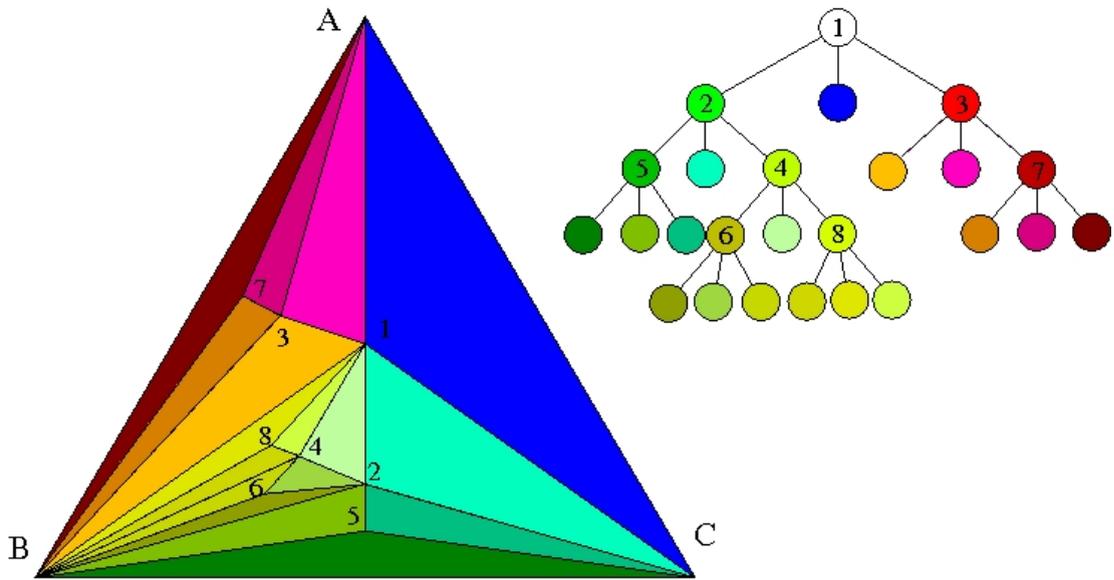


FIGURE 4.4: RANs as random ternary trees.

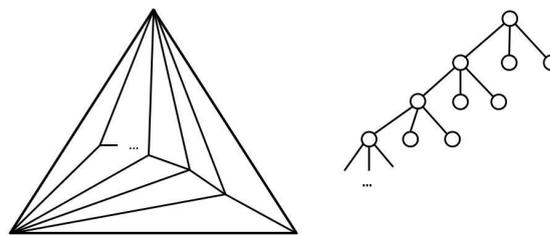


FIGURE 4.5: The height of the random ternary tree cannot be used to lower bound the diameter. The height of the random ternary tree can be arbitrarily large but the diameter is 2.

[93], pp. 284-285⁴. Note that the diameter $d(G_t)$ is at most twice the height of the tree and hence the result follows.

□

The above observation, i.e., the bijection between RANs and random ternary trees cannot be used to lower bound the diameter. A counterexample is shown in Figure 4.5 where the height of the random ternary tree can be made arbitrarily large but the diameter is 2. Albenque and Marckert proved in [297] that if v, u are two i.i.d. uniformly random internal vertices, i.e., $v, u \geq 4$, then the distance $d(u, v)$ tends to $\frac{6}{11} \log n$ with probability 1 as the number of vertices n of the RAN grows to infinity. Finally, it is worth mentioning that the diameter of the RAN grows faster asymptotically than the diameter of the classic preferential attachment model [50] which *whp* grows as $\frac{\log t}{\log \log t}$, see Bollobás and Riordan [79].

⁴ There is a typo in [93]. Specifically it says “ ρ is the unique solution greater than 1 of ...”. However, based on their Theorem 1, they should replace “greater than” with “less than”. Thanks to Abbas Mehrabian for pointing this out.

Chapter 5

Triangle Counting in Large Graphs

5.1 Introduction

In this Chapter we focus on algorithmic techniques for approximate triangle counting in large graphs. Despite the fact that the subgraph of interest is a triangle, our techniques are extendable to counting other types of fixed-size subgraphs. The outline of this Chapter is as follows: in Section 5.2 we present *triangle sparsifiers*, a notion inspired by the seminal work of Benczúr and Karger [64] and Spielman and Srivastava [370]. Specifically, we present a new randomized algorithm for approximately counting the number of triangles in a graph G . The algorithm proceeds as follows: keep each edge independently with probability p , enumerate the triangles in the sparsified graph G' and return the number of triangles found in G' multiplied by p^{-3} . We prove that under mild assumptions on G and p our algorithm returns a good approximation for the number of triangles with high probability.

We illustrate the efficiency of our algorithm on various large real-world datasets where we achieve significant speedups. Furthermore, we investigate the performance of existing sparsification procedures namely the Spielman-Srivastava spectral sparsifier [370] and the Benczúr-Karger cut sparsifier [63, 64] and show that they are not optimal/suitable with respect to triangle counting.

In Section 5.3 we propose a new triangle counting method which provides a $(1 \pm \epsilon)$ multiplicative approximation to the number of triangles in the graph and runs in $O\left(m + \frac{m^{3/2} \log n}{t\epsilon^2}\right)$ time. The key idea of the method is to combine the sampling scheme introduced by Tsourakakis et al. in [394, 396] with the partitioning idea of Alon, Yuster

and Zwick [25] in order to obtain a more efficient sampling scheme. Furthermore, we show that this method can be adapted to the semistreaming model [164] with a constant number of passes and $O\left(m^{1/2} \log n + \frac{m^{3/2} \log n}{t\epsilon^2}\right)$ space. We apply our methods in various networks with several millions of edges and we obtain excellent results both with respect to the accuracy and the running time. Finally, we propose a random projection based method for triangle counting and provide a sufficient condition to obtain an estimate with low variance. Even if such a method is unlikely to be practical it raises some interesting theoretical issues.

In Section 5.4 we present an (almost) optimal algorithm for triangle counting. The proposed randomized algorithm is analyzed via the second moment method and provides tight theoretical guarantees. We discuss various aspects of the proposed algorithm, including an implementation in MAPREDUCE.

5.2 Triangle Sparsifiers

5.2.1 Proposed Algorithm

Algorithm 1 Triangle Sparsifier

Input: Set of edges $E \subseteq \binom{[n]}{2}$
 {Unweighted graph $G([n], E)$ }

Input: Sparsification parameter p

Pick a random subset E' of edges such that the events $\{e \in E'\}$, for all $e \in E$ are independent and the probability of each is equal to p .

$t' \leftarrow$ count triangles on the graph $G'([n], E')$

Return $T \leftarrow \frac{t'}{p^3}$

Our proposed algorithm *Triangle Sparsifier* is shown in Algorithm 1. The algorithm takes an unweighted, simple graph $G(V, E)$, where without loss of generality we assume $V = [n]$, and a sparsification parameter $p \in (0, 1)$ as input. The algorithm first chooses a random subset E' of the set E of edges. The random subset is such that the events

$$\{e \in E'\}, \text{ for all } e \in E,$$

are independent and the probability of each is equal to p . Then, any triangle counting algorithm can be used to count triangles on the sparsified graph with edge set E' . Clearly, the expected size of E' is pm . The output of our algorithm is the number of triangles in the sparsified graph multiplied by $\frac{1}{p^3}$, or equivalently we are counting the number of weighted triangles in G' where each edge has weight $\frac{1}{p}$. It follows immediately that the

expected value $\mathbb{E}[T]$ of our estimate is the number of triangles in G , i.e., t . Our main theoretical result is the following theorem:

Theorem 5.1. *Suppose G is an undirected graph with n vertices, m edges and t triangles. Let also Δ denote the size of the largest collection of triangles with a common edge. Let G' be the random graph that arises from G if we keep every edge with probability p and write T for the number of triangles of G' . Suppose that $\gamma > 0$ is a constant and*

$$\frac{pt}{\Delta} \geq \log^{6+\gamma} n, \quad \text{if } p^2\Delta \geq 1, \quad (5.1)$$

and

$$p^3t \geq \log^{6+\gamma} n, \quad \text{if } p^2\Delta < 1. \quad (5.2)$$

for $n \geq n_0$ sufficiently large. Then

$$\Pr[|T - \mathbb{E}[T]| \geq \epsilon \mathbb{E}[T]] \leq n^{-K}$$

for any constants $K, \epsilon > 0$ and all large enough n (depending on K, ϵ and n_0).

Proof. Write $X_e = 1$ or 0 depending on whether the edge e of graph G survives in G' . Then $T = \sum_{\Delta(e,f,g)} X_e X_f X_g$ where $\Delta(e, f, g) = \mathbf{1}$ (edges e, f, g form a triangle). Clearly $\mathbb{E}[T] = p^3t$.

Refer to Theorem 2.4. We use T in place of Y , $k = 3$.

We have

$$\mathbb{E} \left[\frac{\partial T}{\partial X_e} \right] = \sum_{\Delta(e,f,g)} \mathbb{E}[X_f X_g] = p^2 |\Delta(e)|.$$

We first estimate the quantities $\mathbb{E}_j(T)$, $j = 0, 1, 2, 3$, defined before Theorem 2.4. We get

$$\mathbb{E}_1(T) = p^2 \Delta. \quad (5.3)$$

We also have

$$\mathbb{E} \left[\frac{\partial^2 T}{\partial X_e \partial X_f} \right] = p \mathbf{1}(\exists g : \Delta(e, f, g)),$$

hence

$$\mathbb{E}_2(T) \leq p. \quad (5.4)$$

Obviously, $\mathbb{E}_3(T) \leq 1$.

Hence

$$\mathbb{E}_{\geq 3}(T) \leq 1, \quad \mathbb{E}_{\geq 2}(T) \leq 1,$$

and

$$\mathbb{E}_{\geq 1}(T) \leq \max\{1, p^2\Delta\}, \quad \mathbb{E}_{\geq 0}(T) \leq \max\{1, p^2\Delta, p^3t\}.$$

• CASE 1 ($p^2\Delta < 1$):

We get $\mathbb{E}_{\geq 1}(T) \leq 1$, and from (5.2), $\mathbb{E}[T] \geq \mathbb{E}_{\geq 1}(T)$.

• CASE 2 ($p^2\Delta \geq 1$):

We get $\mathbb{E}_{\geq 1}(T) \leq p^2\Delta$ and, from (5.1), $\mathbb{E}[T] \geq \mathbb{E}_{\geq 1}(T)$.

We get, for some constant $c_3 > 0$, from Theorem 2.4:

$$\Pr\left[|T - \mathbb{E}[T]| \geq c_3\lambda^3(\mathbb{E}[T]\mathbb{E}_{\geq 1}(T))^{1/2}\right] \leq e^{-\lambda+2\log n}. \quad (5.5)$$

Notice that since in both cases we have $\mathbb{E}[T] \geq \mathbb{E}_{\geq 1}(T)$.

We now select λ so that the lower bound inside the probability on the left-hand side of (5.5) becomes $\epsilon\mathbb{E}[T]$. In Case 1 we pick

$$\lambda = \frac{\epsilon^{1/3}}{c_3^{1/3}}(p^3t)^{1/6}$$

while in Case 2

$$\lambda = \frac{\epsilon^{1/3}}{c_3^{1/3}}\left(\frac{pt}{\Delta}\right)^{1/6}$$

to get

$$\Pr[|T - \mathbb{E}[T]| \geq \epsilon\mathbb{E}[T]] \leq \exp(-\lambda + 2\log n) \quad (5.6)$$

Since $\lambda \geq (K+2)\log n$ follows from our assumptions (5.1) and (5.2) if n is sufficiently large, we get $\Pr[|T - \mathbb{E}[T]| \geq \epsilon\mathbb{E}[T]] \leq n^{-K}$, in both cases. \square

Complexity Analysis: The expected running time of edge sampling is sublinear, i.e., $O(pm)$, see Claim 3. The complexity of the counting step depends on which algorithm we use to count triangles¹. For instance, if we use [25] as our triangle counting algorithm, the expected running time of Triangle Sparsifier is $O(pm + (pm)^{\frac{2\omega}{\omega+1}})$, where ω currently is 2.3727 [416]. If we use the node-iterator (or any other standard listing triangle algorithm) the expected running time is $O(pm + p^2 \sum_i d_i^2)$.

Claim 3 (Sparsification in sublinear expected time). The edge sampling can run in $O(pm)$ expected time.

¹We assume for fairness that we use the same algorithm in both the original graph G and the sparsified graph G' to count triangles.

Proof. We do not “toss a p -coin” m times in order to construct E' . This would be very wasteful if p is small. Instead we construct the random set E' with the following procedure which produces the right distribution. Observe that the number X of unsuccessful events, i.e., edges which are not selected in our sample, until a successful one follows a geometric distribution. Specifically, $\Pr[X = x] = (1 - p)^{x-1}p$. To sample from this distribution it suffices to generate a uniformly distributed variable U in $[0, 1]$ and set $X \leftarrow \left\lceil \frac{\ln U}{1-p} \right\rceil$. Clearly the probability that $X = x$ is equal to $\Pr[(1 - p)^{x-1} > U \geq (1 - p)^x] = (1 - p)^{x-1} - (1 - p)^x = (1 - p)^{x-1}p$ as required. This provides a practical and efficient way to pick the subset E' of edges in sublinear expected time $O(pm)$. For more details see [256]. \square

Expected Speedup: The expected speedup with respect to the triangle counting task depends on the triangle counting subroutine that we use. If we use [25] as our subroutine which is the fastest known algorithm, the expected speedup is $p^{-\frac{2\omega}{\omega+1}}$, i.e., currently $p^{-1.407}$ where ω currently is 2.3727 [416]. As already outlined, in practice $p^{-\frac{2\omega}{\omega+1}}$, i.e., currently $p^{-1.41}$, and p^{-2} respectively.

Discussion: This theorem states the important result that the estimator of the number of triangles is concentrated around its expected value, which is equal to the actual number of triangles t in the graph under mild conditions on the triangle density of the graph. The mildness comes from condition (5.1): picking $p = 1$, given that our graph is not triangle-free, i.e., $\Delta \geq 1$, gives that the number of triangles t in the graph has to satisfy $t \geq \Delta \log^{6+\gamma} n$. This is a mild condition on t since $\Delta \leq n$ and thus it suffices that $t \geq n \log^{6+\gamma} n$ (after all, we can always add two dummy connected nodes that connect to every other node, as in Figure 1(a), even if in empirically Δ is smaller than n). The critical quantity besides the number of triangles t , is Δ . Intuitively, if the sparsification procedure throws away the common edge of many triangles, the triangles in the resulting graph may differ significantly from the original. A significant problem is the choice of p for the sparsification. Conditions (5.1) and (5.2) tell us how small we can afford to choose p , but the quantities involved, namely t and Δ , are unknown. We discuss a practical algorithm using a doubling procedure in Section 5.2.2.4. Furthermore, our method justifies significant speedups. For a graph G with $t \geq n^{3/2+\epsilon}$ and $\Delta \sim n$, we get $p = n^{-1/2}$ implying a linear expected speedup if we use a practical exact counting method as the node iterator. Finally, it is worth pointing out that *Triangle Sparsifier* essentially outputs a sparse graph $H(V, E', w)$ with $w = 1/p$ for all edges $e \in E'$ which approximates $G(V, E)$ with respect to the count of triangles (a triangle formed by the edges (e_1, e_2, e_3) in a weighted graph counts for $w(e_1)w(e_2)w(e_3)$ unweighted triangles). As we shall see in Chapter 13, see Figure 13.2, *Triangle Sparsifier* is not recommended for weighted graphs. Finally, it is worth mentioning that the sparsification scheme which has

Description	Availability
SNAP	http://snap.stanford.edu/
UF Sparse Matrix Collection	http://www.cise.ufl.edu/research/sparse
Max Planck	http://socialnetworks.mpi-sws.org/

TABLE 5.1: Dataset sources.

been used for speeding up the computation of linear algebraic decompositions [11, 389] has also been used to count triangles based on spectral properties of real-world networks [390, 393, 395].

5.2.2 Experimental Results

In this Section we present our experimental findings. Specifically, in Section 5.2.2.1 we describe the datasets we used, in Section 5.2.2.2 we give details with respect to the experimental setup and in Section 5.2.2.3 the experimental results.

5.2.2.1 Datasets

The graphs we used with the exceptions of Livejournal-links and Flickr are available on the Web. Table 5.1 summarizes the data resources. We preprocessed the graphs by first making them undirected and removing all self-loops. Furthermore, a common phenomenon was to have multiple edges in the edge file, i.e., a file whose each line corresponds to an edge, despite the fact that the graphs were claimed to be simple. Those multiple edges were removed. Table 5.2 summarizes the datasets we used after the preprocessing.

5.2.2.2 Experimental Setup

The experiments were performed on a single machine, with Intel Xeon CPU at 2.83 GHz, 6144KB cache size and 50GB of main memory. The algorithm was implemented in C++, and compiled using gcc version 4.1.2 and the -O3 optimization flag. Time was measured by taking the user time given by the linux time command. IO times are included in that time since the amount of memory operations performed in setting up the graph is non-trivial. However, we use a modified IO routine that's much faster than the standard C/C++ scanf. Furthermore, as we mentioned in Section 5.2.1 picking a random subset of expected size $p|S|$ from a set S can be done in expected sublinear time [256]. A simple way to do this in practice is to generate the differences between indices of entries retained. This allows us to sample in a sequential way and also results

Name (Abbr.)	Nodes	Edges	Triangle Count
⊙ AS-Skitter (AS)	1,696,415	11,095,298	28,769,868
★Flickr (FL)	1,861,232	15,555,040	548,658,705
★Livejournal-links (LJ)	5,284,457	48,709,772	310,876,909
★Orkut-links (OR)	3,072,626	116,586,585	621,963,073
★Soc-LiveJournal (SL)	4,847,571	42,851,237	285,730,264
★Youtube (YOU)	1,157,822	2,990,442	4,945,382
◇Web-EDU (WE)	9,845,725	46,236,104	254,718,147
◇Web-Google (WG)	875,713	3,852,985	11,385,529
◇Wikipedia 2005/11 (W0511)	1,634,989	18,540,589	44,667,095
◇Wikipedia 2006/9 (W0609)	2,983,494	35,048,115	84,018,183
◇Wikipedia 2006/11 (W0611)	3,148,440	37,043,456	88,823,817
◇Wikipedia 2007/2 (W0702)	3,566,907	42,375,911	102,434,918

TABLE 5.2: Datasets used in our experiments. Abbreviations are included. Symbol ⊙ stands for Autonomous Systems graphs, ★ for online social networks and ◇ for Web graphs. Notice that the networks with the highest triangle counts are online social networks (Flickr, Livejournal, Orkut), verifying the folklore that online social networks are abundant in triangles.

in better cache performance. As a competitor we use the single pass algorithm of [96, § 2.2].

5.2.2.3 Experimental Results

Table 5.2 shows the count of triangles for each graph used in our experiments. Notice that Orkut, Flickr and Livejournal graphs have $\sim 622\text{M}$, 550M and 311M triangles respectively. This confirms the folklore that online social networks are abundant in triangles. Table 5.3 shows the results we obtain for $p = 0.1$ over 5 trials. All running times are reported in seconds. The first column shows the running time for the exact counting algorithm over 5 runs. Standard deviations are negligible for the exact algorithm and therefore are not reported. The second and third column show the error and running time averaged over 5 runs for each dataset (two decimal digits of accuracy). Standard deviations are also included (three decimal digits of accuracy). The last column shows the running time averaged over 5 runs for the 1-pass algorithm as stated in [96, §2.2] and the standard deviations. For each dataset the number of samples needed by the 1-pass algorithm was set to a value that achieves *at most* as good accuracy as the ones achieved by our counting method. Specifically, for any dataset, if $\alpha, \beta(\%)$ are the errors obtained by our algorithm and the Buriol et al. algorithm, we “tune” the number of samples in the latter algorithm in such way that $\alpha \leq \beta \leq \alpha + 1\%$. Even by favoring in this way the 1-pass algorithm of Buriol et al. [96], one can see that the running times achieved by our method are consistently better. However, it is important to outline once again that

our method and other triangle counting methods can be combined. For example, in Section 5.3 we show that Triangle Sparsifiers and other sampling methods can be combined to obtain a superior performance both in practice and theory by improving the sampling scheme of Buriol et al. [96]. As we will see in detail, this is achieved by distinguishing vertices into two subsets according to their degree and using two sampling schemes, one for each subset [257, 311]. We also tried other competitors, but our running times outperform them significantly. For example, even the exact counting method outperforms other approximate counting methods. As we show in Section 5.2.2.4 smaller values of p values work as well and these can be found by a simple procedure.

	Results			
	Exact	Triangle Sparsifier		Buriol et al. [96]
	Avg. time	Avg. err.% (std)	Avg. time (std)	Avg. time (std)
AS	4.45	2.60 (0.022)	0.79 (0.023)	2.72 (0.128)
FL	41.98	0.11 (0.003)	0.96 (0.014)	3.40 (0.175)
LJ	50.83	0.34 (0.001)	2.85 (0.054)	12.40 (0.250)
OR	202.01	0.60 (0.004)	5.60 (0.159)	11.71 (0.300)
SL	38.27	8.27 (0.006)	2.50 (0.032)	8.92 (0.115)
YOU	1.35	1.50 (0.050)	0.30 (0.002)	10.91 (0.130)
WE	8.50	0.70 (0.005)	2.79 (0.090)	6.56 (0.025)
WG	1.60	1.58 (0.011)	0.40 (0.004)	1.85 (0.047)
W0511	32.47	1.53 (0.010)	1.19 (0.020)	3.71 (0.038)
W0609	86.62	0.40 (0.055)	2.07 (0.014)	8.10 (0.040)
W0611	96.11	0.62 (0.008)	2.16 (0.042)	7.90 (0.090)
W0702	122.34	0.80 (0.015)	2.48 (0.012)	11.00 (0.205)

TABLE 5.3: Results of experiments averaged over 5 trials using $p = 0.1$. All running times are reported in seconds. The first column shows the running time for the exact counting algorithm averaged over 5 runs. The second and third column show the error and running time averaged over 5 runs for each dataset (two decimal digits of accuracy). Standard deviations are also included (three decimal digits of accuracy). The last column shows the running time averaged over 5 runs for the 1-pass algorithm as stated in [96, §2.2] and the corresponding standard deviations. The number of samples for each dataset was set to a value that achieves *at most* as good accuracy as the ones achieved by our counting method. See Section 5.2.2.3 for all the details.

5.2.2.4 The “Doubling” Algorithm

As we saw in Section 5.2.1, setting optimally the parameter p requires knowledge of the quantity we want to estimate, i.e., the number of triangles. To overcome this problem we observe that when we have concentration, the squared coefficient of variation $\frac{\text{Var}[T]}{\mathbb{E}[T]^2}$ is “small”. Furthermore, by the Chebyshev inequality and by the median boosting trick [224] it suffices to sample $\{T_1, \dots, T_s\}$ where $s = O\left(\frac{\text{Var}[T]}{\mathbb{E}[T]^2} \frac{1}{\epsilon^2} \ln \frac{1}{\delta}\right)$ in order to obtain a $(1 \pm \epsilon)$ approximation $\mathbb{E}[T] = t$ with probability at least $1 - \delta$. Hence, one can set a desired value for the number of samples s and of the failure probability δ and calculate

the expected error $\epsilon = O\left(\sqrt{\frac{\text{Var}[T]}{\mathbb{E}[T]^2} \frac{1}{s} \ln \frac{1}{\delta}}\right)$. If this value is significantly larger than the desired error threshold then one increases p and repeats the same procedure until the stopping criterion is satisfied. One way one can change p is to use the multiplicative rule $p \leftarrow cp$, where $c > 1$ is a constant. For example, if $c = 2$ then we have a doubling procedure. Notice that we've placed the word doubling in the title of this section in quotes in order to emphasize that one may use any $c > 1$ to change p from one round to the next.

For how many rounds can this procedure run? Let's consider the realistic scenario where one wishes to be optimistic and picks as an initial guess for p a value $p_0 = n^{-\alpha}$ where α is a positive constant, e.g., $\alpha = 1/2$. Let p^* be the minimum value over all possible p with the property that for p^* we obtain a concentrated estimate of the true number of triangles. Clearly, $p^* \leq 1$ and hence the number of rounds performed by our procedure is less than r where $p_0 c^r = 1$. Hence, for any constant $c > 1$ we obtain that the number of rounds performed by our algorithm is $O(\log n)$. Furthermore, note that the running time of the doubling procedure is dominated by the last iteration. To see why, consider for simplicity the scenario where $r + 1$ rounds are needed to deduce concentration, $c = \sqrt{2}$ and the use of the node-iterator algorithm to count triangles in the triangle sparsifier. Then, the total running time shall be $p_0^2 \sum_{v \in V(G)} d(v)^2 \left(1 + 2 + \dots + 2^{r-1} + 2^r\right)$. Finally, observe that $1 + 2 + \dots + 2^{r-1} = O(2^r)$. In practice, this procedure works even for small values of s . An instance of this procedure with $s = 2$, $\delta = 1/100$ and error threshold equal to 3% is shown in Table 5.4.

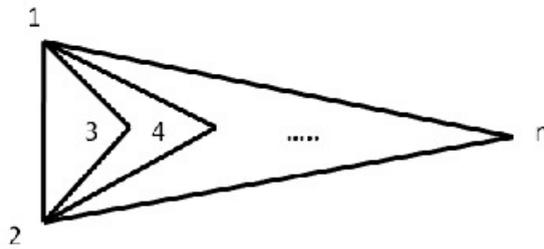
p	$\{T_1, T_2\}$	$\sqrt{\frac{\text{Var}[T]}{\mathbb{E}[T]^2} \frac{1}{s} \ln \frac{1}{\delta}}$	err(%)
0.01	{42, 398, 007 & 50, 920, 488}	0.1960	4.46
0.02	{42, 540, 941 & 43, 773, 753}	0.0307	3.38
0.04	{44, 573, 294 & 43, 398, 549}	0.0287	1.52

TABLE 5.4: Doubling procedure for the Wikipedia 2005 graph with 44,667,095 triangles.

5.2.3 Theoretical Ramifications

In this Section we investigate the performance of the Benczúr-Karger cut sparsifier and the Spielman-Srivastava spectral sparsifier with respect to triangle counting.

Consider the graph G shown in Figure 5.1. The strong edge connectivity of any edge in the graph is 2 and therefore the Benczúr-Karger algorithm does not distinguish the importance of the edge $e = (1, 2)$ with respect to triangle counting. The Spielman-Srivastava sparsifier with probability $1 - o(1)$ throws away the critical edge $e = (1, 2)$

FIGURE 5.1: Graph with linear number $O(n)$ of triangles.

as the number of vertices n tends to infinity. To prove this claim, we need use Foster's theorem 2.9.

Claim 4. The effective resistance $R(1, 2)$ of the edge $(1, 2)$ tends to 0 as n grows to infinity, i.e., $R(1, 2) = o(1)$. Furthermore, all other edges have constant effective resistance.

Proof. Using the in-series and in-parallel network simplification rules [76], the effective conductance of the edge $(1, 2)$ is $1 + \sum_{i=1}^{n-2} \frac{1}{2} = \frac{n}{2}$. Hence, the effective resistance of the edge $e = (1, 2)$ is $2/n$, which also proves the first part of our claim. By Foster's theorem 2.9, the sum of the effective resistances of the edges of G is $n - 1$. Due to symmetry, we obtain that $R_{(1,3)} = R_{(2,3)} = R_{(1,4)} = R_{(2,4)} = \dots = R_{(1,n)} = R_{(2,n)} = R_n$. Therefore we obtain $\frac{2}{n} + 2(n-2)R_n = n - 1 \rightarrow R_n = \frac{n^2 - n - 2}{2n^2 - 4n}$. Asymptotically as $n \rightarrow +\infty$, $R_n \rightarrow \frac{1}{2}$. \square

Clearly, the Spielman-Srivastava sparsifier fails to capture the importance of the edge $(1, 2)$ with respect to triangle counting. Finding an easy-to-compute quantity which allows a sparsification that preserves triangles more efficiently is an interesting problem. It is worth outlining that our analysis does not exclude effective resistances which can be computed very efficiently [259], but the use of them as is typically done in the context of spectral sparsifiers.

5.3 Efficient Triangle Counting in Large Graphs via Degree-based Vertex Partitioning

5.3.1 Proposed Method

Our algorithm combines two approaches that have been taken on triangle counting: sparsify the graph by keeping a random subset of the edges, see Section 5.2, followed by a triple sampling using the idea of vertex partitioning due to Alon, Yuster and Zwick

[29]. In the following, we shall assume that the input is in the form of an edge file, i.e., a file whose each line contains an edge. Notice that given this representation, computing the degrees takes linear time.

5.3.1.1 Edge Sparsification

As we saw in Section 5.2 the following method performs very well in practice: keep each edge with probability p independently. Then for each triangle, the probability of being kept is p^3 . So the expected number of triangles left is $p^3 t$. This is an inexpensive way to reduce the size of the graph as it can be done in one pass over the edge list using $O(mp)$ random variables.

We also proved that from the number of triangles in the sampled graph we can obtain a concentrated estimate around the actual triangle count as long as $p^3 \geq \tilde{\Omega}(\frac{\Delta}{t})^2$. Here, we show a similar bound using more elementary techniques. Suppose we have a set of k triangles such that no two share an edge. For each such triangle we define a random variable X_i which is 1 if the triangle is kept by the sampling and 0 otherwise. Then as the triangles do not have any edges in common, the X_i s are independent and take value 0 with probability $1 - p^3$ and 1 with probability p^3 . So by Chernoff inequality 2.2

$$\Pr \left[\left| \frac{1}{k} \sum_{i=1}^k X_i - p^3 \right| > \epsilon p^3 \right] \leq 2e^{-\epsilon^2 p^3 k/2}.$$

So when $p^3 k \epsilon^2 \geq 4d \log n$ where d is a positive constant, the probability of sparsification returning an ϵ -approximation is at least $1 - n^{-d}$. This is equivalent to $p^3 k \geq (4d \log n)/(\epsilon^2)$ which suggests that in order to sample with small p and hence discard many edges we need like k to be large. To show that such a large set of independent triangles exist, we invoke the Hajnal-Szemerédi Theorem 2.6 on an auxiliary graph H which we construct as follows. For each triangle i ($i = 1, \dots, t$) in G we create a vertex v_i in H . We connect two vertices v_i, v_j in H if and only if they represent triangles i, j respectively which share an edge in G . Notice that the maximum degree in the auxiliary graph H is $O(\Delta)$. Hence, we obtain the following Corollary.

Corollary 5.2. *Given t triangles such that no edge belongs to more than Δ triangles, we can partition the triangles into sets $S_1 \dots S_l$ such that $|S_i| > \Omega(t/\Delta)$ and l is bounded by $O(\Delta)$.*

Combining Corollary 5.2 and the Chernoff bound allows us to prove the next theorem.

²We use the tilde notation to hide polylogarithmic factors $\text{polylog}(n)$.

Theorem 5.3. *If $p^3 \in \Omega(\frac{\Delta \log n}{\epsilon^2 t})$, then with probability $1 - n^{-2}$, the sampled graph has a triangle count that ϵ -approximates t .*

Proof. Consider the partition of triangles given by corollary 5.2 and let $d = 5$. By choice of p we get that the probability that the triangle count in each set is preserved within a factor of $\epsilon/2$ is at least $1 - n^{-d}$. Since there are at most n^3 such sets, an application of the union bounds gives that their total is approximated within a factor of $\epsilon/2$ with probability at least $1 - n^{3-d}$. This gives that the triangle count is approximated within a factor of ϵ with probability at least $1 - n^{3-d}$. Substituting $d = 5$ completes the proof. \square

5.3.1.2 Triple Sampling

Since each triangle corresponds to a triple of vertices, we can construct a set of triples U that include all triangles. From this list, we can then sample triples uniformly at random. Let these samples be numbered from 1 to s . Also, for the i^{th} triple sampled, let X_i be 1 if it is a triangle and 0 otherwise. Since we pick triples randomly from U and t of them are triangles, we have $E(X_i) = \frac{t}{|U|}$ and X_i s are independent. So by Chernoff bound we obtain:

$$\Pr \left[\left| \frac{1}{s} \sum_{i=1}^s X_i - \frac{t}{|U|} \right| > \epsilon \frac{t}{|U|} \right] \leq 2e^{-\epsilon^2 ts / (2|U|)}$$

If $s = \Omega(\frac{|U| \log n}{t \epsilon^2})$, then we have that $|U| \sum_{i=1}^s \frac{X_i}{s}$ approximates t within a factor of ϵ with probability at least $1 - n^{-d}$ for any d of our choice. As $|U| \leq n^3$, this immediately gives an algorithm with runtime $O(n^3 \log n / (t \epsilon^2))$ that approximates t within a factor of ϵ . Slightly more careful bookkeeping can also give tighter bounds on $|U|$ in sparse graphs.

A simple but crucial observation which allows us to decide whether we will sample a triple of vertices or an edge and a vertex is the following. Consider any triple containing vertex u , (u, v, w) . Since $uv, uw \in E$, we have the number of such triples involving u is at most $d(u)^2$. From an edge-vertex sampling point of view, as $vw \in E$, another bound on the number of such triples is m . When $d(u) > m^{1/2}$, the second bound is tighter, and the first is in the other case.

These two cases naturally suggest that low degree vertices with degree at most $m^{1/2}$ be treated separately from high degree vertices with degree greater than $m^{1/2}$. For the number of triangles around low degree vertices, the value of $\sum_u d(u)^2$ is maximized when all edges are concentrated in as few vertices as possible [18]. Since the maximum

degree of such a vertex is $m^{1/2}$, the number of such triangles is upper bounded by $m^{1/2} \cdot (m^{1/2})^2 = m^{3/2}$. Also, as the sum of all degrees is $2m$, there can be at most $2m^{1/2}$ high degree vertices, which means the total number of triangles incident to these high degree vertices is at most $2m^{1/2} \cdot m = 2m^{3/2}$. Combining these bounds give that $|U|$ can be upper bounded by $3m^{3/2}$. Note that this bound is asymptotically tight when G is a complete graph ($n = m^{1/2}$). However, in practice the second bound can be further reduced by summing over the degree of all v adjacent to u , becoming $\sum_{uv \in E} d(v)$. As a result, an algorithm that implicitly constructs U by picking the better one among these two cases by examining the degrees of all neighbors will achieve $|U| \leq O(m^{3/2})$.

This improved bound on U gives an algorithm that ϵ approximates the number of triangles in time:

$$O\left(m + \frac{m^{3/2} \log n}{t\epsilon^2}\right)$$

As our experimental data in Section 4.1 indicate, the value of t is usually $\Omega(m)$ in practice. In such cases, the second term in the above calculation becomes negligible compared to the first one. In fact, in most of our data, just sampling the first type of triples (aka. pretending all vertices are of low degree) brings the second term below the first.

5.3.1.3 Hybrid algorithm

Edge sparsification with a probability of p allows us to only work on $O(mp)$ edges, therefore the total runtime of the triple sampling algorithm after sparsification with probability p becomes:

$$O\left(mp + \frac{\log n (mp)^{3/2}}{\epsilon^2 t p^3}\right) = O\left(mp + \frac{\log n m^{3/2}}{\epsilon^2 t p^{3/2}}\right).$$

As stated above, since the first term in most practical cases are much larger, we can set the value of p to balance these two terms out:

$$pm = \frac{m^{3/2} \log n}{p^{3/2} t \epsilon^2} \Rightarrow p^{5/2} t \epsilon^2 = m^{1/2} \log n \Rightarrow p = \left(\frac{m^{1/2} \log n}{t \epsilon^2}\right)^{2/5}$$

The actual value of p picked would also depend heavily on constants in front of both terms, as sampling is likely much less expensive due to factors such as cache effect and memory efficiency. Nevertheless, our experimental results in section 4 does seem to indicate that this type of hybrid algorithms can perform better in certain situations.

5.3.1.4 Sampling in the Semi-Streaming Model

The previous analysis of triangle counting by Alon, Yuster and Zwick was done in the streaming model [25], where the assumption was constant available space. We show that our sampling algorithm can be done in a slightly weaker model with space usage equaling:

$$O\left(m^{1/2} \log n + \frac{m^{3/2} \log n}{t\epsilon^2}\right)$$

We assume the edges adjacent to each vertex are given in order [164]. We first need to identify high degree vertices, specifically the ones with degree higher than $m^{1/2}$. This can be done by sampling $O(m^{1/2} \log n)$ edges and recording the vertices that are endpoints of one of those edges.

Lemma 5.4. *Suppose $dm^{1/2} \log n$ samples were taken, then the probability of all vertices with degree at least $m^{1/2}$ being chosen is at least $1 - n^{-d+1}$.*

Proof. Consider some vertex v with degree at least $m^{1/2}$. The probability of it being picked in each iteration is at least $m^{1/2}/m = m^{-1/2}$. As a result, the probability of it not picked in $dm^{1/2} \log n$ iterations is:

$$(1 - m^{-1/2})^{dm^{1/2} \log n} = \left[(1 - m^{1/2})^{m^{1/2}}\right]^{d \log n} \leq \left(\frac{1}{e}\right)^{d \log n} = n^{-d}$$

As there are at most n vertices, applying union bound gives that all vertices with degree at least $m^{1/2}$ are sampled with probability at least $1 - n^{-d+1}$. \square

Our proposed method is comprised of the following three steps/passes over the stream.

1. Identifying high degree vertices requires one pass of the graph. Also, note that the number of potential candidates can be reduced to $m^{1/2}$ using another pass over the edge list.
2. For all the low degree vertices, we can read their $O(m^{1/2})$ neighbors and sample from them. For the high degree vertices, we do the following: for each edge,

obtain a random variable y from a binomial distribution equal to the number of edge/vertices pairs that this edge is involved in. Then pick y vertices from the list of high degree vertices randomly. These two sampling procedures can be done together in another pass over the data.

3. Finally, we need to check whether each edge in the sampled triples belong to the edge list. We can store all such queries into a hash table as there are at most $O(\frac{m^{3/2} \log n}{t\epsilon^2})$ edges sampled w.h.p. Then going through the graph edges in a single pass and looking them up in table yields the desired answer.

5.3.2 Experiments

5.3.2.1 Data

The graphs used in our experiments are shown in Table 5.5. Multiple edges and self loops were removed (if any). All graphs with the exceptions of Livejournal-links and Flickr are available on the Web. Table 5.1 summarizes the resources.

Name	Nodes	Edges	Triangle Count	Description
AS-Skitter	1,696,415	11,095,298	28,769,868	Autonomous Systems
Flickr	1,861,232	15,555,040	548,658,705	Person to Person
Livejournal-links	5,284,457	48,709,772	310,876,909	Person to Person
Orkut-links	3,072,626	116,586,585	621,963,073	Person to Person
Soc-LiveJournal	4,847,571	42,851,237	285,730,264	Person to Person
Web-EDU	9,845,725	46,236,104	254,718,147	Web Graph (page to page)
Web-Google	875,713	3,852,985	11,385,529	Web Graph
Wikipedia 2005/11	1,634,989	18,540,589	44,667,095	Web Graph (page to page)
Wikipedia 2006/9	2,983,494	35,048,115	84,018,183	Web Graph (page to page)
Wikipedia 2006/11	3,148,440	37,043,456	88,823,817	Web Graph (page to page)
Wikipedia 2007/2	3,566,907	42,375,911	102,434,918	Web Graph (page to page)
Youtube	1,157,822	2,990,442	4,945,382	Person to Person

TABLE 5.5: Datasets used in our experiments.

5.3.2.2 Experimental Setup and Implementation Details

The experiments were performed on a single machine, with Intel Xeon CPU at 2.83 GHz, 6144KB cache size and 50GB of main memory. The graphs are from real world web-graphs, some details regarding them are in Table 5.1 and in Table 5.5. The algorithm was implemented in C++, and compiled using gcc version 4.1.2 and the -O3 optimization flag. Time was measured by taking the user time given by the linux time command. IO times are included in that time since the amount of memory operations

performed in setting up the graph is non-negligible. However, we use a modified IO routine that's much faster than the standard C/C++ scanf.

A major optimization that we used was to sort the edges in the graph and store the input file in the format as a sequence of neighbor lists per vertex. Each neighbor list begins with the size of the list, followed by the neighbors. This is similar to how software like Matlab stores sparse matrices. The preprocessing time to change the data into this format is not included. It can significantly improve the cache property of the graph stored, and hence the overall performance.

Some implementation details are based on this graph storage format. Specifically, since each triple that we check by definition has 2 edges already in the graph, it suffices to check/query whether the 3rd edge is present in the graph. In order to do this efficiently, rather than querying the existence of an edge upon sampling each triple, we store the entire set of the queries and answer them in one pass through the graph.. Finally, in the next section we discuss certain details behind efficient binomial sampling. Specifically picking a random subset of expected size $p|S|$ from a set S can be done in expected sublinear time, as we already saw in Claim 3.

5.3.2.3 Binomial Sampling in Expected Sublinear time

Most of our algorithms have the following routine in their core: given a list of values, keep each of them with probability p and discard with probability $1 - p$. If the list has length n , this can clearly be done using n random variables. As generating random variables can be expensive, it's preferable to use $O(np)$ random variables in expectation if possible. One possibility is to pick $O(np)$ random elements, but this would likely involve random accesses in the list, or maintaining a list of the indices picked in sorted order. A simple way that we use in our code to perform this sampling is to generate the differences between indices of entries retained [256]. This variable clearly belongs to an exponential distribution, and if x is a uniform random number in $(0, 1)$, taking $\lceil \log_{(1-p)} x \rceil$ as the value of the random variable, see [256]. The primary advantage of doing so is that sampling can be done while accessing the data in a sequential fashion, which results in much better cache performances.

5.3.2.4 Results

The six variants of the code involved in the experiment are first separated by whether the graph was first sparsified by keeping each edge with probability $p = 0.1$. In either case, an exact algorithm based on hybrid sampling with performance bounded by $O(m^{3/2})$

was run. Then two triple based sampling algorithms are also considered. They differ in whether an attempt to distinguish between low and high degree vertices, so the simple version is essentially sampling all 'V' shaped triples off each vertex. Note that no sparsification and exact also generates the exact number of triangles. Errors are measured by the absolute value of the difference between the value produced and the exact number of triangles divided by the exact number. The results on error and running time are averaged over five runs. The results are shown in Tables 5.6, 5.7.

Graph	No Sparsification					
	Exact		Simple		Hybrid	
	err(%)	time	err(%)	time	err(%)	time
AS-Skitter	0.000	4.452	1.308	0.746	0.128	1.204
Flickr	0.000	41.981	0.166	1.049	0.128	2.016
Livejournal-links	0.000	50.828	0.309	2.998	0.116	9.375
Orkut-links	0.000	202.012	0.564	6.208	0.286	21.328
Soc-LiveJournal	0.000	38.271	0.285	2.619	0.108	7.451
Web-EDU	0.000	8.502	0.157	2.631	0.047	3.300
Web-Google	0.000	1.599	0.286	0.379	0.045	0.740
Wiki-2005	0.000	32.472	0.976	1.197	0.318	3.613
Wiki-2006/9	0.000	86.623	0.886	2.250	0.361	7.483
Wiki-2006/11	0.000	96.114	1.915	2.362	0.530	7.972
Wiki-2007	0.000	122.395	0.943	2.728	0.178	9.268
Youtube	0.000	1.347	1.114	0.333	0.127	0.500

TABLE 5.6: Results of experiments averaged over 5 Trials using only triple sampling.

5.3.2.5 Remarks

From Table 5.5 it is evident that social networks are abundant in triangles. For example, the Flickr graph with only ~ 1.9 M vertices has ~ 550 M triangles and the Orkut graph with ~ 3 M vertices has ~ 620 M triangles. Furthermore, from Table 5.6 and Table 5.7 it is clear that none of the variants clearly outperforms the others on all the data. The gain/loss from sparsification is likely due to the fixed sampling rate. Adapting a doubling procedure for the sampling rate as in Section 5.2.2.4 is likely to mitigate this discrepancy. The difference between simple and hybrid sampling are due to the fact that handling the second case of triples has a much worse cache access pattern as it examines vertices that are two hops away. There are alternative implementations of how to handle this situation, which would be interesting for future implementations. A fixed sparsification rate of $p = 10\%$ was used mostly to simplify the setups of the experiments. In practice varying p to look for a rate where the results stabilize is the preferred option.

When compared with previous results on this problem, the error rates and running times of our results are all significantly lower. In fact, on the wiki graphs our exact counting algorithms have about the same order of speed with other approximate triangle counting implementations. This is also why we did not include any competitors in the exposition of the results since our implementation is a highly optimized C/C++ implementation with an emphasis on performance for huge graphs.

Graph	Sparsified ($p = 0.1$)					
	Exact		Simple		Hybrid	
	err(%)	time	err(%)	time	err(%)	time
AS-Skitter	2.188	0.641	3.208	0.651	1.388	0.877
Flickr	0.530	1.389	0.746	0.860	0.818	1.033
Livejournal-links	0.242	3.900	0.628	2.518	1.011	3.475
Orkut-links	0.172	9.881	1.980	5.322	0.761	7.227
Soc-LiveJournal	0.681	3.493	0.830	2.222	0.462	2.962
Web-EDU	0.571	2.864	0.771	2.354	0.383	2.732
Web-Google	1.112	0.251	1.262	0.371	0.264	0.265
Wiki-2005	1.249	1.529	7.498	1.025	0.695	1.313
Wiki-2006/9	0.402	3.431	6.209	1.843	2.091	2.598
Wiki-2006/11	0.634	3.578	4.050	1.947	0.950	2.778
Wiki-2007	0.819	4.407	3.099	2.224	1.448	3.196
Youtube	1.358	0.210	5.511	0.302	1.836	0.268

TABLE 5.7: Results of experiments averaged over 5 trials using sparsification and triple sampling.

As we mentioned earlier in Section 5.1 there exists a lot of interest into signed networks. It is clear that our method applies to this setting as well, by considering individually each possible configuration of a signed triangle. However, we do not include any of our experimental findings here due to the small size of the signed networks available to us via the Stanford Network Analysis library (SNAP).

5.3.3 Theoretical Ramifications

In Section 5.3.3.1 we discuss random projections and triangles, motivated by the simple observation that the inner product of two rows of the adjacency matrix corresponding to two connected vertices forming edge e gives the count of triangles $\Delta(e)$.

5.3.3.1 Random Projections and Triangles

Consider any two vertices $i, j \in V$ which are connected, i.e., $(i, j) \in E$. Observe that the inner product of the i -th and j -th column of the adjacency matrix of graph G gives the number of triangles that edge (i, j) participates in. Viewing the adjacency matrix as a collection of n points in \mathbb{R}^n , a natural question to ask is whether we can use results from the theory of random projections [227] to reduce the dimensionality of the points while preserving the inner products which contribute to the count of triangles. Magen and Zouzias [293] have considered a similar problem, namely random projections which preserve approximately the volume for all subsets of at most k points.

According to Lemma 2.5 projection $x \rightarrow Rx$ from $\mathbb{R}^d \rightarrow \mathbb{R}^k$ approximately preserves all Euclidean distances. However it does not preserve all pairwise inner products. This can easily be seen by considering the set of points $e_1, \dots, e_n \in \mathbb{R}^n = \mathbb{R}^d$ where $e_1 = (1, 0, \dots, 0)$ etc. Indeed, all inner products of the above set are zero, which cannot happen for the points Re_j as they belong to a lower dimensional space and they cannot all be orthogonal. For the triangle counting problem we do not need to approximate all inner products. Suppose $A \in \{0, 1\}^n$ is the adjacency matrix of a simple undirected graph G with vertex set $V(G) = \{1, 2, \dots, n\}$ and write A_i for the i -th column of A . The quantity we are interested in is the number of triangles in G (actually six times the number of triangles)

$$t = \sum_{u,v,w \in V(G)} A_{uv}A_{vw}A_{wu}.$$

If we apply a random projection of the above kind to the columns of A

$$A_i \rightarrow RA_i$$

and write

$$X = \sum_{u,v,w \in V(G)} (RA)_{uw}(RA)_{vw}(RA)_{wu}$$

it is easy to see that $\mathbb{E}[X] = 0$ since X is a linear combination of triple products $R_{ij}R_{kl}R_{rs}$ of entries of the random matrix R and that all such products have expected value 0, no matter what the indices. So we cannot expect this kind of random projection to work.

Therefore we consider the following approach which still has limitations as we will show in the following. Let

$$t = \sum_{u \sim v} A_u^\top A_v, \quad \text{where } u \sim v \text{ means } A_{uv} = 1,$$

and look at the quantity

$$\begin{aligned}
Y &= \sum_{u \sim v} (RA_u)^\top (RA_v) \\
&= \sum_{l=1}^k \sum_{i,j=1}^n \left(\sum_{u \sim v} A_{iu} A_{jv} \right) R_{li} R_{lj} \\
&= \sum_{l=1}^k \sum_{i,j=1}^n \#\{i - * - * - j\} R_{li} R_{lj}.
\end{aligned}$$

This is a quadratic form in the gaussian $N(0, 1)$ variables R_{ij} . By simple calculation for the mean value and diagonalization for the variance we see that if the X_j are independent $N(0, 1)$ variables and

$$Z = X^\top B X,$$

where $X = (X_1, \dots, X_n)^\top$ and $B \in \mathbb{R}^{n \times n}$ is *symmetric*, that

$$\begin{aligned}
\mathbb{E}[Z] &= \text{Tr } B \\
\text{Var}[Z] &= \text{Tr } B^2 = \sum_{i,j=1}^n (B_{ij})^2.
\end{aligned}$$

Hence $\mathbb{E}[Y] = \sum_{l=1}^k \sum_{i=1}^n \#\{i - * - * - i\} = k \cdot t$ so the mean value is the quantity we want (multiplied by k). For this to be useful we should have some concentration for Y near $\mathbb{E}[Y]$. We do not need exponential tails because we have only one quantity to control. In particular, a statement of the following type

$$\Pr[|Y - \mathbb{E}[Y]| > \epsilon \mathbb{E}[Y]] < 1 - c_\epsilon,$$

where $c_\epsilon > 0$ would be enough. The simplest way to check this is by computing the standard deviation of Y . By Chebyshev's inequality it suffices that the standard deviation be much smaller than $\mathbb{E}[Y]$. According to the formula above for the variance of a quadratic form we get

$$\begin{aligned}
\text{Var}[Y] &= \sum_{l=1}^k \sum_{i,j=1}^n \#\{i - * - * - i\}^2 \\
&= C \cdot k \cdot \#\{x - * - * - * - * - * - x\} \\
&= C \cdot k \cdot (\text{number of circuits of length 6 in } G).
\end{aligned}$$

Algorithm 2 Colorful Triangle Sampling**Input:** Unweighted graph $G([n], E)$ **Input:** Number of colors $N = 1/p$ Let $f : V \rightarrow [N]$ have uniformly random values $E' \leftarrow \{\{u, v\} \in E \mid f(u) = f(v)\}$ $T \leftarrow$ number of triangles in the graph (V, E') **return** T/p^2

Therefore, to have concentration it is sufficient that

$$\text{Var}[Y] = o(k \cdot (\mathbb{E}[Y])^2). \quad (5.7)$$

Observe that (5.7) is a sufficient -and not necessary- condition. Furthermore, (5.7) is certainly not always true as there are graphs with many 6-circuits and no triangles at all (the circuits *may* repeat vertices or edges).

5.4 Colorful Triangle Counting

In this Section we present a new sampling approach to approximating the number of triangles in a graph $G(V, E)$, that significantly improves existing sampling approaches. Furthermore, it is easily implemented in parallel. The key idea of our algorithm is to correlate the sampling of edges such that if two edges of a triangle are sampled, the third edge is always sampled. This decreases the degree of the multivariate polynomial that expresses the number of sampled triangles. This Section is organized as follows: in Section 5.4.1 we present our randomized algorithm. In Section 5.4.2 we present our main theoretical results, we analyze our algorithm and we discuss some of its important properties. In Section 5.4.3 we present an implementation of our algorithm in the popular MAPREDUCE framework.

5.4.1 Algorithm

Our algorithm, summarized as Algorithm 2, samples each edge with probability p , where $N = 1/p$ is integer, as follows. Let $f : [n] \rightarrow [N]$ be a random coloring of the vertices of $G([n], E)$, such that for all $v \in [n]$ and $i \in [N]$, $\Pr[f(v) = i] = p$. We call an edge *monochromatic* if both its endpoints have the same color. Our algorithm samples exactly the set E' of monochromatic edges, counts the number T of triangles in $([n], E')$ (using any exact or approximate triangle counting algorithm), and multiplies this count by p^{-2} .

Work presented in Sections 5.2, 5.3 has used a related sampling idea, the difference being that edges were sampled *independently* with probability p . Some intuition why this sampling procedure is less efficient than what we propose can be obtained by considering the case where a graph has t edge-disjoint triangles. With independent edge sampling there will be no triangles left (with probability $1 - o(1)$) if $p^3 t = o(1)$. Using our colorful sampling idea there will be $\omega(1)$ triangles in the sample with probability $1 - o(1)$ as long as $p^2 t = \omega(1)$. This means that we can choose a smaller sample, and still get accurate estimates from it.

5.4.2 Analysis

We wish to pick p as small as possible but at the same time have a strong concentration of the estimate around its expected value. How small can p be? In Section 5.4.2.1 we present a second moment argument which gives a sufficient condition for picking p . Our main theoretical result, stated as Theorem 5.5 in Section 5.4.2.2, provides a sufficient condition to this question. In Section 5.4.2.3 we analyze the complexity of our method. Finally, in Section 5.4.2.4 we discuss several aspects of our work.

5.4.2.1 Second Moment Method

Using the second moment method we are able to obtain the following strong theoretical guarantee:

Theorem 5.5. *Let n, t, Δ, T denote the number of vertices in G , the number of triangles in G , the maximum number of triangles an edge of G is contained in and the number of monochromatic triangles in the randomly colored graph respectively. Also let $N = \frac{1}{p}$ the number of colors used. If $p \geq \max\left(\frac{\Delta \log n}{t}, \sqrt{\frac{\log n}{t}}\right)$, then $T \sim \mathbb{E}[T]$ with probability $1 - \frac{1}{\log n}$.*

Proof. By Chebyshev's inequality 2.1, if $\text{Var}[T] = o(\mathbb{E}[T]^2)$ then $T \sim \mathbb{E}[T]$ with probability $1 - o(1)$ [26]. Let X_i be a random variable for the i -th triangle, $i = 1, \dots, t$, such that $X_i = 1$ if the i -th triangle is monochromatic. The number of monochromatic triangles T is equal to the sum of these indicator variables, i.e., $T = \sum_{i=1}^t X_i$. By the linearity of expectation and by the fact that $\Pr[X_i = 1] = p^2$ we obtain that $\mathbb{E}[T] = p^2 t$. It is easy to check that the only case where two indicator variables are dependent is when they share an edge. In this case the covariance is non-zero and for any $p > 0$, $\text{Cov}[X_i, X_j] = p^3 - p^4 < p^3$. We write $i \sim j$ if and only if X_i, X_j are dependent.

We obtain the following upper bound on the variance of T , where δ_e is the number of triangles edge e is contained in and $\Delta = \max_{e \in E(G)} \delta_e$:

$$\text{Var} [T] \leq \mathbb{E} [T] + \sum_{i \sim j} \text{Cov} [X_i \wedge X_j] < p^2 t + p^3 \sum_e \delta_e^2 \leq p^2 t + 3p^3 t \Delta$$

We pick p large enough to obtain $\text{Var} [T] = o(\mathbb{E} [T]^2)$. It suffices:

$$p^4 t^2 \gg p^2 t + 3p^3 t \Delta \Rightarrow p^2 t \gg 1 + 3p \Delta \quad (5.8)$$

We consider two cases, determined by which of the two terms of the right hand side is larger:

• CASE 1 ($p\Delta < 1/3$):

Since the right hand side of Inequality (5.8) is constant, it suffices that $p^2 t = \omega(n)$ where $\omega(n)$ is some slowly growing function. We pick $\omega(n) = \log n$ and hence $p \geq \sqrt{\frac{\log n}{t}}$.

• CASE 2 ($p\Delta \geq 1/3$):

In this case the right hand side of Inequality (5.8) is $\Theta(p\Delta)$ and therefore it suffices to pick $\frac{p^2 t}{\Delta} = \log n$.

Combining the above two cases we get that if

$$p \geq \max\left(\frac{\Delta \log n}{t}, \sqrt{\frac{\log n}{t}}\right)$$

inequality (5.8) is satisfied and hence by Chebyshev's inequality $T \sim \mathbb{E} [T]$ with probability $1 - \frac{1}{\log n}$.

□

Extremal Cases and Tightness of Theorem 5.5

Given the assumptions of Theorem 5.5, is the condition on p tight? The answer is affirmative as shown in Figure 5.2. Specifically, in Figure 5.2(a) G consists of t/Δ “books” of triangles, each of size Δ . This shows that p has to be at least $\omega(n) \frac{\Delta}{t}$ to hope for concentration, where $\omega(n)$ is some growing function of n . Similarly, when G consists of t disjoint triangles as shown in Figure 5.2(b) p has to be at least $\omega(n)t^{-1/2}$.

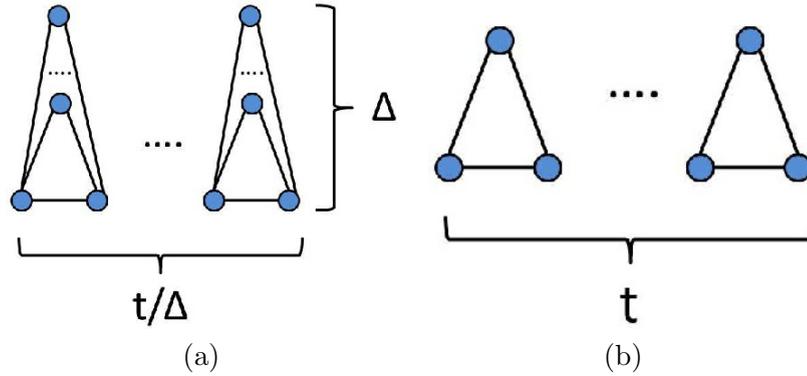


FIGURE 5.2: Conditions of Theorem 5.5 are tight. In order to hope for concentration p has to be greater than (a) $\frac{\Delta}{t}$ and (b) $t^{-1/2}$.

Therefore, unless we know more about G , we cannot hope for milder conditions on p , i.e., Theorem 5.5 provides an optimal condition on p .

5.4.2.2 Concentration via the Hajnal-Szemerédi Theorem

Here, we present a different approach to obtaining concentration, based on partitioning the set of triangles/indicator variables in sets containing many independent random indicator variables and then taking a union bound. Our theoretical result is the following theorem:

Theorem 5.6. *Let t_{\max} be the maximum number of triangles a vertex v is contained in. Also, let n, t, p, T be defined as above and ϵ a small positive constant. If $p^2 \geq \frac{16t_{\max} \log n}{\epsilon^2 t}$, then $\Pr [|T - \mathbb{E}[T]| > \epsilon \mathbb{E}[T]] \leq n^{-1}$.*

Proof. Let X_i be defined as above, $i = 1, \dots, t$. Construct an auxiliary graph H as follows: add a vertex in H for every triangle in G and connect two vertices representing triangles t_1 and t_2 if and only if they have a common vertex. The maximum degree of H is $3t_{\max} = O(\delta^2)$, where $\delta = O(n)$ is the maximum degree in the graph. Invoke the Hajnal-Szemerédi Theorem on H : we can partition the vertices of H (triangles of G) into sets S_1, \dots, S_q such that $|S_i| > \Omega(\frac{t}{t_{\max}})$ and $q = \Theta(t_{\max})$. Let $k = \frac{t}{t_{\max}}$. Note that the set of indicator variables X_i corresponding to any set S_j is independent. Applying the Chernoff bound for each set $S_i, i = 1, \dots, q$ we obtain

$$\Pr \left[\left| \frac{1}{k} \sum_{i=1}^k X_i - p^2 \right| > \epsilon p^2 \right] \leq 2e^{-\epsilon^2 p^2 k/2}$$

If $p^2 k \epsilon^2 \geq 4d' \log n$, then $2e^{-\epsilon^2 p^2 k/2}$ is upper bounded by $n^{-d'}$, where $d' > 0$ is a constant. Since $q = O(n^3)$ by taking a union bound over all sets S_i we see that the triangle count is approximated within a factor of ϵ with probability at least $1 - n^{3-d'}$. Setting $d' = 4$ completes the proof. \square

It's worth noting that for any constant $K > 0$ the above proof gives that if $p^2 \geq \frac{4(K+3)t_{\max} \log n}{\epsilon^2 t}$ then $\Pr[|T - \mathbb{E}[T]| > \epsilon \mathbb{E}[T]] \leq n^{-K}$.

5.4.2.3 Complexity

The running time of our procedure of course depends on the subroutine we use on the second step, i.e., to count triangles in the edge set E' . Let $d(i)$ denote the degree of vertex i . Assuming we use node iterator, i.e., the exact method that examines each vertex independently and counts the number of edges among its neighbors, our algorithm runs in $O(n + m + p^2 \sum_{i \in [n]} d^2(i))$ expected time³ by efficiently storing the graph and retrieving the neighbors of v colored with the same color as v in $O(1 + p d(v))$ expected time. Note that this implies that the speedup with respect to the counting task is $1/p^2$.

5.4.2.4 Discussion

Despite the fact that the second moment argument gave us strong conditions on p , the use of the Hajnal-Szemerédi theorem, see Theorem 2.6 and [208], has the potential of improving the Δ factor. The condition we provide on p is *sufficient* to obtain concentration. Note –see Figure 5.3– that it was necessary to partition the triangles into vertex disjoint rather than edge disjoint triangles since we need mutually independent variables per chromatic class in order to apply the Chernoff bound. If we were able to remove the dependencies in the chromatic classes defined by edge disjoint triangles, then the overall result could probably be improved. It's worth noting that for $p = 1$ we obtain that $t \geq n\omega(n)$, where $\omega(n)$ is any slowly growing function of n . This is –to the best of our knowledge– the mildest condition on the triangle density needed for a randomized algorithm to obtain concentration. Finally, notice that when $t \leq \frac{\Delta^2 \log n}{t_{\max}}$ and $t_{\max} \geq 1$ Theorem 5.6 yields a better bound than Theorem 5.5. The same holds when $t > \Delta^2 \log n$ and $t_{\max} \leq 1$. The latter scenario is far more restrictive and both Theorem 5.5 and Theorem 5.6 give for instance the same bound $p \geq \sqrt{\frac{\log n}{t}}$ for the graph of Figure 5.2(b).

³We assume that uniform sampling of a color takes constant time. If not, then we obtain the term $O(n \log(\frac{1}{p}))$ for the vertex coloring procedure.

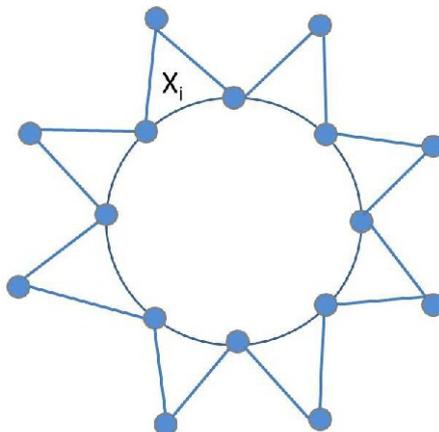


FIGURE 5.3: Consider the indicator variable X_i corresponding to the i -th triangle. Note that $\Pr[X_i | \text{rest are monochromatic}] = p \neq \Pr[X_i] = p^2$. The indicator variables are *pairwise* but not mutually independent.

Furthermore, the powerful theorem of Kim and Vu 2.4 that was used in Section 5.2 is not immediately applicable here: let Y_e be an indicator variable for each edge e such that $Y_e = 1$ if and only if e is monochromatic, i.e., both its endpoints receive the same color. Note that the number of triangles is a Boolean polynomial $T = \frac{1}{3} \sum_{\Delta(e,f,g)} (Y_e Y_f + Y_f Y_g + Y_e Y_g)$ but the Boolean variables are not independent as the Kim-Vu [249] theorem requires. It is worth noting that the degree of the polynomial is two. Essentially, this is the reason for which our method obtains better results than work in Section 5.2 where the degree of the multivariate polynomial is three. Finally, it is worth noting that using a simple doubling procedure as the one outlined in Section 5.2, we can pick p effectively in practice despite the fact that it depends on the quantity t which we want to estimate by introducing an extra logarithm in the running time.

Finally, from an experimentation point of view, it is interesting to see how well the upper bound $3\Delta t$ matches the sum $\sum_{e \in E(G)} \delta_e^2$, where δ_e is the number of triangles edge e is contained in, and the typical values for Δ and t_{\max} in real-world graphs. The following table shows these numbers for five graphs ⁴ taken from the SNAP library [5]. We see that Δ and t_{\max} are significantly less than their upperbounds and that typically $3\Delta t$ is significantly larger than $\sum_{e \in E(G)} \delta_e^2$ except for the collaboration network of Arxiv Astro Physics. The results are shown in Table 5.8.

5.4.3 A MapReduce Implementation

MAPREDUCE [130] has become the *de facto* standard in academia and industry for analyzing large scale networks. For a brief overview of MAPREDUCE see Section 1.1.6.1.

⁴AS:Autonomous Systems, Oregon: Oregon route views, Enron: Email communication network, ca-HepPh and AstroPh:Collaboration networks. Self-edges were removed.

Name	Vertices(n)	Edges(m)	Triangle Count(t)	Δ	t_{\max}	$\sum_{e \in E(G)} \delta_e^2$	$3\Delta t$
AS	7,716	12,572	6,584	344	2,047	595,632	6,794,688
Oregon	11,492	23,409	19,894	537	3,638	2,347,560	32,049,234
Enron	36,692	183,831	727,044	420	17,744	75,237,684	916,075,440
ca-HepPh	12,008	118,489	3,358,499	450	39,633	1.8839×10^9	4.534×10^9
AstroPh	18,772	198,050	1,351,441	350	11,269	148,765,753	1.419×10^9

TABLE 5.8: Values for the variables involved in our formulae for five real-world graphs. Typically, Δ and t_{\max} are significantly less than the obvious upper bounds $n - 2$ and $\binom{n-1}{2}$ respectively. Furthermore, $3\Delta t$ is significantly larger than $\sum_{e \in E(G)} \delta_e^2$.

Algorithm 3 MAPREDUCE Colorful Triangle Counting $G(V, E), p = 1/N$

Map: Input $\langle e = (u, f(u), v, f(v)); 1 \rangle$ {Let f be a uniformly at random coloring of the vertices with N colors}

if $f(u) = f(v)$ then emit $\langle f(u); (u, v) \rangle$

Reduce: Input $\langle c; E_c = \{(u, v)\} \subseteq E \rangle$ { Every edge $(u, v) \in E_c$ has color c , i.e., $f(u) = f(v)$ }

Scale each triangle by $\frac{1}{p^2}$.

Recent work by Suri and Vassilvitskii [377] proposes two algorithms for counting triangles. The first is an efficient MAPREDUCE implementation of the node iterator algorithm, see also [352] and the second is based on partitioning the graph into overlapping subsets so that each triangle is present in at least one of the subsets.

Our method is amenable to being implemented in MAPREDUCE and the skeleton of such an implementation is shown in Algorithm 2⁵. We implicitly assume that in a first round vertices have received a color uniformly at random from the N available colors and that we have the coloring information for the endpoints of each edge. Each mapper receives an edge together with the colors of its edgepoints. If the edge is monochromatic, then it's emitted with the color as the key and the edge as the value. Edges with the same color are shipped to the same reducer where locally a triangle counting algorithm is applied. The total count is scaled appropriately. Trivially, the following lemma holds by the linearity of expectation and the fact that the endpoints of any edge receive a *given* color c with probability p^2 .

Lemma 5.7. *The expected size to any reduce instance is $O(p^2m)$ and the expected total space used at the end of the map phase is $O(pm)$.*

⁵It's worth pointing out for completeness reasons that in practice one would not scale the triangles after the first reduce. It would emit the count of monochromatic triangles which would be summed up in a second round and scaled by $1/p^2$.

Chapter 6

Dense Subgraphs

6.1 Introduction

Given a graph $G = (V, E)$ and a subset of vertices $S \subseteq V$, let $G[S] = (S, E[S])$ be the subgraph induced by S , and let $e[S]$ be the size of $E[S]$. The *edge density* of the set S is defined as $\delta(S) = e[S]/\binom{|S|}{2}$. Finding a dense subgraph of G would in principle require to find a set of vertices $S \subseteq V$ that maximizes $\delta(S)$. However, the direct maximization of δ is not a meaningful problem, as even a single edge achieves maximum density. Therefore, effort has been devoted to define alternative density functions whose maximization allows for extracting subgraphs having large δ and, at the same time, non-trivial size. Different choices of the density function lead to different variants of the dense-subgraph problem. Some variants can be solved in polynomial time, while others are **NP**-hard, or even inapproximable.

In this Chapter we introduce a new framework for finding dense subgraphs. We focus on a special case of the framework which we refer to as **optimal quasi-cliques**. Table 6.1 offers a preview of the results that will follow. Specifically, it compares our **optimal quasi-cliques** with **densest subgraphs** on some popular graphs.¹ The results in the table clearly show that **optimal quasi-cliques** have much larger edge density than **densest subgraphs**, smaller diameters and larger triangle densities. Moreover, **densest subgraphs** are usually quite large-sized: in the graphs we report in Table 6.1, the **densest subgraphs** contain always more than the 30% of the vertices in the input graph. For instance, in the Football graph, the **densest subgraph** corresponds to the whole graph, with edge density < 0.1 and diameter 4, while the extracted **optimal quasi-clique** is a 12-vertex subgraph with edge density 0.73 and diameter 2. The Jazz graph contains a perfect clique of 30 vertices: our

¹Densest subgraphs are extracted here with the exact Goldberg's algorithm [199]. As far as optimal quasi-cliques, we optimize f_α with $\alpha = \frac{1}{3}$ and use our local-search method.

	densest subgraph				optimal quasi-clique			
	$\frac{ S }{ V }$	δ	D	τ	$\frac{ S }{ V }$	δ	D	τ
Dolphins	0.32	0.33	3	0.04	0.12	0.68	2	0.32
Football	1	0.09	4	0.03	0.10	0.73	2	0.34
Jazz	0.50	0.34	3	0.08	0.15	1	1	1
Celeg. N.	0.46	0.13	3	0.05	0.07	0.61	2	0.26

TABLE 6.1: Difference between densest subgraph and optimal quasi-clique on some popular graphs. $\delta = e[S]/\binom{|S|}{2}$ is the edge density of the extracted subgraph, D is the diameter, and $\tau = t[S]/\binom{|S|}{3}$ is the triangle density.

method finds this clique achieving perfect edge density, diameter, and triangle density scores. By contrast, the densest subgraph contains 100 vertices, and has edge density 0.34 and triangle density 0.08.

6.2 A General Framework

Let $G = (V, E)$ be a graph, with $|V| = n$ and $|E| = m$. For a set of vertices $S \subseteq V$, let $e[S]$ be the number of edges in the subgraph induced by S . We define the following function.

Definition 6.1 (Edge-surplus). Let $S \subseteq V$ be a subset of the vertices of a graph $G = (V, E)$, and let $\alpha > 0$ be a constant. We define *edge-surplus* as:

$$f_\alpha(S) = g(e[S]) - \alpha h(|S|),$$

where functions g, h are both strictly increasing. We also define $f_\alpha(\emptyset) = 0$.

We note that the first term $g(e[S])$ encourages subgraphs abundant in edges whereas the second term $-\alpha h(|S|)$ penalizes large subgraphs. Our framework for finding dense subgraphs is based on the following optimization problem.

Problem 1 (optimal (g, h, α) -edge-surplus). Given a graph $G = (V, E)$, a positive real α and a pair of functions g, h , find a subset of vertices $S^* \subseteq V$ such that $f_\alpha(S^*) \geq f_\alpha(S)$, for all sets $S \subseteq V$. We refer to the set S^* as the *optimal (g, h, α) -edge-surplus* of the graph G .

The edge-surplus definition subsumes numerous popular existing density measures by choosing appropriately g, h, α . Three important cases follow.

- By setting $g(x) = h(x) = \log x, \alpha = 1$, the optimal (g, h, α) -edge-surplus problem becomes equivalent to maximizing $\log e[S] - \log |S| = \log \frac{e[S]}{|S|}$. This is equivalent to the popular **densest subgraph**.

- By setting $g(x) = \log x, h(x) = \log \left(\frac{x(x-1)}{2} \right), \alpha = 1$ the optimal (g, h, α) -edge-surplus problem becomes equivalent to maximizing $\frac{e[S]}{\binom{|S|}{2}}$.
- By setting $g(x) = x, h(x) = \frac{x(x-1)}{2}$ and restricting $\alpha \in (0, 1)$ we obtain the following optimal (g, h, α) -edge-surplus problem: $\max_{\emptyset \neq S \subseteq V} e[S] - \alpha \binom{|S|}{2}$. We call this problem OQC-PROBLEM. We notice that it turns the quasi-clique condition into an objective. To the best of our knowledge, this optimization problem does not appear in the existing literature.

The densest subgraph problem on the one hand is polynomially time solvable but results typically in very large subgraphs. Also, maximizing the edge density $\delta(S)$ results in trivial subgraphs such as edges or triangles which achieve the maximum possible density value 1. We wish to better understand the OQC-PROBLEM problem. We start by discussing properties of the objective. **Understanding the objective:** Consider the case $\alpha = 0$. Clearly, the optimal solution is the whole graph. When $0 < \alpha < 1$ the problem in general is NP-hard. We discuss the case where $\alpha > 1/2$. It is straightforward to check that by setting $\alpha = 1 - \frac{1}{\Omega(n^2)}$, e.g., $\alpha = 1 - n^{-3}$, one solves the maximum clique problem. Furthermore, assuming that finding a hidden clique of order $O(n^{1/2-\delta})$ where $\delta > 0$ in a random binomial graph $G \sim G(n, 1/2)$ is hard, then one can see that our problem is hard. To see why, notice that for any set of S vertices the expected score is $\left(\frac{1}{2} - \alpha\right) \binom{n^{1/2-\delta}}{2}$ and for the hidden clique $(1 - \alpha) \binom{n^{1/2-\delta}}{2}$. Therefore, if we could optimally solve the OQC-PROBLEM problem, then by setting $\alpha > 1/2$, we could solve in expectation the hidden clique problem.

Theorem 6.2. *The OQC-PROBLEM is NP-hard.*

Proving that the problem is NP-hard for any $\alpha \in (0, 1)$ remains open. When $\alpha = 1$ all cliques receive score equal to 0, independent of their size. When $\alpha > 1$ then the problem stops being interesting as the optimal solution will be any single edge. An important property of the edge-surplus abstraction is that it allows us to model scenarios in numerous practical situations where we wish to find a dense subgraph with bounds on its size. For instance, by relaxing the monotonicity property of function $h(\cdot)$ the k -densest subgraph problem can be modeled as an optimal (g, h, α) -edge-surplus problem by setting $g(x) = x$ and

$$h(x) = \begin{cases} 0 & x = k \\ +\infty & \text{otherwise.} \end{cases}$$

By choosing $h(x)$ to penalize severely undesired size, a good algorithm will avoid outputting a subgraph of undesired size.

Finally, we outline that we cannot make any general statement on the optimal (g, h, α) -edge-surplus problem, since certain cases are polynomially time solvable whereas others **NP**-hard. The following theorem provides a family of optimal (g, h, α) -edge-surplus problems that are efficiently solvable

Theorem 6.3. *Let $g(x) = x$ and $h(x)$ be a concave function. Then the optimal (g, h, α) -edge-surplus problem is in **P**.*

Proof. The optimal (g, h, α) -edge-surplus problem becomes $\max_{\emptyset \neq S \subseteq V} e[S] - \alpha h(|S|)$ where $h(x)$ is a concave function. The claim follows directly from combining the following facts.

Fact 1 The function defined by the map $S \mapsto e[S]$ is a supermodular function.

Fact 2 The function $h(|S|)$ is submodular given that h is concave. Since $\alpha > 0$, the function $-\alpha h(|S|)$ is supermodular.

Fact 3 Combining the above facts with the fact that the sum of two supermodular functions is supermodular, we obtain $f_\alpha(S)$ is a supermodular function.

Fact 4 Maximizing supermodular functions is strongly polynomially time solvable [357].

□

However, we outline that the scenario where $g(x) = x$ and $h(x)$ is concave cannot be useful in real applications as the output subgraph will be large.

6.3 Optimal Quasi-cliques

Scalable Algorithms. The first efficient algorithm we propose is an adaptation of the greedy algorithm by Asashiro et al. [41], which has been shown to provide a $\frac{1}{2}$ -approximation for the densest subgraph problem [104]. The outline of our algorithm, called GREEDYOQC, is shown as Algorithm 4. The algorithm iteratively removes the vertex with the smallest degree. The output is the subgraph produced over all iterations that maximizes the objective function f_α . The algorithm can be implemented in $\mathcal{O}(n + m)$ time: the trick consists in keeping a list of vertices for each possible degree and updating the degree of any vertex v during the various iterations of the algorithm simply by moving v to the appropriate degree list.

The GREEDYOQC algorithm provides an additive approximation guarantee for the OQC-PROBLEM, as shown next.

Theorem 6.4. *Let \bar{S} be the set of vertices outputted by the GREEDYOQC algorithm and let S^* be the optimal vertex set. Consider also the specific iteration of the algorithm*

Algorithm 4 GREEDYOQC**Input:** Graph $G(V, E)$ **Output:** Subset of vertices $\bar{S} \subseteq V$ $S_n \leftarrow V$ **for** $i \leftarrow n$ **downto** 1 **do** Let v be the vertex with the smallest degree in $G[S_i]$ $S_{i-1} \leftarrow S_i \setminus \{v\}$ **end for** $\bar{S} \leftarrow \arg \max_{i=1, \dots, n} f_\alpha(S_i)$

where a vertex within S^* is removed for the first time and let S_I denote the vertex set currently kept in that iteration. It holds that:

$$f_\alpha(\bar{S}) \geq f_\alpha(S^*) - \frac{\alpha}{2} |S_I| (|S_I| - |S^*|).$$

Proof. Given a subset of vertices $S \subseteq V$ and a vertex $u \in S$, let $d_S(u)$ denote the degree of u in $G[S]$.

We start the analysis by considering the first vertex belonging to S^* removed by the algorithm from the current vertex set. Let v denote such a vertex, and let also S_I denote the set of vertices still present just before the removal of v . By the optimality of S^* , we obtain:

$$\begin{aligned} f_\alpha(S^*) &\geq f_\alpha(S^* \setminus \{u\}), \quad \forall u \in S^* \\ &\Leftrightarrow e[S^*] - \alpha \binom{|S^*|}{2} \geq (e[S_I] - d_{S^*}(u)) - \alpha \binom{|S^*| - 1}{2}, \quad \forall u \in S^* \\ &\Leftrightarrow d_{S^*}(u) \geq \alpha(|S^*| - 1), \quad \forall u \in S^*. \end{aligned}$$

As the algorithm greedily removes vertices with the smallest degree in each iteration, it is easy to see that $d_V(u) \geq d_{S_I}(u) \geq d_{S^*}(u) \geq \alpha(|S^*| - 1)$, $\forall u$. Therefore, noticing also

that $S^* \subseteq S_I$, it holds that:

$$\begin{aligned}
f_\alpha(S_I) &= e[S_I] - \alpha \binom{|S_I|}{2} + \alpha \binom{n}{2} \\
&= \frac{1}{2} \left(\sum_{u \in S^*} d_{S^*}(u) + \sum_{u \in S^*} (d_{S_I}(u) - d_{S^*}(u)) + \right. \\
&\quad \left. + \sum_{u \in S_I \setminus S^*} d_{S_I}(u) \right) - \alpha \binom{|S_I|}{2} + \alpha \binom{n}{2} \\
&\geq \frac{1}{2} \left(\sum_{u \in S^*} d_{S^*}(u) + \sum_{u \in S_I \setminus S^*} d_{S_I}(u) \right) - \alpha \binom{|S_I|}{2} + \alpha \binom{n}{2} \\
&= e[S^*] + \frac{1}{2} \sum_{u \in S_I \setminus S^*} d_{S_I}(u) - \alpha \binom{|S_I|}{2} + \alpha \binom{n}{2} \\
&\geq e[S^*] + \frac{1}{2} (|S_I| - |S^*|) \alpha (|S^*| - 1) - \alpha \binom{|S_I|}{2} + \alpha \binom{n}{2} \\
&= f_\alpha(S^*) - \frac{\alpha}{2} |S_I| (|S_I| - |S^*|).
\end{aligned}$$

As the final output of the algorithm is the best over all iterations, we finally obtain:

$$f_\alpha(\bar{S}) \geq f_\alpha(S_I) \geq f_\alpha(S^*) - \frac{\alpha}{2} |S_I| (|S_I| - |S^*|).$$

□

The above result can be interpreted as follows. Assuming that $|S_I|$ is $\mathcal{O}(|\bar{S}|)$, the additive approximation factor proved in Theorem 6.4 becomes $f_\alpha(\bar{S}) \geq f_\alpha(S^*) - \frac{\alpha}{2} |\bar{S}| (|\bar{S}| - |S^*|)$. Thus, the error achieved by the GREEDYOQC algorithm is guaranteed to be bounded by an additive factor proportional to the size of the optimal quasi-clique outputted. As optimal quasi-cliques are typically small graphs, this results in an approximation guarantee that is very tight in practice.

Finally, we present a local search heuristic for solving the OQC-PROBLEM. The algorithm, called LOCALSEARCHOQC, performs local operations and outputs a vertex set S that is guaranteed to be locally optimal, i.e., if any single vertex is added to or removed from S , then the objective function decreases.

The outline of LOCALSEARCHOQC is shown as Algorithm 5. The algorithm initially selects a random vertex and then it keeps adding vertices to the current set S while the objective improves. When no vertex can be added, the algorithm tries to find a vertex in S whose removal may improve the objective. As soon as such a vertex is encountered, it is removed from S and the algorithm re-starts from the adding phase. The process

Algorithm 5 LOCALSEARCHOQC**Input:** Graph $G = (V, E)$; maximum number of iterations T_{MAX} **Output:** Subset of vertices $\bar{S} \subseteq V$ $S \leftarrow \{v\}$, where v is chosen uniformly at random $b_1, b_2 \leftarrow \text{TRUE}$, $t \leftarrow 1$.**while** b_1 and $t \leq T_{MAX}$ **do****while** b_2 **do**If there exists $u \in V \setminus S$ such that $f_\alpha(S \cup \{u\}) \geq f_\alpha(S)$ then let $S \leftarrow S \cup \{u\}$ otherwise set $b_2 \leftarrow \text{FALSE}$ **end while**If there exists $u \in S$ such that $f_\alpha(S \setminus \{u\}) \geq f_\alpha(S)$ then let $S \leftarrow S \setminus \{u\}$ otherwise, set $b_1 \leftarrow \text{FALSE}$ $t \leftarrow t + 1$ **end while** $\bar{S} \leftarrow \arg \max_{\hat{S} \in \{S, V \setminus S\}} f_\alpha(\hat{S})$

continues until a local optimum is reached or the number of iterations exceeds T_{\max} . The time complexity of LOCALSEARCHOQC is $\mathcal{O}(T_{\max} m)$.

In order to enhance the performance of the LOCALSEARCHOQC algorithm, one may use the following heuristic [197]. Let v^* be the vertex that maximizes the ratio $\frac{t(v^*)}{d(v^*)}$, where $t(v^*)$ is the number of triangles of v^* and $d(v^*)$ its degree (we approximate the number of triangles in which each vertex participates with the technique described in [257]). Given vertex v^* , we use as a seed the set $\{v^* \cup N(v^*)\}$, where $N(v^*) = \{u : (u, v^*) \in E\}$ is the neighborhood of v^* .

Parameter Selection. Finally, a natural question that arises whenever a parameter exists, is how to choose an appropriate value. No doubt, there exist different possible, principled ways which lead to different choices of α . For instance, setting α to be the graph edge density $\delta(G)$ results into a normalized version of our criterion. However, since real-world networks are sparse, we do not encourage this for practical purposes. Instead, we provide a simple criterion to pick α .

Let us consider two disjoint sets of vertices S_1, S_2 in the graph G . Assume that $G[S_1 \cup S_2]$ is disconnected, i.e., $G[S_1]$ and $G[S_2]$ form two separate connected components. Also, without any loss of generality, assume that $f_\alpha(S_1) \leq f_\alpha(S_2)$. As our goal is to favor small dense subgraphs, a natural condition to satisfy is $f_\alpha(S_1 \cup S_2) \leq f_\alpha(S_1) \leq f_\alpha(S_2)$, i.e., we require for our objective to prefer the set S_1 (or S_2) rather than the larger set $S_1 \cup S_2$. Therefore, we obtain:

$$e[S_1] + e[S_2] - \alpha \binom{|S_1| + |S_2|}{2} + \alpha \binom{n}{2} \leq e[S_1] - \alpha \binom{|S_1|}{2} + \alpha \binom{n}{2},$$

which, considering that $e[S_2] \leq \binom{|S_2|}{2}$, leads to:

$$\alpha \geq \frac{\binom{|S_2|}{2}}{\binom{|S_1|+|S_2|}{2} - \binom{|S_1|}{2}} = \frac{|S_2| - 1}{2|S_1| + |S_2| - 1}.$$

Let us now assume for simplicity that $|S_1| = |S_2| = k$; then the above condition becomes: $\alpha \geq \frac{k-1}{3k-1}$. As $\frac{k-1}{3k-1} < \frac{1}{3}$, it suffices choosing $\alpha \geq \frac{1}{3}$ to have the condition satisfied.

On the other hand, it is easy to see that, when α is close to 1, f_α tends to be maximized even by subgraphs of trivial structure (e.g., single edges), which is clearly something that we want to avoid. By combining the two arguments above, we conclude that a good choice for α is a value around $\frac{1}{3}$, which is the value we adopt in our experiments.

Multiplicative Approximation, a 0.796-approximation algorithm for a shifted objective. We also design a multiplicative approximation algorithm for a shifted objective which works for any $\alpha > 0$. Notice that this shifting is not necessary since the optimal objective value is positive in the interesting range of $0 < \alpha < 1$ as a single edge results in a positive score $1 - \alpha$. Our algorithm is based on semidefinite programming and in particular on the techniques developed by Goemans-Williamson [198]. Our algorithm is a β -approximation algorithm, where $\beta > 0.796$, for a *shifted* objective. Specifically, we *shift* our objective by a constant $c = c(n)$ to make it non-negative.

Observation 1. Consider $f'_\alpha(S) = f_\alpha(S) + \alpha \binom{n}{2}$. Then $f'_\alpha(S) \geq 0$ for any $S \subseteq V$ since $\binom{n}{2} - \binom{s}{2} \geq 0$ for any $S \subseteq V$ and $e[S] \geq 0$. Furthermore we have $f'_\alpha(S_1) \geq f'_\alpha(S_2)$ if and only if $f_\alpha(S_1) \geq f_\alpha(S_2)$.

All the guarantees we obtain in this Section refer to the shifted objective f'_α . Therefore, from now on we will abuse slightly the notation and use f_α to denote f'_α . We formulate our maximization problem as an integer program. We introduce a variable $x_i \in \{-1, +1\}$ for each vertex $i \in V = \{1, \dots, n\}$ and an extra variable x_0 which expresses whether a vertex belongs to S or not:

$$\text{It is } i \in S \text{ if and only if } x_0 x_1 = 1.$$

Notice that the term $\frac{1+x_0x_i+x_0x_j+x_ix_j}{4}$ equals 1 if and only if both i, j belong in S , otherwise it equals 0. Furthermore, the term $\binom{n}{2}$ enters the objective as $\frac{1}{2} \sum_{i \neq j} 1$. Therefore, we get the following integer program:

$$\begin{aligned}
\max \quad & \sum_{e=(i,j)} \frac{1 + x_i x_j + x_i x_j + x_i x_j}{4} + \\
& \frac{\alpha}{2} \sum_{i \neq j} \left(1 - \frac{1 + x_i x_j + x_i x_j + x_i x_j}{4} \right) \\
\text{subject to } & x_i \in \{-1, +1\}, \text{ for all } i \in \{0, 1, \dots, n\}.
\end{aligned} \tag{6.1}$$

We relax the integrality constraint and we allow the variables to be vectors in the unit sphere in \mathbb{R}^{n+1} . By using the variable transformation $y_{ij} = x_i x_j$, we obtain the following semidefinite programming relaxation:

$$\begin{aligned}
\max \quad & \alpha \sum_{e=(i,j)} \frac{1 + y_{0i} + y_{0j} + y_{ij}}{4} + \\
& \frac{1}{2} \sum_{i \neq j} \left(1 - \frac{1 + y_{0i} + y_{0j} + y_{ij}}{4} \right) \\
\text{subject to } & y_{ii} = 1, \text{ for all } i \in \{0, 1, \dots, n\} \\
& \text{and } Y \succeq 0, Y \text{ symmetric.}
\end{aligned} \tag{6.2}$$

The above SDP can be solved within an additive error of δ of the optimum in polynomial time by interior point algorithms or the ellipsoid method [24]. In what follows, we refer to the optimal value of the integer program as IP^* and of the semidefinite program as SDP^* . Our algorithm, SDP-OQC , is shown as Algorithm 6.

Theorem 6.5. *Algorithm SDP-OQC is a β -approximation algorithm for f_α where $\beta > 0.796$ with probability at least $1 - \mathcal{O}(n^{-1})$.*

Proof. First, notice that we can rewrite the objective as

$$\begin{aligned}
& \sum_{e=(i,j)} \frac{1 + y_{0i} + y_{0j} + y_{ij}}{4} + \\
& \frac{\alpha}{4} \left(\sum_{i \neq j} \frac{1 - y_{0i}}{2} + \sum_{i \neq j} \frac{1 - y_{0j}}{2} + \sum_{i \neq j} \frac{1 - y_{ij}}{2} \right).
\end{aligned}$$

Algorithm 6 SDP-OQC**Input:** $G = (V, E)$ **1. Relaxation**

Solve the semidefinite program (6.2)

Compute a Cholesky decomposition of the resulting Y Let v_0, v_1, \dots, v_n be the resulting vectors**2. Randomized Rounding**Randomly choose a unit length vector $r \in \mathbb{R}^{n+1}$ Set $S = \{i \in [n] : \text{sgn}(v_i r) = \text{sgn}(v_0 r)\}$ **3. Boosting the success probability**Repeat steps 1-2 for $t = 1, \dots, T$ Output the best solution over $T = c_{\epsilon, \alpha, \beta} \log n$ runs% Here $\epsilon > 0$ and $c_{\epsilon, \alpha, \beta} \geq \frac{2(\alpha+1)}{3\epsilon\alpha\beta} + 1$.

Let $Y = [v_0 v_1 \dots v_n]^T [v_0 v_1 \dots v_n]$ be the Cholesky decomposition of matrix Y . We analyze the randomized rounding step which is equivalent to considering a random hyperplane \mathcal{H} that goes through the origin and placing in set S all the vertices whose corresponding vector for the Cholesky decomposition v_i is on the same side of \mathcal{H} with v_0 . Our goal now is to lower bound the expectation of our objective upon this randomized rounding. The expectation of the terms of the form $\frac{1-y_{ij}}{2}$ is equal to the probability that the two vectors v_i, v_j are on different sides of the random hyperplane. As in Goemans-Williamson [198] this probability is

$$\Pr[\text{sgn}(v_i r) \neq \text{sgn}(v_j r)] = \frac{\arccos(v_i v_j)}{\pi}.$$

Furthermore, again as in Goemans-Williamson [198], for any $0 \leq \theta \leq \pi$ we have

$$\frac{\theta}{\pi} \geq 0.87856 \frac{1 - \cos \theta}{2} > \beta \frac{1 - \cos \theta}{2}.$$

Now, we lower bound the expectation of the first term in our objective. To do so, we need to compute the probability that $\text{sgn}(v_i r) = \text{sgn}(v_j r) = \text{sgn}(v_0 r)$. Consider the following events:

$$A : \quad \text{sgn}(v_i r) = \text{sgn}(v_j r) = \text{sgn}(v_0 r)$$

$$B_i : \quad \text{sgn}(v_i r) \neq \text{sgn}(v_j r) = \text{sgn}(v_0 r)$$

$$B_j : \quad \text{sgn}(v_j r) \neq \text{sgn}(v_i r) = \text{sgn}(v_0 r)$$

$$B_0 : \quad \text{sgn}(v_0 r) \neq \text{sgn}(v_j r) = \text{sgn}(v_i r)$$

Notice that $\Pr[B_i] = \Pr[\text{sgn}(v_j r) = \text{sgn}(v_0 r)] - \Pr[A]$, and that similar equations hold for indices $(j, 0)$. Furthermore, by the pigeonhole principle $\Pr[A] + \Pr[B_i] + \Pr[B_j] +$

$\Pr[B_0] = 1$. Hence, by solving for $\Pr[A]$ and by using elementary calculus we obtain the following lower bound:

$$\begin{aligned} \Pr[A] &= 1 - \frac{1}{2\pi} (\arccos(v_0v_i) + \arccos(v_0v_j) + \arccos(v_iv_j)) \\ &\geq \frac{\beta}{4} (1 + v_0v_i + v_0v_j + v_iv_j). \end{aligned}$$

Combining the above lower bounds, we obtain

$$\begin{aligned} \mathbb{E}[f_\alpha(S)] &\geq \beta \left(\sum_{e=(i,j)} \frac{1 + y_{0i} + y_{0j} + y_{ij}}{4} \right) + \\ &\quad \frac{\alpha\beta}{4} \left(\sum_{i \neq j} \frac{1 - y_{0i}}{2} + \sum_{i \neq j} \frac{1 - y_{0j}}{2} + \sum_{i \neq j} \frac{1 - y_{ij}}{2} \right) \\ &= \beta \text{SDP}^* \geq \beta \text{IP}^*. \end{aligned}$$

Now, we boost the probability of success of our randomized algorithm. Specifically, we analyze part 3 of our algorithm. First, we lower bound IP^* . Consider adding each vertex with probability $\frac{1}{2}$ to S . The expected value of the objective is $\frac{m}{4} + \alpha \left(\binom{n}{2} - \binom{n/2}{2} \right) \approx \frac{m}{4} + \alpha \frac{3n^2}{8} \geq \frac{3\alpha}{2} \binom{n}{2}$. Hence,

$$\text{IP}^* \geq \frac{3\alpha}{2} \binom{n}{2}.$$

Also notice that the objective is upper bounded always by $(\alpha + 1) \binom{n}{2}$. Define a constant γ as

$$\gamma = \frac{\mathbb{E}[f_\alpha(S)]}{(\alpha + 1) \binom{n}{2}}.$$

Note that $\gamma \leq 1$. We obtain the following lower bound on γ :

$$(\alpha + 1) \binom{n}{2} \geq \mathbb{E}[f_\alpha(S)] = \gamma(\alpha + 1) \binom{n}{2} \geq \beta \text{IP}^* \geq \beta \frac{3\alpha}{2} \binom{n}{2},$$

and hence $1 \geq \gamma \geq \frac{3\alpha\beta}{2(\alpha+1)}$.

Let $p = \Pr[W < (1 - \epsilon)\mathbb{E}[f_\alpha(S)]]$, where W is the actual objective value achieved by our randomized algorithm. We obtain the following (generous) upper bound on p as follows:

$$\mathbb{E}[f_\alpha(S)] \leq p(1 - \epsilon)\mathbb{E}[f_\alpha(S)] + (1 - p)(\alpha + 1) \binom{n}{2},$$

which by solving for p gives

$$p \leq 1 - \frac{\epsilon\gamma}{1 - \gamma + \epsilon\gamma} \leq 1 - \frac{\epsilon \frac{3\alpha\beta}{2(\alpha+1)}}{1 - \frac{(1-\epsilon)3\alpha\beta}{2(\alpha+1)}} = 1 - q.$$

Let $c_{\epsilon,\alpha,\beta} \geq \frac{2(\alpha+1)}{3\epsilon\alpha\beta} + 1$. Running the algorithm $c_{\epsilon,\alpha,\beta} \log n \geq \frac{1}{q} \log n$ times gives that the success probability is $1 - o(1)$, i.e.,

$$\Pr[f_\alpha(S) \geq (1 - \epsilon)\mathbb{E}[f_\alpha(S)]] \geq 1 - (1 - q)^{\frac{\log n}{q}} \approx 1 - \mathcal{O}(n^{-1}).$$

□

6.4 Problem variants

We present two variants of our basic problem, that have many practical applications: finding top- k optimal quasi-cliques (Section 6.4.1) and finding an optimal quasi-clique that contains a given set of *query vertices* (Section 6.4.2).

6.4.1 Top- k optimal quasi-cliques

The top- k version of our problem is as follows: given a graph $G = (V, E)$ and a constant k , find top- k *disjoint* optimal quasi-cliques. This variant is particularly useful in scenarios where finding a single dense subgraph is not sufficient, rather a set of $k > 1$ dense components is required.

From a formal viewpoint, the problem would require to find k subgraphs for which the sum of the various objective function values computed on each subgraph is maximized. Due to its intrinsic hardness, however, here we heuristically tackle the problem in a greedy fashion: we find one dense subgraph at a time, we remove all the vertices of the subgraph from the graph, and we continue until we find k subgraphs or until we are left with an empty graph. Note that this iterative approach allows us to automatically fulfil a very common requirement of finding top- k subgraphs that are pairwise disjoint.

6.4.2 Constrained optimal quasi-cliques

The constrained optimal quasi-cliques variant consists in finding an optimal quasi-clique that contains a set of pre-specified *query vertices*. This variant is inspired by the *community-search problem* [368], which has many applications, such as finding thematic groups, organizing social events, tag suggestion. Next, we formalize the problem, prove

that it is **NP**-hard, and adapt our scalable algorithms (i.e., GREEDYOQC and LOCALSEARCHOQC) for this variant.

Let $G = (V, E)$ be a graph, and $Q \subseteq V$ be a set of query vertices. We want to find a set of vertices $S \subseteq V$, so that the induced subgraph contains the query vertices Q and maximizes our objective function f_α . Formally, we define the following problem.

Problem 2 (CONSTRAINED-OQC-PROBLEM). Given a graph $G = (V, E)$ and set $Q \subseteq V$, find $S^* \subseteq V$ such that $f_\alpha(S^*) = \max_{Q \subseteq S \subseteq V} f_\alpha(S)$.

It is easy to see that, when $Q = \emptyset$, the CONSTRAINED-OQC-PROBLEM reduces to the OQC-PROBLEM. The following hardness result is immediate from Theorem 1 in [399].

Theorem 6.6. *The CONSTRAINED-OQC-PROBLEM is **NP**-hard.*

The GREEDYOQC algorithm can be adapted to solve the CONSTRAINED-OQC-PROBLEM simply by ignoring the nodes $u \in Q$ during the execution of the algorithm, so as to never remove vertices of Q .

Similarly, our LOCALSEARCHOQC algorithm can solve the CONSTRAINED-OQC-PROBLEM with a couple of simple modifications: the set S is initialized to the set of query vertices Q , while, during the iterative phase of the algorithm, we never allow a vertex $u \in Q$ to leave S .

6.5 Experimental evaluation

In this section we present our empirical evaluation, first on publicly available real-world graphs (Section 6.5.1), whose main characteristics are shown in Table 6.2, and then on synthetic graphs where the ground truth is known (Section 6.5.2).

We compare our optimal quasi-cliques with densest subgraphs. The latter is chosen as our baseline given its popularity. For extracting optimal quasi-cliques we use our scalable algorithms, GREEDYOQC and LOCALSEARCHOQC. As far as the semidefinite-programming algorithm presented in Section 6.3, we recall that it has been introduced mainly to show theoretical properties of the problem tackled in this paper. From a practical viewpoint, we have been able to run it only on the smallest datasets. We do not present the results we obtained with the semidefinite-programming algorithm presented in Section 6.3, which we implemented using SDTP3 [381] in MATLAB for the following reasons: (i) it does not scale to large networks; (ii) the results are inferior to the results of the scalable algorithms for the small graphs we tested it on. For instance, for the *polbooks*

	Vertices	Edges	Description
Dolphins	62	159	Biological Network
Polbooks	105	441	Books Network
Adjnoun	112	425	Adj. and Nouns in 'David Copperfield'
Football	115	613	Games Network
Jazz	198	2 742	Musicians Network
Celegans N.	297	2 148	Biological Network
Celegans M.	453	2 025	Biological Network
Email	1 133	5 451	Email Network
AS-22july06	22 963	48 436	Auton. Systems
Web-Google	875 713	3 852 985	Web Graph
Youtube	1 157 822	2 990 442	Social Network
AS-Skitter	1 696 415	11 095 298	Auton. Systems
Wikipedia 2005	1 634 989	18 540 589	Web Graph
Wikipedia 2006/9	2 983 494	35 048 115	Web Graph
Wikipedia 2006/11	3 148 440	37 043 456	Web Graph

TABLE 6.2: Graphs used in our experiments.

network, the corresponding edge density is 0.19 compared to 0.67 and 0.61 that we obtain using GREEDYOQC and LOCALSEARCHOQC respectively; (*iii*) it serves mainly as a theoretical contribution.

Following our discussion in Section 6.3, we run our algorithms with $\alpha = \frac{1}{3}$. For LOCALSEARCHOQC, we set $T_{\max} = 50$. For finding densest subgraphs, we use the Goldberg's exact algorithm [199] for small graphs, while for graphs whose size does not allow the Goldberg's algorithm to terminate in reasonable time we use Charikar's approximation algorithm [104].

All algorithms are implemented in JAVA, and all experiments are performed on a single machine with Intel Xeon CPU at 2.83GHz and 50GB RAM.

6.5.1 Real-world graphs

We experiment with the real graphs in Table 6.2. The results are shown in Table 6.3. We compare optimal quasi-cliques outputted by the GREEDYOQC and LOCALSEARCHOQC algorithms with densest subgraphs extracted with Charikar's algorithm. Particularly, we use Charikar's method to be able to handle the largest graphs. For consistence, Table 6.3 reports on results achieved by Charikar's method also for the smallest graphs. We recall that the results in Table 6.1 refer instead to the exact Goldberg's method. However, a comparison of the two tables on their common rows shows that Charikar's algorithm, even though it is approximate, produces almost identical results with the results produced by Goldberg's algorithm.

Table 6.3 clearly confirms the preliminary results reported in the Introduction: **optimal quasi-cliques** have larger edge and triangle densities, and smaller diameter than **densest subgraphs**. Particularly, the edge density of **optimal quasi-cliques** is evidently larger on all graphs. For instance, on **Football** and **Youtube**, the edge density of **optimal quasi-cliques** (for both the **GREEDYOQC** and **LOCALSEARCHOQC** algorithms) is about 9 times larger than the edge density of **densest subgraphs**, while on **Email** the difference increases up to 20 times (**GREEDYOQC**) and 14 times (**LOCALSEARCHOQC**). Still, the triangle density of the **optimal quasi-cliques** outputted by both **GREEDYOQC** and **LOCALSEARCHOQC** is one order of magnitude larger than the triangle density of **densest subgraphs** on 11 out of 15 graphs.

Comparing our two algorithms, we can see that **LOCALSEARCHOQC** performs generally better than **GREEDYOQC**. Indeed, the edge density achieved by **LOCALSEARCHOQC** is higher than that of **GREEDYOQC** on 10 out of 15 graphs, while the diameter of the **LOCALSEARCHOQC** **optimal quasi-cliques** is never larger than the diameter of the **GREEDYOQC** **optimal quasi-cliques**.

Concerning efficiency, all algorithms are linear in the number of edges of the graph. Charikar's and **GREEDYOQC** algorithm are somewhat slower than **LOCALSEARCHOQC**, but mainly due to bookkeeping. **LOCALSEARCHOQC** algorithm's running times vary from milliseconds for the small graphs (e.g., 0.004s for **Dolphins**, 0.002s for **Celegans N.**), few seconds for the larger graphs (e.g., 7.94s for **Web-Google** and 3.52s for **Youtube**) and less than one minute for the largest graphs (e.g., 59.27s for **Wikipedia 2006/11**).

Top- k optimal quasi-cliques. Figure 6.1 evaluates top- k **optimal quasi-cliques** and top- k **densest subgraphs** on the **AS-Skitter** and **Wikipedia 2006/11** graphs using the iterative method described in Section 6.4.1. Similar results hold for the other graphs but are omitted due to space constraints.

For each graph we show two scatterplots. The x axis in logarithmic scale reports the size of each of the top- k dense components, while the y axes show the edge density and the diameter, respectively. In all figures, **optimal quasi-cliques** correspond to blue filled circles (**LOCALSEARCHOQC**) or red diamonds (**GREEDYOQC**), while **densest subgraphs** correspond to green circles. It is evident that **optimal quasi-cliques** are significantly better in terms of both edge density and diameter also in this top- k variant. The edge density is in the range 0.4 – 0.7 and the diameter is always 2 or 3, except for a 56-vertex clique in **Wikipedia 2006/11** with diameter 1. On the contrary, the **densest subgraphs** are large graphs, with diameter ranging typically from 3 to 5, with significantly smaller edge densities: besides few exceptions, the edge density of **densest subgraphs** is always around 0.1 or even less.

	S			δ			D			τ		
	densest	opt. quasi-clique		densest	opt. quasi-clique		densest	opt. quasi-clique		densest	opt. quasi-clique	
	subgraph	GREEDY	LS	subgraph	GREEDY	LS	subgraph	GREEDY	LS	subgraph	GREEDY	LS
Dolphins	19	13	8	0.27	0.47	0.68	3	3	2	0.05	0.12	0.32
Polbooks	53	13	16	0.18	0.67	0.61	6	2	2	0.02	0.28	0.24
Adjnoun	45	16	15	0.20	0.48	0.60	3	3	2	0.01	0.10	0.12
Football	115	10	12	0.09	0.89	0.73	4	2	2	0.03	0.67	0.34
Jazz	99	59	30	0.35	0.54	1	3	2	1	0.08	0.23	1
Celeg. N.	126	27	21	0.14	0.55	0.61	3	2	2	0.07	0.20	0.26
Celeg. M.	44	22	17	0.35	0.61	0.67	3	2	2	0.07	0.26	0.33
Email	289	12	8	0.05	1	0.71	4	1	2	0.01	1	0.30
AS-22july06	204	73	12	0.40	0.53	0.58	3	2	2	0.09	0.19	0.20
Web-Google	230	46	20	0.22	1	0.98	3	2	2	0.03	0.99	0.95
Youtube	1874	124	119	0.05	0.46	0.49	4	2	2	0.02	0.12	0.14
AS-Skitter	433	319	96	0.41	0.53	0.49	2	2	2	0.10	0.19	0.13
Wiki '05	24555	451	321	0.26	0.43	0.48	3	3	2	0.02	0.06	0.10
Wiki '06/9	1594	526	376	0.17	0.43	0.49	3	3	2	0.10	0.06	0.11
Wiki '06/11	1638	527	46	0.17	0.43	0.56	3	3	2	0.31	0.06	0.35

TABLE 6.3: Densest subgraphs extracted with Charikar’s method vs. optimal quasi-cliques extracted with the proposed GREEDYOQC algorithm (GREEDY) and LOCAL-SEARCHOQC algorithm (LS). $\delta = e[S]/\binom{|S|}{2}$ is the edge density of the extracted subgraph S , D is the diameter, and $\tau = t[S]/\binom{|S|}{3}$ is the triangle density.

6.5.2 Synthetic graphs

Experiments on synthetic graphs deal with the following task: a (small) clique is planted in two different types of random graphs, and the goal is to check if the dense subgraph algorithms are able to recover those cliques. Two different random-graph models are used as host graphs for the cliques: (i) Erdős-Rényi and (ii) random power-law graphs. In the former model, each edge exists with probability p independently of the other edges. To generate a random power-law graph, we follow the Chung-Lu model [114]: we first generate a degree sequence (d_1, \dots, d_n) that follows a power law with a pre-specified slope and we connect each pair of vertices i, j with probability proportional to $d_i d_j$.

We evaluate our algorithms by measuring how “close” are the returned subgraphs to the planted clique. In particular, we use the measures of *precision* P and *recall* R , defined as

$$P = \frac{\#\{\text{returned vertices from hidden clique}\}}{\text{size}\{\text{subgraph returned}\}}, \text{ and}$$

$$R = \frac{\#\{\text{returned vertices from hidden clique}\}}{\text{size}\{\text{hidden clique}\}}.$$

Next we discuss the results obtained. For the Erdős-Rényi model we also provide a theoretical justification of the outcome of the two tested algorithms.

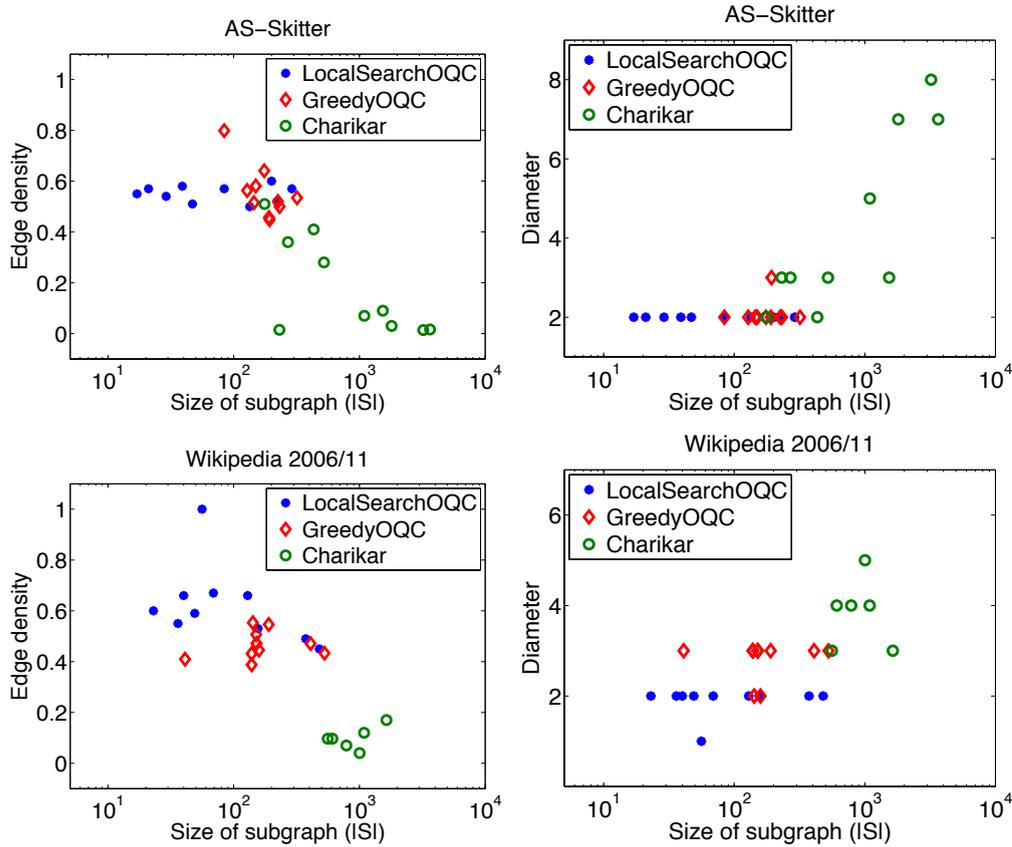


FIGURE 6.1: Edge density and diameter of the top-10 subgraphs found by our GREEDYOQC and LOCALSEARCHOQC methods, and Charikar’s algorithm, on the AS-skitter graph (top) and the Wikipedia 2006/11 graph (bottom).

Erdős-Rényi graphs. We plant a clique of 30 vertices on Erdős-Rényi graphs with $n = 3000$ and edge probabilities $p \in \{0.5, 0.1, 0.008\}$. Those values of p are selected to represent very dense, medium-dense, and sparse graphs.

We report in Table 6.4 the results of running our LOCALSEARCHOQC and GREEDYOQC algorithms for extracting optimal quasi-cliques, as well as the Goldberg’s algorithm for extracting densest subgraphs. We observe that our two algorithms, LOCALSEARCHOQC and GREEDYOQC, produce *identical* results, thus we refer to both of them as **optimal quasi-cliques algorithms**. We see that the algorithms produce two kinds of results: they either find the hidden clique, or they miss it and return the whole graph. In the very dense setting ($p = 0.5$) all algorithms miss the clique, while in the sparse setting ($p = 0.008$) all algorithms recover it. However, at the middle-density setting ($p = 0.1$) only the **optimal quasi-cliques algorithms** find the clique, while the Goldberg’s algorithm misses it.

To better understand the results shown on Table 6.4, we provide a theoretical explanation of the behavior of the algorithms depending on their objective. Assume that h is the size of the hidden clique. If $np \geq h - 1$ the densest subgraph criterion always returns the

Erdős-Rényi parameters		densest subgraph			optimal quasi-clique		
n	p	$ S $	P	R	$ S $	P	R
3 000	0.5	3 000	0.01	1.00	3 000	0.01	1.00
3 000	0.1	3 000	0.01	1.00	30	1.00	1.00
3 000	0.008	30	1.00	1.00	30	1.00	1.00

TABLE 6.4: Subgraphs returned by Goldberg’s max-flow algorithm and by our two algorithms (GREEDYOQC, LOCALSEARCHOQC) on Erdős-Rényi graphs with 3 000 vertices and three values of p , and with a planted clique of 30 vertices.

whole graph. In our experiments, this happens with $p = 0.5$ and $p = 0.1$. On the other hand, if $np < h - 1$ the densest subgraph corresponds to the hidden clique, and therefore the Goldberg’s algorithm cannot miss it.

Now consider our objective function, i.e., the edge-surplus function f_α , and let us discard for simplicity the constant $\binom{n}{2}$, because this does not affect the validity of the following reasoning. The expected score for the hidden clique is $\mathbb{E}[f_\alpha(H)] = f_\alpha(H) = (1 - \alpha)\binom{h}{2}$. The expected score for the whole network is $\mathbb{E}[f_\alpha(V)] = (p\binom{n}{2} + (1 - p)\binom{h}{2}) - \alpha\binom{n}{2}$. We obtain the following two cases: (A) when $p > \alpha$, we have $\mathbb{E}[f_\alpha(V)] \geq f_\alpha(H)$. (B) when $p < \alpha$, we have $f_\alpha(H) \geq \mathbb{E}[f_\alpha(V)]$. This rough analysis explains our findings.²

Power-law graphs. We plant a clique of 15 vertices in random power-law graphs of again 3 000 vertices, with power-law exponent varying from 2.2 to 3.1. We select these values since most real-world networks have power-law exponent in this range [315]. For each exponent tested, we generate five random graphs, and all the figures we report are averages over these five trials.

Again, we compare our GREEDYOQC and LOCALSEARCHOQC algorithms with the Goldberg’s algorithm. The LOCALSEARCHOQC algorithm is run seeded with one of the vertices of the clique. The justification of this choice is that we can always re-run the algorithm until it finds such a vertex with high probability.³

The precision and recall scores of the three competing algorithms as a function of the power-law exponent are shown in Figure 6.2. As the exponent increases the host graph becomes sparser and both algorithms have no difficulty in finding the hidden clique. However, for exponent values ranging between 2.2 and 2.6 the optimal quasi-cliques are significantly better than the densest subgraphs. Indeed, in terms of precision, the Goldberg’s algorithm is outperformed by both our algorithms. In terms of recall, our LOCALSEARCHOQC is better than Goldberg’s, while our GREEDYOQC performs

²The analysis can be tightened via Chernoff bounds, but we avoid this here due to space constraints.

³If the hidden clique is of size $\mathcal{O}(n^\epsilon)$, for some $0 \leq \epsilon < 1$, a rough calculation shows that it suffices to run the algorithm a sub-linear number of times (i.e., $\mathcal{O}((1 - \gamma)n^{1-\epsilon})$ times) in order to obtain one of the vertices of the clique as a seed with probability at least $1 - \gamma$.

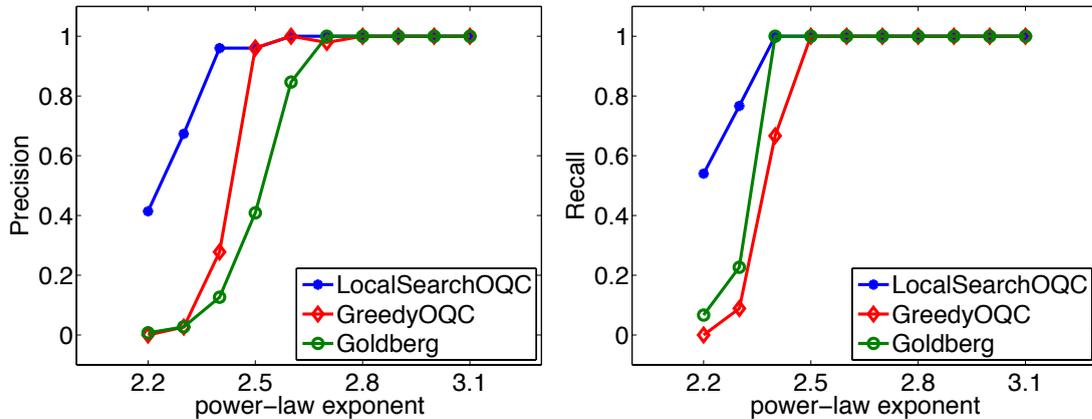


FIGURE 6.2: Precision and recall for our method and Goldberg’s algorithm vs. the power-law exponent of the host graph.

slightly worse. An explanation for this is that the GREEDYOQC algorithm detects other high-density subgraphs, but not exactly the planted clique. As an example, with power-law exponent 2.3, GREEDYOQC finds a subgraph with 23 vertices and edge density 0.87.

Stability with respect to α . We also test the sensitivity of our density measure with respect to the parameter α . We use again the planted-clique setting, and we test the ability of our algorithms to recover the clique as we vary the parameter α . We omit detailed plots, due to space constraints, but we report that the behavior of both algorithms is extremely stable with respect to α . Essentially, the algorithms again either find the clique or miss it, depending on the graph-generation parameters, as we saw in the previous section, namely, the probability p of the Erdős-Rényi graphs, or the exponent of the power-law graphs. Moreover, in all cases, the performance of our algorithms, measured by precision and recall as in the last experiment, does not depend on α .

6.6 Applications

In this section we show experiments concerning our constrained optimal quasi-cliques variant introduced in Section 6.4.2. To this end, we focus on two applications that can be commonly encountered in real-world scenarios: finding thematic groups and finding highly correlated genes from a microarray dataset. For the sake of brevity of presentation, we show next results for only one of our scalable algorithms, particularly the LOCALSEARCHOQC algorithm.

Abiteboul, Bernstein, Brodie, Carey, Ceri, Crof,
 DeWitt, Ehrenfeucht, Franklin, Gawlick, Gray, Haas,
 Halevy, Hellerstein, Ioannidis, Jagadish, Kanellakis,
 Kersten, Lesk, Maier, Molina, Naughton, Papadimitriou,
 Pazzani, Pirahesh, Schek, Sellis, Silberschatz, Snodgrass,
 Stonebraker, Ullman, Weikum, Widom, Zdonik

FIGURE 6.3: Authors returned by our LOCALSEARCHOQC algorithm when queried with Papadimitriou and Abiteboul. The set includes well-known database scientists. The induced subgraph has 34 vertices and 457 edges. The edge density is 0.81, the diameter is 3, the triangle density is 0.66.

Alt, Blum, Garey, Guibas, Johnson,
 Karp, Mehlhorn, Papadimitriou, Preparata,
 Tarjan, Welzl, Widgerson, Yannakakis,

FIGURE 6.4: Authors returned by our LOCALSEARCHOQC algorithm when queried with Papadimitriou and Blum. The set includes well-known theoretical computer scientists. The induced subgraph has 13 vertices and 38 edges. The edge density is 0.49, the diameter is 3, the triangle density is 0.14.

6.6.1 Thematic groups

Motivation. Suppose that a set of scientists Q wants to organize a workshop. How do they invite other scientists to participate in the workshop so that the set of all the participants, including Q , have similar interests?

Setup. We use a co-authorship graph extracted from the DBLP dataset. The dataset contains publications in all major computer-science journals. There is an undirected edge between two authors if they have coauthored a journal article. Taking the largest connected component gives a graph of 226K vertices and 1.4M edges.

We evaluate the results of our algorithm qualitatively, in a sanity check form rather than a strict and quantitative way, which is not even well-defined. We perform the following two queries: $Q_1 = \{\text{Papadimitriou, Abiteboul}\}$ and $Q_2 = \{\text{Papadimitriou, Blum}\}$.

Results. Papadimitriou is one of the most prolific computer scientists and has worked on a wide range of areas. With query Q_1 we invoke his interests in database theory given that Abiteboul is an expert in this field. As we can observe from Figure 6.3, the optimal quasi-clique outputted contains database scientists. On the other hand, with query Q_2 we invoke Papadimitriou's interests in theory, given that Blum is a Turing-award theoretical computer scientist. As we can see in Figure 6.4, the returned optimal quasi-clique contains well-known theoretical computer scientists.

p53, BRCA1, ARID1A, ARID1B, ZNF217, FGFR1, KRAS,
NCOR1, PIK3CA, APC, MAP3K13, STK11, AKT1, RB1

FIGURE 6.5: Genes returned by our LOCALSEARCHOQC algorithm when queried with p53. The induced subgraph is a clique with 14 vertices.

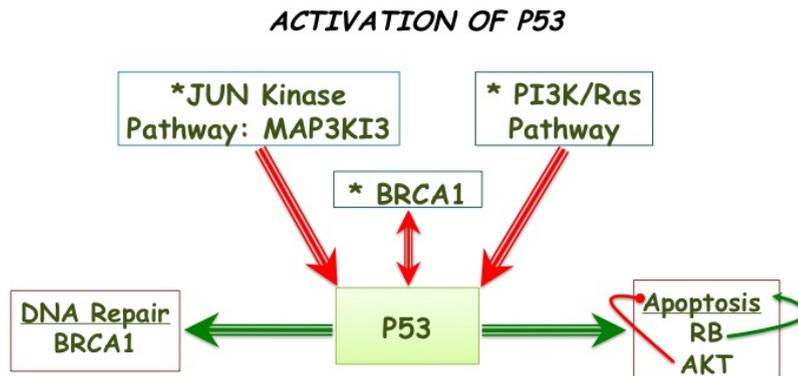


FIGURE 6.6: A tumorigenesis pathway consistent with our findings.

6.6.2 Correlated genes

Motivation. Detecting correlated genes has several applications. For instance, clusters of genes with similar expression levels are typically under similar transcriptional control. Furthermore, genes with similar expression patterns may imply co-regulation or relationship in functional pathways. Detecting gene correlations has played a key role in discovering unknown types of breast cancer [366]. Here, we wish to illustrate that optimal quasi-cliques provide a useful graph theoretic framework for gene co-expression network analysis [267], without delving deeply into biological aspects of the results.

Setup. We use the publicly available breast-cancer dataset of van de Vijner et al. [402], which consists of measurements across 295 patients of 24 479 probes. Upon running a standard probe selection algorithm based on Singular Value Decomposition (SVD), we obtain a 295×1000 matrix. The graph G in input to our LOCALSEARCHOQC algorithm is derived using the well-established approach defined in [267]: each gene corresponds to a vertex in G , while an edge between any pair of genes i, j is drawn if and only if the modulus of the Pearson’s correlation coefficient $|\rho(i, j)|$ exceeds a given threshold θ ($\theta = 0.99$ in our setting). We perform the following query, along the lines of the previous section: “find highly correlated genes with the tumor protein 53 (p53)”. We select p53 since it is known to play a key role in tumorigenesis.

Results. The output of our algorithm is a clique consisting of 14 genes shown in Figure 6.5. A potential explanation of our finding is the pathway depicted in Figure 6.6,

which shows that the activation of the p53 signaling can be initiated by signals coming from the PI3K/AKT pathway. Both PI3KCA and AKT1 that are detected by our method are key players of this pathway. Furthermore, signals from the JUN kinase pathway can also trigger the p53-cascade; MAP3K13 is a member of this pathway.

One of the results of p53 signaling is apoptosis, a process promoted by RB. The latter can also regulate the stability and the apoptotic function of p53. Finally, our output includes BRCA1, which is known to physically associate with p53 and affect its actions [413].

Chapter 7

Structure of the Web Graph

7.1 Introduction

In this Chapter we design scalable algorithms which allow us to understand better the structure of information, social and economic networks, including the Web graph, a prominent information network. Our goal is to shed light on important properties of the Web graph and other large real-world networks such as the effective diameter, number and sizes of connected components and temporal patterns of evolution. Specifically, we perform experiments on the *YahooWeb* graph with 60 billion edges. We reveal facts about the structure of the Web, including the well known small-world phenomenon, the multimodal shape of the radius distribution and time-varying patterns.

Analyzing networks of this scale is a challenge. For this purpose we design HADI, a MAPREDUCE algorithm which scales to large-scale networks. Our algorithm relies on approximate counting of distinct elements of a multiset. We use the algorithm of Flajolet and Martin [168] as our “black-box”, which is historically the first algorithm proposed for this problem. Since then, many other algorithms have been proposed, including the Hyperloglog counters [169]. It is worth mentioning that the current state-of-the-art method [75] which outperforms HADI is based on the latter counters [169], which provide an exponential space improvement over the original Flajolet-Martin counters [168].

The rest of the Chapter is organized as follows: Section 7.2 defines the notions of effective radius and diameter and reviews briefly related work. Section 7.3 provides an explanation of how the number of distinct elements of multisets is related to finding the diameter of a graph. Section 7.4 presents algorithmic tools for finding the distribution of radii and the diameter of large-scale networks. Section 7.5 provides wall-clock times. Section 7.6 presents findings on the structure of real-world networks.

7.2 Related Work

Definitions: First, we review basic graph theoretic definitions [83] and then introduce the notions of effective radius and diameter. Let $G(V, E)$ be a directed graph. The radius/eccentricity of a vertex v is the greatest shortest-path distance between v and any other vertex. The radius $r(G)$ is the minimum radius of any vertex. The diameter $d(G)$ is the maximum radius of any vertex. Since the radius and the diameter are susceptible to outliers (e.g., long chains), we follow the literature [279] and define the *effective* radius and diameter as follows.

Definition 7.1 (Effective Radius). For a node v in a graph G , the effective radius $r_{eff}(v)$ of v is the 90th-percentile of all the shortest distances from v .

Definition 7.2 (Effective Diameter). The effective diameter $d_{eff}(G)$ of a graph G is the minimum number of hops in which 90% of all connected pairs of nodes can reach each other.

In Section 7.6 we use the following three radius-based Plots:

1. **Static Radius Plot** (or just “Radius plot”) of graph G shows the distribution (count) of the effective radius of nodes at a specific time.
2. **Temporal Radius Plot** shows the distributions of effective radius of nodes at several times.
3. **Radius-Degree Plot** shows the scatter-plot of the effective radius $r_{eff}(v)$ versus the degree d_v for each node v .

Computing Radius and Diameter: Typical algorithms to compute the radius and the diameter of a graph include Breadth First Search (BFS) and Floyd’s algorithm [125] when no negative cycles are present. Both approaches are prohibitively slow for large-scale graphs, requiring $O(n^2 + nm)$ and $O(n^3)$ time respectively. For the same reason, related BFS or all-pair shortest-path based algorithms like [46, 166, 292, 364] can not handle large-scale graphs.

A sampling approach starts BFS from a subset of nodes, typically chosen at random as in [90]. Despite its practicality, this approach has no obvious solution for choosing the representative sample for BFS. An interesting approach has been proposed by Cohen [119], but according to practitioner’s experience [74] it appears not to be as scalable as the ANF algorithm [327]. The latter is closely related to our work since it is a sequential

algorithm based on Flajolet-Martin sketches [168]. We review its key idea in the next Section.

Distinct Elements in Multisets: Let $A = \{a_1, \dots, a_m\}$ be a multiset where $a_i \in [n]$ for all $i = 1, \dots, m$. Let $m_i = |\{j : a_j = i\}|$. For each $k \geq 0$ define $F_k = \sum_{i=1}^n m_i^k$. The numbers F_k are called frequency moments of the multiset and provide useful statistics. We notice that F_0 is the number of distinct elements in A , $F_1 = m$. Historically, Morris was the first to show that F_1 can be approximated with $O(\log \log m) = O(\log \log n)$ [309]. Flajolet and Martin designed an algorithm that needs $O(\log n)$ bits of memory to approximate F_0 [168]. Since then, many excellent researches have looked into this problem, see [27, 48, 49, 67, 119, 142, 169, 219, 418]. Recently, Kane, Nelson and Woodruff provided an optimal algorithm for estimating F_0 [234].

7.3 Distinct Elements and the Diameter

In this Section, we sketch the key idea of [327], which shows how one can use a space efficient algorithm for estimating distinct elements in a multiset to estimate the diameter and radius distribution of a graph. Assume that for each vertex v in the graph, we maintain the number of neighbors reachable from v within h hops. As h increases, the number of neighbors increases until it stabilizes. The diameter is h where the number of neighbors within $h + 1$ does not increase for every node.

To generate the Radius plot, we need to calculate the effective radius of every node. In addition, the effective diameter is useful for tracking the evolution of networks. Assume we have a ‘set’ data structure that supports two functions: $add()$ for adding an item, and $size()$ for returning the count of distinct items. With the set, radii of vertices can be computed as follows:

1. For each vertex i , create a set S_i and initialize it by adding i to it.
2. For each vertex i , continue updating S_i by adding 1,2,3,...-step neighbors of i to S_i . When the size of S_i stabilizes for first time, then the vertex i reached its radius. Iterate until all vertices reach their radii.

Although simple and clear, the above algorithm requires $\Omega(n^2)$ space, since there are n vertices and each vertex requires $\Omega(n)$ space. This is prohibitive in practice. We describe the ANF algorithm which serves as the basis of our investigation into the structure of real-world networks [327]. We use the Flajolet-Martin algorithm [168, 327] for counting the number of distinct elements in a multiset. The main idea of the Flajolet-Martin

algorithm is as follows. We maintain a bitstring $BITMAP[0 \dots L - 1]$ of length L which encodes the set. For each item we add, we do the following:

1. Pick an $index \in [0 \dots L - 1]$ with probability $1/2^{index+1}$.
2. Set $BITMAP[index]$ to 1.

Let R denote the index of the leftmost ‘0’ bit in $BITMAP$. It is clear that 2^R should be a good estimate of the number of distinct elements. The main result of Flajolet-Martin is that the unbiased estimate of the size of the set is given by

$$\frac{1}{\varphi} 2^R \quad (7.1)$$

where $\varphi = 0.77351 \dots$. A more concentrated estimate can be obtained by using multiple bitstrings and averaging the R . If we use K bitstrings R_1 to R_K , the size of the set can be estimated by

$$\frac{1}{\varphi} 2^{\frac{1}{K} \sum_{l=1}^K R_l} \quad (7.2)$$

The application of the Flajolet-Martin algorithm to radius and diameter estimation is straight-forward. We maintain K Flajolet-Martin (FM) bitstrings $b(h, i)$ for each vertex i and the current hop number h . $b(h, i)$ encodes the number of vertices reachable from vertex i within h hops, and can be used to estimate radii and diameter as shown below. The bitstrings $b(h, i)$ are iteratively updated until the bitstrings of all vertices stabilize. At the h -th iteration, each vertex receives the bitstrings of its neighboring vertices, and updates its own bitstrings $b(h - 1, i)$ handed over from the previous iteration:

$$b(h, i) = b(h - 1, i) \text{ BIT-OR } \{b(h - 1, j) | (i, j) \in E\} \quad (7.3)$$

where ‘BIT-OR’ denotes bitwise-OR function. After h iterations, a vertex i has K bitstrings that encode the *neighborhood function* $N(h, i)$, that is, the number of vertices within h hops from the vertex i . $N(h, i)$ is estimated from the K bitstrings by

$$N(h, i) = \frac{1}{0.77351} 2^{\frac{1}{K} \sum_{l=1}^K b_l(i)} \quad (7.4)$$

where $b_l(i)$ is the position of leftmost ‘0’ bit of the l^{th} bitstring of vertex i . The iterations continue until the bitstrings of all vertices stabilize, which is a necessary condition that

the current iteration number h exceeds the diameter $d(G)$. After the iterations finish at h_{max} , we calculate the effective radius for every node and the diameter of the graph, as follows:

- $r_{eff}(i)$ is the smallest h such that $N(h, i) \geq 0.9 \cdot N(h_{max}, i)$.
- $d_{eff}(G)$ is the smallest h such that $N(h) = \sum_i N(h, i) = 0.9 \cdot N(h_{max})$. If $N(h) > 0.9 \cdot N(h_{max}) > N(h - 1)$, then $d_{eff}(G)$ is linearly interpolated from $N(h)$ and $N(h - 1)$. That is, $d_{eff}(G) = (h - 1) + \frac{0.9 \cdot N(h_{max}) - N(h - 1)}{N(h) - N(h - 1)}$.

The parameter K is typically set to 32 [168], and $MaxIter$ is set to 256 since real graphs have relatively small effective diameter.

7.4 Hadoop Algorithms

In this Section we present a way to implement ANF [327] on the top of both a MAPREDUCE system and a parallel SQL DBMS. Our algorithm is HADI, a parallel radius and diameter estimation algorithm. It is important to notice that HADI is a disk-based algorithm. HADI saves two pieces of information to a distributed file system (such as HDFS (Hadoop Distributed File System) in the case of HADOOP):

- **Edge** has a format of $(srcid, dstid)$.
- **Bitstrings** has a format of $(nodeid, bitstring_1, \dots, bitstring_K)$.

Section 7.4.1 presents HADI-naive which gives the big picture and explains why this kind of an implementation should not be used in practice. Section 7.4.2 presents HADI-plain, a significantly improved implementation. Section 7.4.3 presents the optimized version of HADI-optimized, which scales almost linearly as a function of the number of machines allocated. Section 7.4.4 presents the space and time complexity of HADI. Finally, Section 7.4.5 shows how HADI can run on the top of a relational database management system (RDBMS).

7.4.1 HADI-naive in MapReduce

Data: The edge file is saved as a sparse adjacency matrix in HDFS. Each line of the file contains a nonzero element of the adjacency matrix of the graph, in the format of $(srcid, dstid)$. Also, the bitstrings of each node are saved in a file in the format of $(nodeid,$

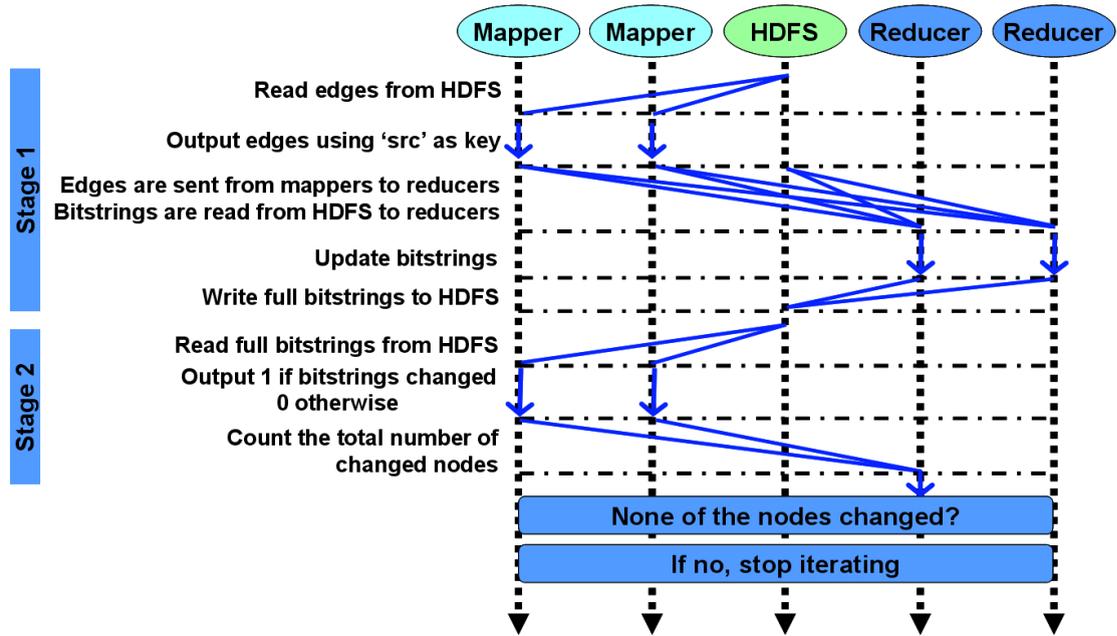


FIGURE 7.1: One iteration of HADI-naive. *Stage 1*. Bitstrings of all vertex are sent to every reducer. *Stage 2*. Sums up the count of changed nodes.

$flag, bitstring_1, \dots, bitstring_K$). The $flag$ variable records whether a bitstring changed or not.

Main Program Flow The main idea of HADI-naive is to use the bitstrings file as a logical “cache” to machines which contain edge files. The bitstring update operation in Equation (7.3) of Section 7.3 requires that the machine which updates the bitstrings of node i should have access to (a) all edges adjacent from i , and (b) all bitstrings of the adjacent nodes. To meet the requirement (a), it is needed to reorganize the edge file so that edges with a same source id are grouped together. That can be done by using an identity mapper which outputs the given input edges in $(srcid, dstid)$ format. The most simple yet naive way to meet the requirement (b) is sending the bitstrings to every reducer which receives the reorganized edge file.

Thus, HADI-naive iterates over two-stages of MAPREDUCE. The first stage updates the bitstrings of each node and sets the ‘Changed’ flag if at least one of the bitstrings of the node is different from the previous bitstring. The second stage counts the number of changed vertex and stops iterations when the bitstrings stabilized, as illustrated in the swim-lane diagram of Figure 7.1.

Although conceptually simple and clear, HADI-naive is unnecessarily expensive, because it ships all the bitstrings to all reducers.

7.4.2 HADI-plain in MapReduce

HADI-plain improves HADI-naive by *copying only the necessary bitstrings to each reducer*. The details follow.

Data: As in HADI-naive, the edges are saved in the format of $(srcid, dstid)$, and bitstrings are saved in the format of $(nodeid, flag, bitstring_1, \dots, bitstring_K)$ in files over HDFS. The initial bitstrings generation can be performed in completely parallel way. The *flag* of each node records the following information:

- *Effective Radii* and *Hop Numbers* to calculate the effective radius.
- *Changed* flag to indicate whether at least a bitstring has been changed or not.

Main Program Flow: As mentioned in the beginning, HADI-plain copies only the necessary bitstrings to each reducer. The main idea is to replicate bitstrings of node j exactly x times where x is the in-degree of node j . The replicated bitstrings of node j is called the *partial bitstring* and represented by $\hat{b}(h, j)$. The replicated $\hat{b}(h, j)$'s are used to update $b(h, i)$, the bitstring of node i where (i, j) is an edge in the graph. HADI-plain iteratively runs three-stage MAPREDUCE jobs until all bitstrings of all vertex stop changing. Algorithm 7, 8, and 9 shows HADI-plain, and Figure 7.2 shows the swim-lane. We use h for denoting the current iteration number which starts from $h=1$. Output(a, b) means to output a pair of data with the key a and the value b .

Stage 1 We generate (key, value) pairs, where the key is the node id i and the value is the partial bitstrings $\hat{b}(h, j)$'s where j ranges over all the neighbors adjacent from node i . To generate such pairs, the bitstrings of node j are grouped together with edges whose *dstid* is j . Notice that at the very first iteration, bitstrings of vertex do not exist; they have to be generated on the fly, and we use the *Bitstring Creation Command* for that. The `NewFMBitstring()` function generates K FM bitstrings [168]. Notice also that line 19 of Algorithm 7 is used to propagate the bitstrings of one's own node. These bitstrings are compared to the newly updated bitstrings at Stage 2 to check convergence.

Stage 2 Bitstrings of node i are updated by combining partial bitstrings of itself and vertex adjacent from i . For the purpose, the mapper is the Identity mapper (output the input without any modification). The reducer combines them, generates new bitstrings, and sets *flag* by recording (a) whether at least a bitstring changed or not, and (b) the current iteration number h and the neighborhood value $N(h, i)$. This h and $N(h, i)$ are used to calculate the effective radius of vertex after all bitstrings converge. Notice that only the last neighborhood $N(h_{last}, i)$ and other neighborhoods $N(h', i)$ that satisfy

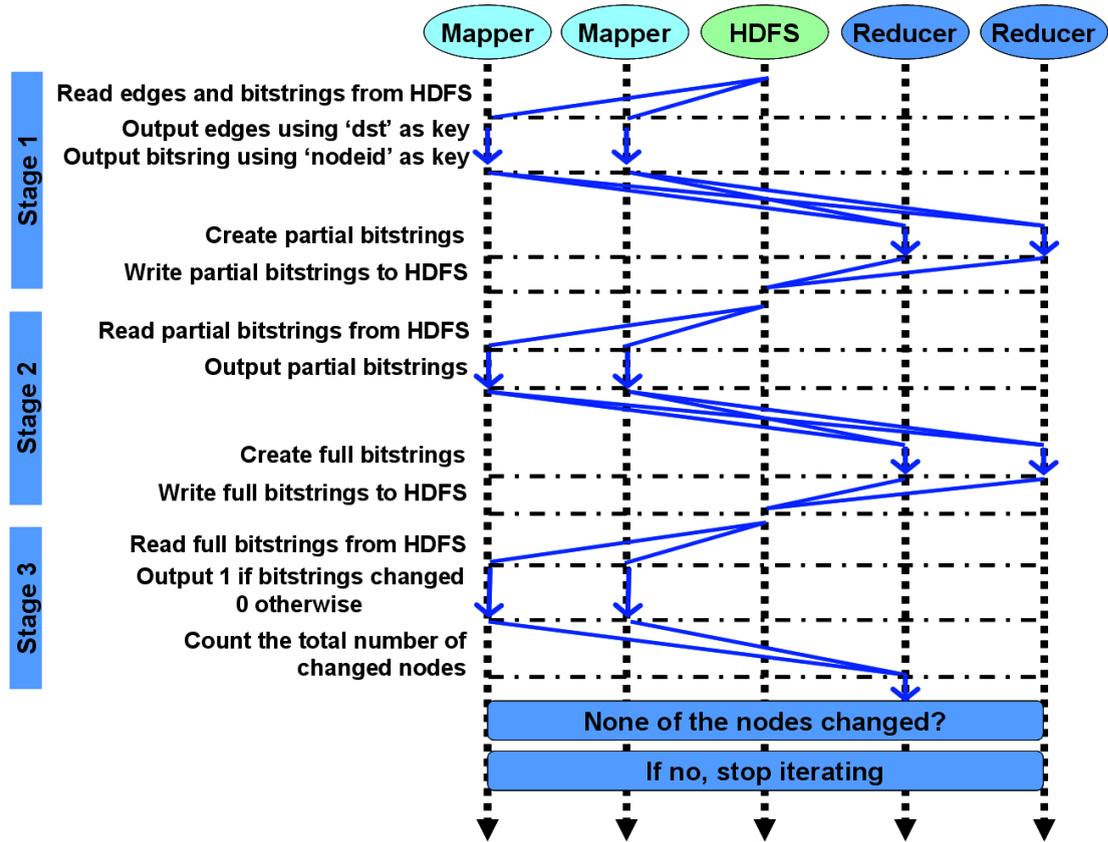


FIGURE 7.2: One iteration of HADI-plain. *Stage 1.* Edges and bitstrings are matched to create partial bitstrings. *Stage 2.* Partial bitstrings are merged to create full bitstrings. *Stage 3.* Sums up the count of changed nodes, and compute $N(h)$, the neighborhood function. Computing $N(h)$ is not drawn in the figure for clarity.

$N(h', i) \geq 0.9 \cdot N(h_{last}, i)$ need to be saved to calculate the effective radius. The output of Stage 2 is fed into the input of Stage 1 at the next iteration.

Stage 3 We calculate the number of changed vertex and sum up the neighborhood value of all vertex to calculate $N(h)$. We use only two unique keys(key_for_changed and key_for_neighborhood), which correspond to the two calculated values. The analysis of line 2 can be done by checking the *flag* field and using Equation (7.4) in Section 7.3. The variable *changed* is set to 1 or 0, based on whether the bitmask of node k changed or not.

When all bitstrings of all vertex converge, a MAPREDUCE job to finalize the effective radius and diameter is performed and the program finishes. Compared to HADI-naive, the advantage of HADI-plain is clear: bitstrings and edges are evenly distributed over machines so that the algorithm can handle as much data as possible, given sufficiently many machines.

Algorithm 7 HADI Stage 1. Output is a set of partial bitstrings $B' = \{(i, b(h-1, j))\}$

Input: Edge data $E = \{(i, j)\}$

Input: Current bitstring $B = \{(i, b(h-1, i))\}$ or Bitstring Creation Command $BC = \{(i, cmd)\}$

```

1: Stage1-Map(key  $k$ , value  $v$ )
2: if  $(k, v)$  is of type B or BC then
3:   Output( $k, v$ )
4: else if  $(k, v)$  is of type E then
5:   Output( $v, k$ )
6: end if
7: Stage1-Reduce(key  $k$ , values  $V[]$ )
8: SRC  $\leftarrow []$ 
9: for  $v \in V$  do
10:  if  $(k, v)$  is of type BC then
11:     $\hat{b}(h-1, k) \leftarrow \text{NewFMBitstring}()$ 
12:  else if  $(k, v)$  is of type B then
13:     $\hat{b}(h-1, k) \leftarrow v$ 
14:  else if  $(k, v)$  is of type E then
15:    Add  $v$  to SRC
16:  end if
17: end for
18: for  $src \in SRC$  do
19:   Output( $src, \hat{b}(h-1, k)$ )
20: end for
21: Output( $k, \hat{b}(h-1, k)$ )

```

7.4.3 HADI-optimized in MapReduce

HADI-optimized further improves HADI-plain. It uses two orthogonal ideas: “block operation” and “bit shuffle encoding”. Both try to address some subtle performance issues. Specifically, HADOOP has the following two major bottlenecks:

- Materialization: at the end of each map/reduce stage, the output is written to the disk, and it is also read at the beginning of next reduce/map stage.
- Sorting: at the *Shuffle* stage, data is sent to each reducer and sorted before they are handed over to the *Reduce* stage.

HADI-optimized addresses these two issues.

Block Operation: Our first optimization is the block encoding of the edges and the bitstrings. The main idea is to group w by w sub-matrix into a super-element in the adjacency matrix E, and group w bitstrings into a super-bitstring. Now, HADI-plain is performed on these super-elements and super-bitstrings, instead of the original edges

Algorithm 8 HADI Stage 2. Output is a set of full bitstring $B = \{(i, b(h, i))\}$

Input: Partial bitstring $B = \{(i, \hat{b}(h-1, j))\}$

```

1: Stage2-Map(key  $k$ , value  $v$ )
   { Identity Mapper }
2: Output( $k, v$ )
3: Stage2-Reduce(key  $k$ , values  $V[]$ )
4:  $b(h, k) \leftarrow 0$ 
5: for  $v \in V$  do
6:    $b(h, k) \leftarrow b(h, k)$  BIT-OR  $v$ 
7: end for
8: Update flag of  $b(h, k)$ 
9: Output( $k, b(h, k)$ )
10: if  $(k, v)$  is of type BC then
11:    $\hat{b}(h-1, k) \leftarrow \text{NewFMBitstring}()$ 
12: else if  $(k, v)$  is of type B then
13:    $\hat{b}(h-1, k) \leftarrow v$ 
14: else if  $(k, v)$  is of type E then
15:   Add  $v$  to SRC
16: end if
17: for  $src \in SRC$  do
18:   Output( $src, \hat{b}(h-1, k)$ )
19: end for
20: Output( $k, \hat{b}(h-1, k)$ )

```

Algorithm 9 HADI Stage 3. Output is the number of changed nodes, Neighborhood $N(h)$

Input: Full bitstring $B = \{(i, b(h, i))\}$

```

1: Stage3-Map(key  $k$ , value  $v$ )
2: Analyze  $v$  to get (changed,  $N(h, i)$ )
3: Output(key_for_changed, changed)
4: Output(key_for_neighborhood,  $N(h, i)$ )
5: Stage3-Reduce(key  $k$ , values  $V[]$ )
6: Changed  $\leftarrow 0$ 
7:  $N(h) \leftarrow 0$ 
8: for  $v \in V$  do
9:   if  $k$  is key_for_changed then
10:    Changed  $\leftarrow$  Changed +  $v$ 
11:   else if  $k$  is key_for_neighborhood then
12:     $N(h) \leftarrow N(h) + v$ 
13:   end if
14: end for
15: Output(key_for_changed, Changed)
16: Output(key_for_neighborhood,  $N(h)$ )

```

and bitstrings. Of course, appropriate decoding and encoding is necessary at each stage. Figure 7.3 shows an example of converting data into block-format.

By this block operation, the performance of HADI-plain changes as follows:

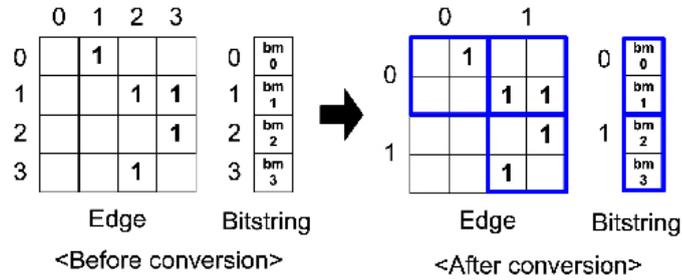


FIGURE 7.3: Converting the original edge and bitstring to blocks. The 4-by-4 edge and length-4 bitstring are converted to 2-by-2 super-elements and length-2 super-bitstrings. Notice the lower-left super-element of the edge is not produced since there is no nonzero element inside it.

- *Input size* decreases in general, since we can use fewer bits to index elements inside a block.
- *Sorting time* decreases, since the number of elements to sort decreases.
- *Network traffic* decreases since the result of matching a super-element and a super-bitstring is a bitstring which can be at maximum *block_width* times smaller than that of HADI-plain.
- *Map and Reduce functions* take more time, since the block must be decoded to be processed, and be encoded back to block format.

For reasonable-size blocks, the performance improvement are significant.

Bit Shuffle Encoding: In our effort to decrease the input size, we propose an encoding scheme that can compress the bitstrings. Recall that in HADI-plain, we use K (e.g., 32, 64) bitstrings for each node, to increase the accuracy of our estimator. Since HADI requires $O(K(m+n)\log n)$ space, the amount of data increases when K is large. For example, the YahooWeb graph spans 120 GBytes (with 1.4 billion nodes, 6.6 billion edges). However the required disk space for just the bitstrings is $32 \cdot (1.4B + 6.6B) \cdot 8$ byte = 2 Tera bytes (assuming 8 byte for each bitstring), which is more than 16 times larger than the input graph.

The main idea of Bit Shuffle Encoding is to carefully reorder the bits of the bitstrings of each node, and then use Run Length Encoding. By construction, the leftmost part of each bitstring is almost full of one's, and the rest is almost full of zeros. Specifically, we make the reordered bit strings to contain long sequences of 1's and 0's: we get all the first bits from all K bitstrings, then get the second bits, and so on. As a result we get a single bit-sequence of length $K \cdot |\text{bitstring}|$, where most of the first bits are '1's, and most of the last bits are '0's. Then we encode only the length of each bit sequence, achieving good space savings (and, eventually, time savings, through fewer I/Os).

7.4.4 Analysis

HADI depends on four parameters: the number M of machines, the number of vertices n and edges m , and the diameter d . The time complexity is dominated by the shuffling time needed for the shuffling during Stage1. Specifically, HADI takes $O(\frac{m+n}{M} \log \frac{m+n}{M})$ time to run. It requires $O((m+n) \log n)$ space units and is required $O((m+n) \log n)$.

7.4.5 Hadi in SQL

Using relational database management systems (RDBMS) for graph mining is a promising research direction, especially given the findings of [334]. We mention that HADI can be implemented on the top of an Object-Relational DBMS (parallel or serial): it needs repeated joins of the edge table with the appropriate table of bit-strings, and a user-defined function for bit-OR-ing. We sketch a potential implementation of HADI in a RDBMS.

Data: In parallel RDBMS implementations, data is saved in tables. The edges are saved in the table E with attributes src (source node id) and dst (destination node id). Similarly, the bitstrings are saved in the table B with

Main Program Flow: The main flow comprises iterative execution of SQL statements with appropriate user defined functions. The most important and expensive operation is updating the bitstrings of nodes. Observe that the operation can be concisely expressed as a SQL statement:

```
SELECT INTO B_NEW E.src, BIT-OR(B.b)
FROM E, B
WHERE E.dst=B.id
GROUP BY E.src
```

The SQL statement requires BIT-OR(), a UDF function that implements the bit-OR-ing of the Flajolet-Martin bitstrings. The RDBMS implementation iteratively runs the SQL until B_NEW is same as B . B_NEW created at an iteration is used as B at the next iteration.

7.5 Wall-clock Times

In this section, we perform experiments to answer the following questions:

- Q1: How fast is HADI?

- Q2: How does it scale up with the graph size and the number of machines?
- Q3: How do the optimizations help performance?

7.5.1 Experimental Setup

We use both real and synthetic graphs in our experiments. These datasets are shown in Table 7.1.

- YahooWeb: web pages and their hypertext links indexed by Yahoo! Altavista search engine in 2002.
- Patents: U.S. patents, citing each other (from 1975 to 1999).
- LinkedIn: people connected to other people (from 2003 to 2006).
- Kronecker: Synthetic Kronecker graphs [279] using a chain of length two as the seed graph.

Graph	Nodes	Edges	File	Description
YahooWeb	1.4 B	6.6 B	116G	page-page
LinkedIn	7.5 M	58 M	1G	person-person
Patents	6 M	16 M	264M	patent-patent
Kronecker	177 K	1,977 M	25G	synthetic
	120 K	1,145M	13.9G	
	59 K	282 M	3.3G	
Erdős-Rényi	177 K	1,977 M	25G	random $G_{n,p}$
	120 K	1,145 M	13.9G	
	59 K	282 M	3.3G	

TABLE 7.1: Datasets (B: Billion, M: Million, K: Thousand, G: Gigabytes)

In order to test the scalability of HADI, we use synthetic graphs, namely Kronecker and Erdős-Rényi graphs. HADI runs on *M45*, one of the fifty most powerful supercomputers in the world. *M45* has 480 hosts (each with 2 quad-core Intel Xeon 1.86 GHz, running RHEL5), with 3Tb aggregate RAM, and over 1.5 Peta-byte disk size.

Finally, we use the following notations to indicate different optimizations of HADI:

- HADI-BSE: HADI-plain with bit shuffle encoding.
- HADI-BL: HADI-plain with block operation.
- HADI-OPT: HADI-plain with both bit shuffle encoding and block operation.

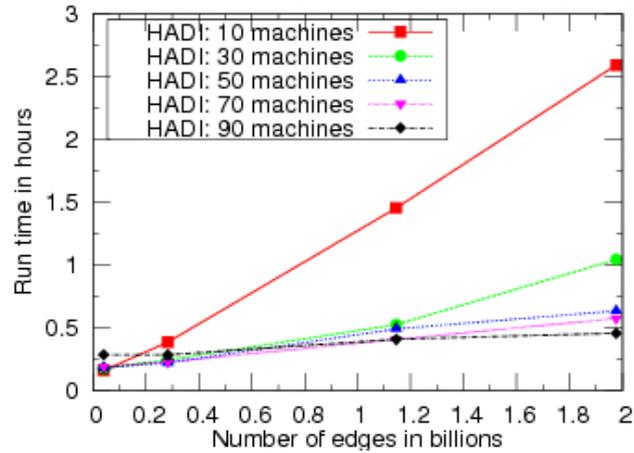


FIGURE 7.4: Running time versus number of edges with HADI-OPT on Kronecker graphs for three iterations. Notice the excellent scalability: linear on the graph size (number of edges).

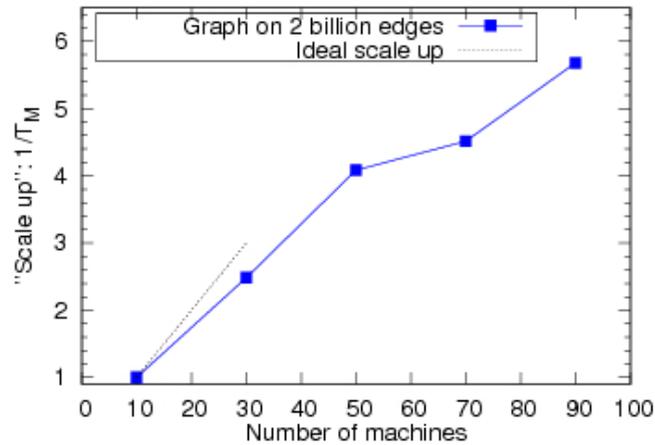


FIGURE 7.5: “Scale-up” (throughput $1/T_M$) versus number of machines M , for the Kronecker graph (2B edges). Notice the near-linear growth in the beginning, close to the ideal(dotted line).

7.5.2 Running Time and Scale-up

Figure 7.4 gives the wall-clock time of HADI-OPT versus the number of edges in the graph. Each curve corresponds to a different number of machines used (from 10 to 90). HADI has excellent scalability, with its running time being linear on the number of edges. The rest of the HADI versions (HADI-plain, HADI-BL, and HADI-BSE), were slower, but had a similar linear trend.

Figure 7.5 gives the throughput $1/T_M$ of HADI-OPT. We also tried HADI with one machine; however it didn’t complete, since the machine would take so long that it would often fail in the meanwhile. For this reason, we do not report the typical scale-up score

$s = T_1/T_M$ (ratio of time with 1 machine, over time with M machine), and instead we report just the inverse of T_M . HADI scales up near-linearly with the number of machines M , close to the ideal scale-up.

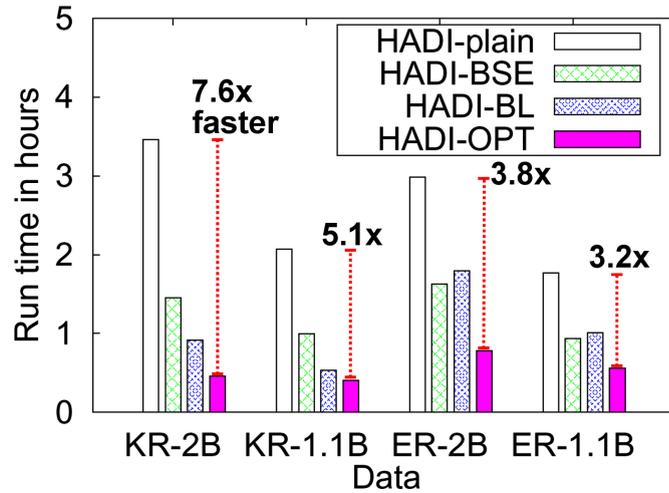


FIGURE 7.6: Run time of HADI with/without optimizations for Kronecker and Erdős-Rényi graphs with several billions of edges, on the M45 HADOOP cluster using 90 machines for 3 iterations. HADI-OPT is up to $7.6\times$ faster than HADI-plain.

7.5.3 Effect of Optimizations

Among the optimizations that we mentioned earlier, which one helps the most, and by how much? Figure 7.6 plots the running time of different graphs versus different HADI optimizations. For the Kronecker graphs, we see that block operation is more efficient than bit shuffle encoding. Here, HADI-OPT achieves $7.6\times$ better performance than HADI-plain. For the Erdős-Rényi graphs, however, we see that block operations do not help more than bit shuffle encoding, because the adjacency matrix has no block structure, while Kronecker graphs do. Also notice that HADI-BLK and HADI-OPT run faster on Kronecker graphs than on Erdős-Rényi graphs of the same size. Again, the reason is that Kronecker graphs have fewer nonzero blocks (i.e., “communities”) by their construction, and the “block” operation yields more savings.

7.6 Structure of Real-World Networks

HADI reveals new patterns in massive graphs which we present in this section. We distinguish these new patterns into *static* (Section 7.6.1) and *temporal* (Section 7.6.2).

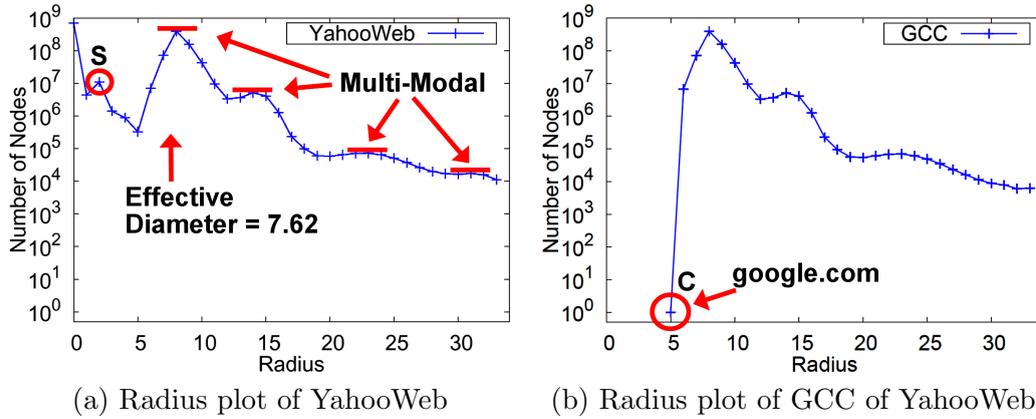


FIGURE 7.7: (a) Radius plot (Count versus Radius) of the YahooWeb graph. Notice the effective diameter is surprisingly small. Also notice the peak (marked 'S') at radius 2, due to star-structured disconnected components. (b) Radius plot of GCC (Giant Connected Component) of YahooWeb graph. The *only* vertex with radius 5 (marked 'C') is `google.com`.

7.6.1 Static Patterns

7.6.1.1 Diameter

What is the diameter of the Web? Albert et al. [22] computed the diameter on a directed Web graph with approximately 0.3 million vertices and conjectured that it should be around 19 for a 1.4 billion-vertex Web graph as shown in the upper line of Figure 7.8. Broder et al. [90] used their sampling approach from approximately 200 million-vertices and reported 16.15 and 6.83 as the diameter for the directed and the undirected cases, respectively. What should the effective diameter be, for a significantly larger crawl of the Web, with billions of vertices? Figure 7.7 gives the surprising answer:

Observation 2 (Small Web). The effective diameter of the YahooWeb graph (year: 2002) is surprisingly small, between 7 and 8.

The previous results from Albert et al. [22] and Broder et al. [90] also consider the undirected version of the Web graph. We compute the average diameter and show the comparison of diameters of different graphs in Figure 7.8. We first observe that the average diameters of all graphs are relatively small (< 20) for both the directed and the undirected cases. We also observe that the Albert et al.'s conjecture for the diameter of the directed graph is over-pessimistic: both the sampling approach and HADI output smaller values for the diameter of the directed graph. For the diameter of the undirected graph, we observe the constant/shrinking diameter pattern [280].

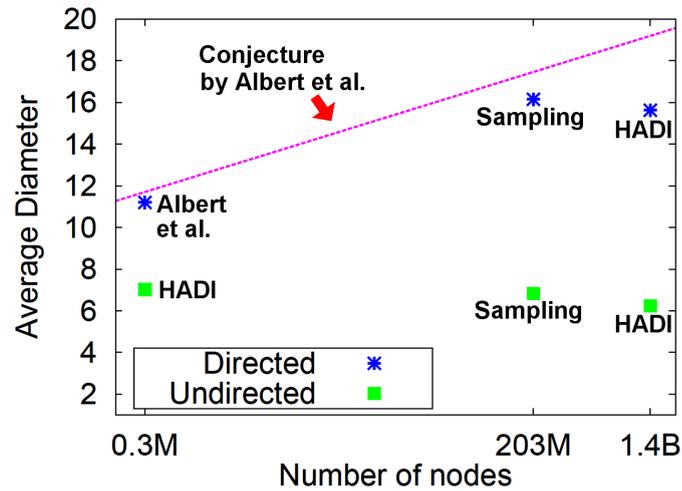


FIGURE 7.8: Average diameter vs. number of vertices in lin-log scale for the three different Web graphs, where M and B stand for millions and billions respectively. (0.3M): web pages inside nd.edu at 1999, from Albert et al.'s work. (203M): web pages crawled by Altavista at 1999, from Broder et al.'s work (1.4B): web pages crawled by Yahoo at 2002 (YahooWeb in Table 7.1). Notice the relatively small diameters for both the directed and the undirected cases.

7.6.1.2 Shape of Distribution

Figure 7.7 shows that the radii distribution in the Web Graph is multimodal. In other relatively smaller networks, we observe a bimodal structure. As shown in the Radius plot of U.S. Patent and LinkedIn network in Figure 7.9, they have a peak at zero, a dip at a small radius value (9, and 4, respectively) and another peak very close to the dip.

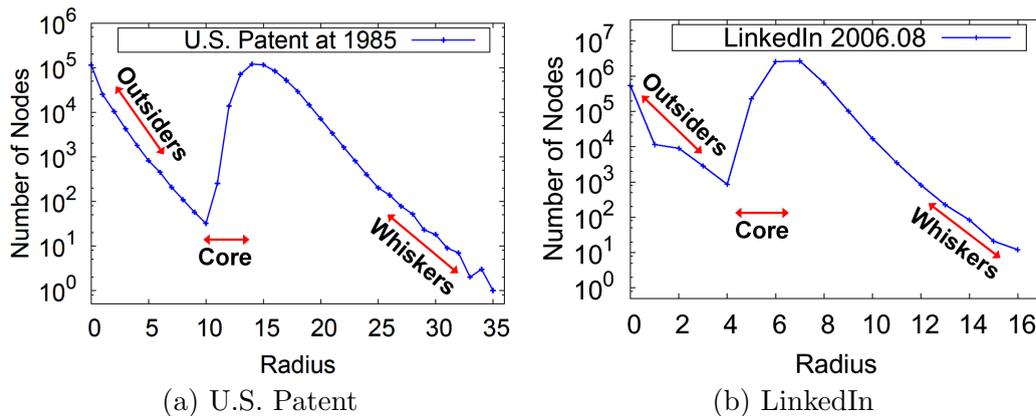


FIGURE 7.9: Static Radius Plot (Count versus Radius) of U.S. Patent and LinkedIn graphs. Notice the bimodal structure with 'outsiders' (vertices in the DCs), 'core' (central vertices in the GCC), and 'whiskers' (vertices connected to the GCC with long paths).

Observation 3 (Multi-modal and Bi-modal). The radius distribution of the Web graph has a multimodal structure. Smaller networks have a bimodal structure.

A natural question to ask with respect to the bimodal structure is what are the common properties of the vertices that belong to the first peak; similarly, for the vertices in the first dip, and the same for the vertices of the second peak. After investigation, the former are vertices that belong to disconnected components (DCs); vertices in the dip are usually core vertices in the giant connected component (GCC), and the vertices at the second peak are the vast majority of well connected vertices in the GCC. Figure 7.10 exactly shows the radii distribution for the vertices of the GCC (in blue), and the vertices of the few largest remaining components.

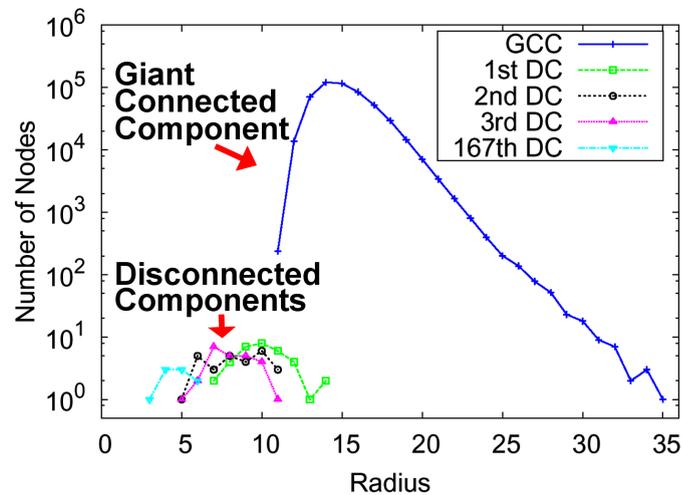


FIGURE 7.10: Radius plot (Count versus radius) for several connected components of the U.S. Patent data in 1985. In blue: the distribution for the giant connected component; rest colors: several disconnected components.

In Figure 7.10, we clearly see that the second peak of the bimodal structure came from the giant connected component. But, where does the first peak around radius 0 come from? We can get the answer from the distribution of connected component of the same graph in Figure 7.11. Since the ranges of radius are limited by the size of connected components, we see the first peak of Radius plot came from the disconnected components whose size follows a power law.

Now we can explain the three important areas of Figure 7.9: ‘outsiders’ are the vertices in the disconnected components, and responsible for the first peak and the negative slope to the dip. ‘Core’ are the central vertices with the smallest radii from the giant connected component. ‘Whiskers’ [277] are the vertices connected to the GCC with long paths.

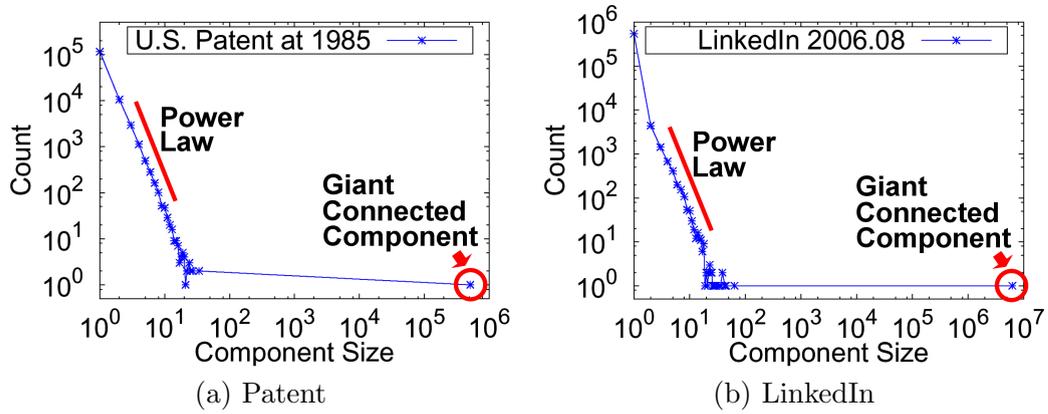


FIGURE 7.11: Size distribution of connected components. Notice the size of the disconnected components (DCs) follows a power-law which explains the first peak around radius 0 of the radius plots in Figure 7.9.

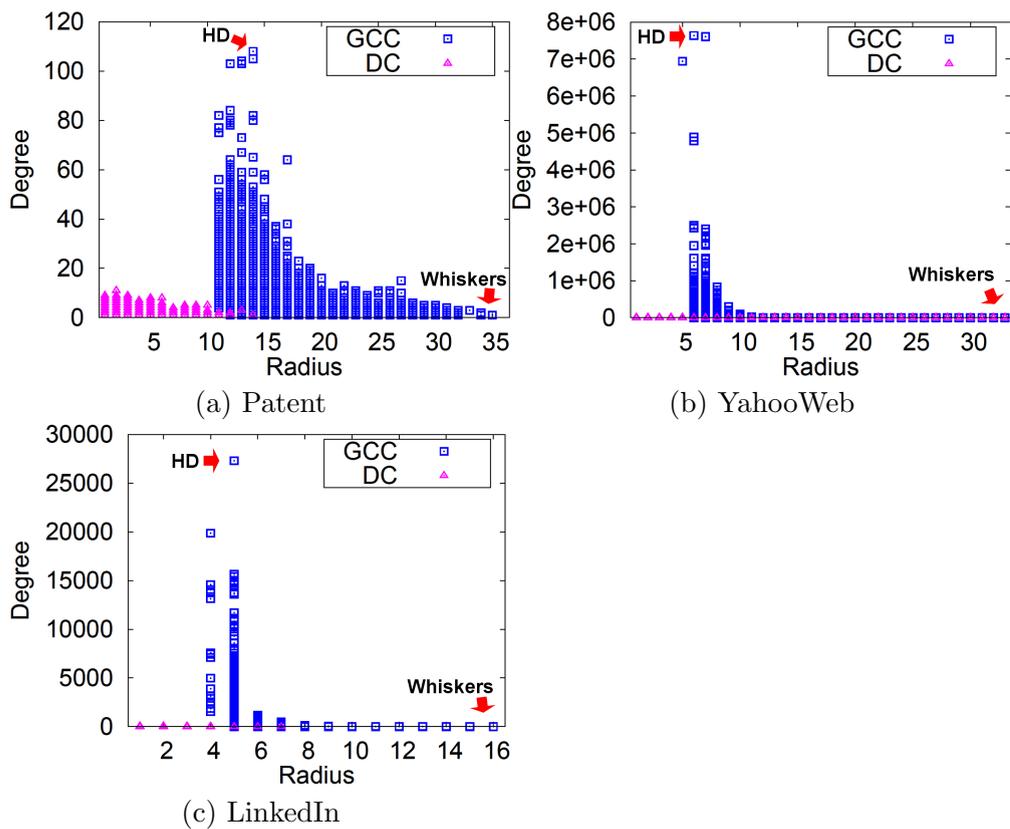


FIGURE 7.12: Radius-Degree plots of real-world graphs. HD represents the vertex with the highest degree. Notice that HD belongs to core vertices inside the GCC, and whiskers have small degree.

7.6.1.3 Radius plot of GCC

Figure 7.7(b) shows that all vertices of the GCC of the YahooWeb graph have radius 6 or more except for `google.com` that has radius one.

7.6.1.4 “Core” and “Whisker” vertices

Figure 7.12 shows the Radius-Degree plot of Patent, YahooWeb and LinkedIn graphs. The Radius-Degree plot is a scatterplot with one dot for every vertex plotting the degree of the vertex versus its radius. The points corresponding to vertices in the GCC are colored with blue, while the rest is in magenta. We observe that the highest degree vertices belong to the set of core vertices inside the GCC but are not necessarily the ones with the smallest radius. Finally, the whisker vertices have small degree and belong to chain subgraphs.

7.6.2 Temporal Patterns

Here we study the radius distribution as a function of time. We know that the diameter of a graph typically grows with time, spikes at the ‘gelling point’, and then shrinks [280, 301]. Indeed, this holds for our datasets as shown in Figure 7.13.

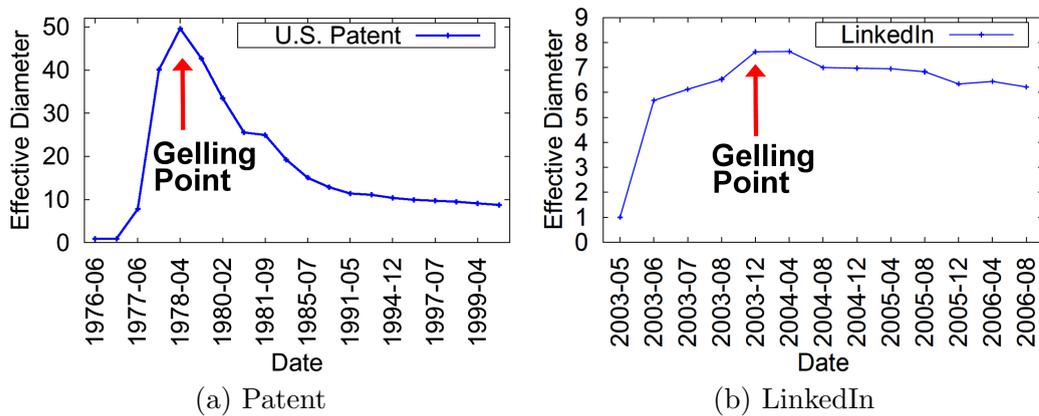


FIGURE 7.13: Evolution of the effective diameter of real graphs. The diameter increases until a ‘gelling’ point, and starts to decrease after the point.

Figure 7.14 shows our findings. The radius distribution expands to the right until it reaches the gelling point. Then, it contracts to the left. Finally, the decreasing segments of several, real radius plots seem to decay exponentially, that is

$$\text{count}(r) \propto \exp(-cr) \quad (7.5)$$

for every time tick *after* the gelling point. $count(r)$ is the number of vertices with radius r , and c is a constant. For the Patent and LinkedIn graphs, the absolute correlation coefficient was at least 0.98.

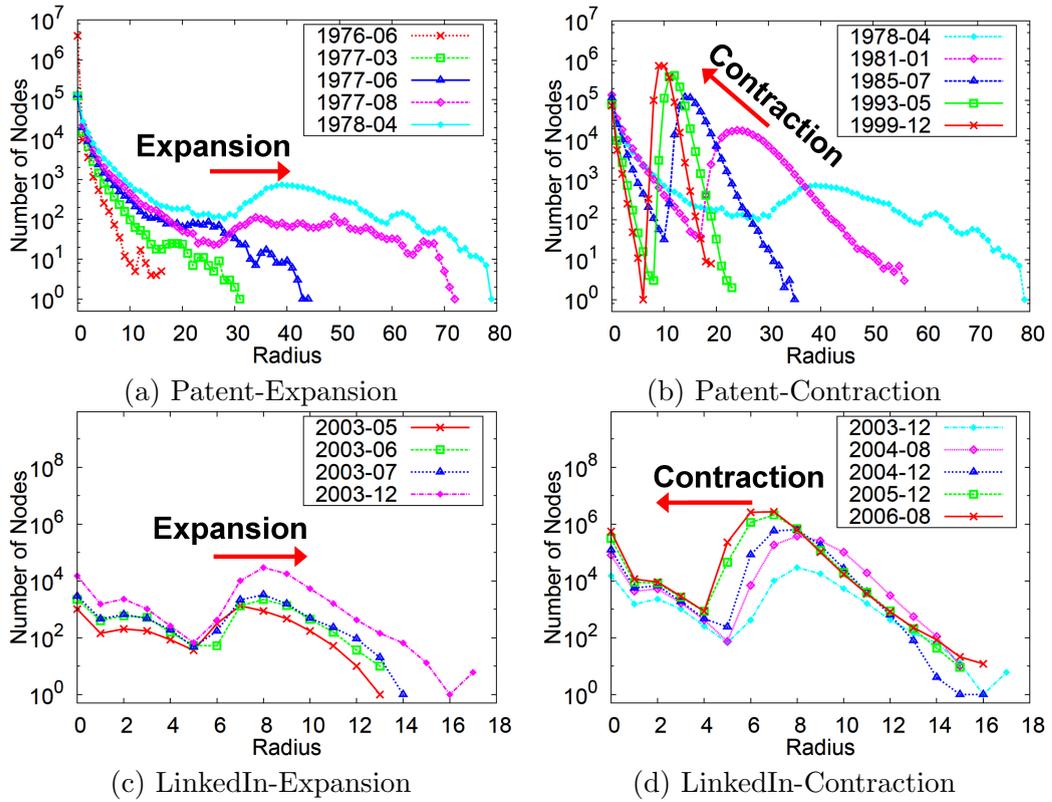


FIGURE 7.14: Radius distribution over time. “Expansion”: the radius distribution moves to the right until the gelling point. “Contraction”: the radius distribution moves to the left after the gelling point.

Chapter 8

FENNEL: Streaming Graph Partitioning for Massive Scale Graphs

8.1 Introduction

Big volumes of data are typically managed and analyzed on large distributed systems [130]. Specifically, the data is partitioned across a large number of cheap, commodity machines which are typically connected by gigabit Ethernet. Minimizing the communication between the machines is critical, since the communication cost often dominates computation cost. A key problem towards minimizing communication cost for big graph data is the *balanced graph partitioning* (BGPA) problem: partition the vertex set of the graph in a given number of machines in such a way that each partition is balanced and the number edges cut is minimized. The *balanced graph partitioning* problem is a classic NP-hard problem [258] for which several approximation algorithms have been designed.

# Clusters (k)	Fennel		Best competitor		Hash Partition		METIS	
	λ	ρ	λ	ρ	λ	ρ	λ	ρ
2	6.8%	1.1	34.3%	1.04	50%	1	11.98%	1.02
4	29%	1.1	55.0%	1.07	75%	1	24.39%	1.03
8	48%	1.1	66.4%	1.10	87.5%	1	35.96%	1.03

TABLE 8.1: Fraction of edges cut λ and the normalized maximum load ρ for Fennel, the previously best-known heuristic (linear weighted degrees [373]) and hash partitioning of vertices for the Twitter graph with approximately 1.5 billion edges. Fennel and best competitor require around 40 minutes, METIS more than $8\frac{1}{2}$ hours.

In practice, systems aim at providing good partitions in order to enhance their performance, e.g., [295, 346]. It is worth emphasizing that the balanced graph partitioning problem appears in various guises in numerous domains [244].

Another major challenge in the area of big graph data is efficient processing of dynamic graphs. For example, new accounts are created and deleted every day in online services such as Facebook, Skype and Twitter. Furthermore, graphs created upon post-processing datasets such as Twitter posts are also dynamic, see for instance [36]. It is crucial to have efficient graph partitioners of dynamic graphs. For example, in the Skype service, each time a user logs in, his/her online contacts get notified. It is expensive when messages have to be sent across different graph partitions since this would typically involve using network infrastructure. The balanced graph partitioning problem in the dynamic setting is known as *streaming graph partitioning* [373]. Vertices (or edges) arrive and the decision of the placement of each vertex (edge) has to be done “on-the-fly” in order to incur as little computational overhead as possible.

It is worth noting that the state-of-the-art work on graph partitioning seems to roughly divide in two main lines of research. Rigorous mathematically work and algorithms that do not scale to massive graphs, e.g., [260], and heuristics that are used in practice [244, 245, 344, 373]. Our work contributes towards bridging the gap between theory and practice.

The remainder of the Chapter is organized as follows. In Section 8.2, we introduce our graph partitioning framework and present our main theoretical result. In Section 8.3, we present our scalable, streaming algorithm. In Section 8.4, we evaluate our method versus the state-of-the-art work on a broad set of real-world and synthetic graphs, while in Section 8.5 we provide our experimental results in the Apache Giraph.

8.2 Proposed Framework

Notation. Throughout this Chapter we use the following notation. Let $G(V, E)$ be a simple, undirected graph. Let the number of vertices and edges be denoted as $|V| = n$ and $|E| = m$. For a subset of vertices $S \subseteq V$, let $e(S, S)$ be the set of edges with both end vertices in the set S , and let $e(S, V \setminus S)$ be the set of edges whose one end-vertex is in the set S and the other is not. For a given vertex v let $t_S v$ be the number of triangles (v, w, z) such that $w, z \in S$. We define a *partition* of vertices $\mathcal{P} = (S_1, \dots, S_k)$ to be a family of pairwise disjoint sets vertices, i.e., $S_i \subseteq V$, $S_i \cap S_j = \emptyset$ for every $i \neq j$. We call S_i to be a cluster of vertices. Finally, for a graph $G = (V, E)$ and a partition

$\mathcal{P} = (S_1, S_2, \dots, S_k)$ of the vertex set V , let $\partial e(\mathcal{P})$ be the set of edges that cross partition boundaries, i.e. $\partial e(\mathcal{P}) = \cup_{i=1}^k e(S_i, V \setminus S_i)$.

Graph Partitioning Framework. We formulate a graph partitioning framework that is based on accounting for the cost of internal edges and the cost of edges cut by a partition of vertices in a single global objective function.

The size of individual partitions. We denote with $\sigma(S_i)$ the size of the cluster of vertices S_i , where σ is a mapping to the set of real numbers. Special instances of interest are (1) *edge cardinality* where the size of the cluster i is proportional to the total number of edges with at least one end-vertex in the set S_i , i.e. $|e(S_i, S_i)| + |e(S_i, V \setminus S_i)|$, (2) *interior-edge cardinality* where the size of cluster i is proportional to the number of internal edges $|e(S_i, S_i)|$, and (3) *vertex cardinality* where the size of partition i is proportional to the total number of vertices $|S_i|$. The edge cardinality of a cluster is an intuitive measure of cluster size. This is of interest for computational tasks over input graph data where the computational complexity within a cluster of vertices is linear in the number of edges with at least one vertex in the given cluster. For example, this is the case for iterative computations such as solving the power iteration method. The vertex cardinality is a standard measure of the size of a cluster and for some graphs may serve as a proxy for the edge cardinality, e.g. for the graphs with bounded degrees.

The global objective function. We define a global objective function that consists of two elements: (1) the inter-partition cost $c_{\text{OUT}} : N^k \rightarrow \mathbb{R}_+$ and (2) the intra-partition cost $c_{\text{IN}} : N^k \rightarrow \mathbb{R}_+$. These functions are assumed to be increasing and super-modular (or convex, if extended to the set of real numbers). For every given partition of vertices $\mathcal{P} = (S_1, S_2, \dots, S_k)$, we define the global cost function as

$$f(\mathcal{P}) = c_{\text{OUT}}(|e(S_1, V \setminus S_1)|, \dots, |e(S_k, V \setminus S_k)|) \\ + c_{\text{IN}}(\sigma(S_1), \dots, \sigma(S_k)).$$

It is worth mentioning some particular cases of interest. Special instance of interest for the inter-partition cost is the linear function in the total number of cut edges $|\partial e(\mathcal{P})|$. This case is of interest in cases where an identical cost is incurred per each edge cut, e.g. in cases where messages are exchanged along cut edges and these messages are transmitted through some common network bottleneck. For the intra-partition cost, a typical goal is to balance the cost across different partitions and this case is accommodated by defining $c_{\text{IN}}(\sigma(S_1), \dots, \sigma(S_k)) = \sum_{i=1}^k c(\sigma(S_i))$, where $c(x)$ is a convex increasing function such that $c(0) = 0$. In this case, the intra-partition cost function, being defined as a sum of convex functions of individual cluster sizes, would tend to balance the cluster sizes, since the minimum is attained when sizes are equal.

We formulate the graph partitioning problem as follows.

Optimal k -Graph Partitioning

Given a graph $G = (V, E)$, find a partition $\mathcal{P}^* = \{S_1^*, \dots, S_k^*\}$ of the vertex set V , such that $f(\mathcal{P}^*) \geq f(\mathcal{P})$, for all partitions \mathcal{P} such that $|\mathcal{P}| = k$.

We refer to the partition \mathcal{P}^* as the optimal k graph partition of the graph G .

Streaming setting. The streaming graph partitioning problem can be defined as follows. Let $G = (V, E)$ be an input graph and let us assume that we want to partition the graph into k disjoint subsets of vertices. The vertices arrive in some order, each one with the set of its neighbors. We consider three different stream orders, as in [373].

- Random: Vertices arrive according to a random permutation.
- BFS: This ordering is generated by selecting a vertex uniformly at random and performing breadth first search starting from that vertex.
- DFS: This ordering is identical to the BFS ordering, except that we perform depth first search.

A k -partitioning streaming algorithm has to decide whenever a new vertex arrives to which cluster it is going to be placed. A vertex is never moved after it has been assigned to a cluster. The formal statement of the problem follows.

Classic Balanced Graph Partitioning. We consider the traditional instance of a graph partitioning problem that is a special case of our framework by defining the inter-partition cost to be equal to the total number of edges cut and the intra-partition cost defined in terms of the vertex cardinalities.

The starting point in the existing literature, e.g., [258, 260], is to admit hard cardinality constraints, so that $|S_i^*| \leq \nu \frac{n}{k}$ for $i = 1, \dots, k$, where $\nu \geq 1$ is a fixed constant. This set of constraints makes the problem significantly hard. Currently, state-of-the-art work depends on the impressive ARV barrier [37] which results in a $O(\sqrt{\log n})$ approximation factor. The typical formulation is the following:

$$\begin{aligned} & \text{minimize}_{\mathcal{P}=(S_1, \dots, S_k)} && |\partial e(\mathcal{P})| \\ & \text{subject to} && |S_i| \leq \nu \frac{n}{k}, \forall i \in \{1, \dots, k\} \end{aligned}$$

Our approach: Just Relax! The idea behind our approach is to *relax* the hard cardinality constraints by introducing a term in the objective $c_{\text{IN}}(\mathcal{P})$ whose minimum is achieved when $|S_i| = \frac{n}{k}$ for all $i \in \{1, \dots, k\}$. Therefore, our framework is based on a well-defined global graph partitioning objective function, which allows for a principled design of approximation algorithms and heuristics as shall be demonstrated in Section 8.3. Our graph partitioning method is based on solving the following optimization problem:

$$\boxed{\text{minimize}_{\mathcal{P}=(S_1, \dots, S_k)} \quad |\partial e(\mathcal{P})| + c_{\text{IN}}(\mathcal{P})} \quad (8.1)$$

Intra-partition cost: With the goal in mind to favor balanced partitions, we may define the intra-partition cost function by $c_{\text{IN}}(\mathcal{P}) = \sum_{i=1}^k c(|S_i|)$ where $c(x)$ is an increasing function chosen to be *super-modular*, so that the following increasing returns property holds $c(x+1) - c(x) \geq c(y+1) - c(y)$, for every $0 \leq y \leq x$.

We shall focus our attention to the following family of functions $c(x) = \alpha x^\gamma$, for $\alpha > 0$ and $\gamma \geq 1$. By the choice of the parameter γ , this family of cost functions allows us to control how much the imbalance of cluster sizes is accounted for in the objective function. In one extreme case where $\gamma = 1$, we observe that the objective corresponds to minimizing the number of cut-edges, thus entirely ignoring any possible imbalance of the cluster sizes. On the other hand, by taking larger values for the parameter γ , the more weight is put on the cost of partition imbalance, and this cost may be seen to approximate hard constraints on the imbalance in the limit of large γ . Parameter α is also important. We advocate a principled choice of α independently of whether it is suboptimal compared to other choices. Specifically, we choose $\alpha = m \frac{k^{\gamma-1}}{n^\gamma}$. This provides us a proper scaling, since for this specific choice of α , our optimization problem is equivalent to minimizing a natural normalization of the objective function $\frac{\sum_{i=1}^k e(S_i, V \setminus S_i)}{m} + \frac{1}{k} \sum_{i=1}^k \left(\frac{|S_i|}{\frac{n}{k}} \right)^\gamma$.

An equivalent maximization problem. We note that the optimal k graph partitioning problem admits an equivalent formulation as a maximization problem. It is of interest to consider this alternative formulation as it allows us to make a connection with the concept of graph modularity, which we do later in this section. For a graph $G = (V, E)$ and $S \subseteq V$, we define the function $h : 2^V \rightarrow \mathbb{R}$ as:

$$h(S) = |e(S, V \setminus S)| - c(|S|)$$

where $h(\emptyset) = h(\{v\}) = 0$ for every $v \in V$. Given $k \geq 1$ and a partition $\mathcal{P} = \{S_1, \dots, S_k\}$ of the vertex set V , we define the function g as

$$g(\mathcal{P}) = \sum_{i=1}^k h(S_i).$$

Now, we observe that maximizing the function $g(\mathcal{P})$ over all possible partitions \mathcal{P} of the vertex set V such that $|\mathcal{P}| = k$ corresponds to the k graph partitioning problem. Indeed, this follows by noting that

$$\begin{aligned} g(\mathcal{P}) &= \sum_{i=1}^k |e(S_i, S_i)| - c(|S_i|) \\ &= (m - \sum_{i=1}^k |e(S_i, V \setminus S_i)|) - c(|S_i|) \\ &= m - f(\mathcal{P}). \end{aligned}$$

Thus, maximizing function $g(\mathcal{P})$ corresponds to minimizing function $f(\mathcal{P})$, which is precisely the objective of our k graph partitioning problem.

Modularity: We note that when the function $c(x)$ is taken from the family $c(x) = \alpha x^\gamma$, for $\alpha > 0$ and $\gamma = 2$, our objective has a combinatorial interpretation. Specifically, our problem is equivalent to maximizing the function

$$\sum_{i=1}^k [|e(S_i, S_i)| - p \binom{|S_i|}{2}]$$

where $p = \alpha/2$. In this case, each summation element admits the following interpretation: it corresponds to the difference between the realized number of edges within a cluster and the expected number of edges within the cluster under the null-hypothesis that the graph is an Erdős-Rényi random graph with parameter p . This is intimately related to the concepts of graph modularity [196, 314, 316] and quasi-cliques [399]. Recall that an approximation algorithm for a maximization problem is meaningful as a notion if the optimum solution is positive. However, in our setting our function g does not result as it can easily be seen in a non-negative optimum. For instance, if G is the empty graph on n vertices, any partition of G in k parts results in a negative objective (except when $k = n$ when the objective becomes 0). Therefore, we need to shift our objective in order to come up with a multiplicative approximation algorithm. We define the following shifted objective, along the lines of Chapter 6.

Definition 8.1. $k \geq 1$. Also let $\mathcal{P}^* = \{S_1^*, \dots, S_k^*\}$ be a partition of the vertex set V . We define the function g as

$$\tilde{g}(\mathcal{P}) = \alpha \binom{n}{2} + \sum_{i=1}^k f(S_i).$$

Claim 5. For any partition \mathcal{P} , $\tilde{g}(\mathcal{P}) \geq 0$.

Proof. The proof follows directly from the fact that for any positive-valued s_1, s_2, \dots, s_k such that $\sum_{i=1}^k s_i = n$, the following holds $n^2 \geq s_1^2 + \dots + s_k^2$. \square

Due to the combinatorial interpretation of the objective function, we design a semidefinite programming approximation algorithm. Before that, we see how random partitioning performs.

Random Partitioning: Suppose each vertex is assigned to one of k partitions uniformly at random. This simple graph partition is a faithful approximation of hash partitioning of vertices that is commonly used in practice. In expectation, each of the k clusters will have $\frac{n}{k}$ vertices. Let S_1, \dots, S_k be the resulting k clusters. How well does this simple algorithm perform with respect to our objective? Let \mathcal{P}^* be an optimal partition for the optimal quasi-clique problem, i.e. $g(\mathcal{P}^*) \geq g(\mathcal{P})$, for every partition \mathcal{P} of the vertex set V into k partitions. Notice that \mathcal{P}^* is also an optimal partition for the optimal quasi-clique problem with shifted objective function. Now, note that an edge $e = (u, v)$ has probability $\frac{1}{k}$ that both its endpoints belong to the same cluster. By the linearity of expectation, we obtain by simple calculations:

$$\begin{aligned} \mathbb{E}[\tilde{g}(S_1, \dots, S_k)] &= \frac{|E|}{k} + \alpha \frac{k-1}{k} \binom{n}{2} \\ &\geq \frac{1}{k} \left(|E| + \alpha \binom{n}{2} \right) \\ &\geq \frac{1}{k} \tilde{g}(\mathcal{P}^*) \end{aligned}$$

where last inequality comes from the simple upper bound $\tilde{g}(\mathcal{P}) \leq |E| + \alpha \binom{n}{2}$ for any partition \mathcal{P} .

An SDP Rounding Algorithm

We define a vector variable x_i for each vertex $i \in V$ and we allow x_i to be one of the unit vectors e_1, \dots, e_k , where e_j has only the j -th coordinate 1.

$$\begin{aligned}
 & \text{maximize} && \sum_{e=(i,j)} x_i x_j + \alpha \sum_{i<j} (1 - x_i x_j) \\
 & \text{subject to} && x_i \in \{e_1, \dots, e_k\}, \forall i \in \{1, \dots, n\}
 \end{aligned}
 \tag{8.2}$$

We obtain the following semidefinite programming relaxation:

$$\begin{aligned}
 & \text{maximize} && \sum_{e=(i,j)} y_{ij} + \alpha \sum_{i<j} (1 - y_{ij}) \\
 & \text{subject to} && y_{ii} = 1, \forall i \in \{1, \dots, n\} \\
 & && y_{ij} \geq 0, \forall i \neq j \\
 & && Y \succeq 0, Y \text{ symmetric}
 \end{aligned}
 \tag{8.3}$$

The above SDP can be solved within an additive error of δ of the optimum in time polynomial in the size of the input and $\log(\frac{1}{\delta})$ by interior point algorithms or the ellipsoid method [24]. In what follows, we refer to the optimal value of the integer program as OPT_{IP} and of the semidefinite program as OPT_{SDP} . Our algorithm is the following:

SDP-Relax

- *Relaxation:* Solve the semidefinite program [24] and compute a Cholesky decomposition of Y . Let v_0, v_1, \dots, v_n be the resulting vectors.
- *Randomized Rounding:* Randomly choose $t = \lceil \log k \rceil$ unit length vectors $r_i \in \mathbb{R}^k, i = 1, \dots, t$. These t random vectors define $2^t = k$ possible regions in which the vectors v_i can fall: one region for each distinct possibility of whether $r_j v_i \geq 0$ or $r_j v_i < 0$. Define a cluster by adding all vertices whose vector v_i fall in a given region.

Theorem 8.2. *Algorithm SDP-Relax is a $\Omega(\frac{\log k}{k})$ approximation algorithm for the Optimal Quasi-Clique Partitioning.*

Proof. Let C_k be the score of the partition produced by our randomized rounding. Define $A_{i,j}$ to be the event that vertices i and j are assigned to the same partition. Then,

$$\mathbb{E}[C_k] = \sum_{e=(i,j)} \Pr[A_{i,j}] + \alpha \sum_{i<j} (1 - \Pr[A_{i,j}])$$

As in Goemans-Williamson [198], given a random hyperplane with normal vector r that goes through the origin, the probability of $\text{sgn}(v_i^T r) = \text{sgn}(v_j^T r)$, i.e., i and j fall on the same side of the hyperplane, is $1 - \frac{\arccos(v_i^T v_j)}{\pi}$. Since we have t independent hyperplanes

$$\Pr[A_{i,j}] = \left(1 - \frac{\arccos(v_i^T v_j)}{\pi}\right)^t.$$

Let us define, for $t \geq 1$,

$$f_1(\theta) = \frac{\left(1 - \frac{\theta}{\pi}\right)^t}{\cos(\theta)}, \quad \theta \in [0, \frac{\pi}{2})$$

and

$$\rho_1 = \min_{0 \leq \theta < \frac{\pi}{2}} f_1(\theta).$$

Similarly, define for $t \geq 1$,

$$f_2(\theta) = \frac{1 - \left(1 - \frac{\theta}{\pi}\right)^t}{1 - \cos(\theta)}, \quad \theta \in [0, \frac{\pi}{2})$$

and

$$\rho_2 = \min_{0 \leq \theta < \frac{\pi}{2}} f_2(\theta).$$

We wish to find a ρ such that $\mathbb{E}[C_k] \geq \rho \text{OPT}_{\text{SDP}}$. Since, $\text{OPT}_{\text{SDP}} \geq \text{OPT}_{\text{IP}}$, this would then imply $\mathbb{E}[C_k] \geq \rho \text{OPT}_{\text{IP}}$.

We note that

$$f_1'(\theta) = \frac{-\frac{t}{\pi} \left(1 - \frac{\theta}{\pi}\right)^{t-1} \cos(\theta) + \left(1 - \frac{\theta}{\pi}\right)^t \sin(\theta)}{\cos^2(\theta)}.$$

It follows that $f_1'(\theta(t)) = 0$ is equivalent to

$$t = (\pi - \theta(t)) \tan(\theta(t)).$$

Notice that the following two hold

$$\lim_{t \rightarrow \infty} \theta(t) = \frac{\pi}{2}$$

and

$$\frac{\pi}{2} \tan(\theta(t)) \leq t \leq \pi \tan(\theta(t)). \quad (8.4)$$

We next show that $f_1(\theta(t)) \geq \frac{1}{\pi}t2^{-t}$, by the the following series of relations

$$\begin{aligned} f_1(\theta(t)) &= \frac{\left(1 - \frac{\theta(t)}{\pi}\right)^t}{\cos(\theta(t))} = \sqrt{1 + \tan^2(\theta(t))} \left(1 - \frac{\theta(t)}{\pi}\right)^t \\ &\geq \sqrt{1 + \tan^2(\theta(t))} 2^{-t} \geq \sqrt{1 + \left(\frac{t}{\pi}\right)^2} 2^{-t} \\ &\geq \frac{1}{\pi}t2^{-t} \end{aligned}$$

where the second equality is by the fact $\cos(\theta) = \frac{1}{\sqrt{1 + \tan^2(\theta)}}$, the first inequality is by the fact $\theta(t) \leq \frac{\pi}{2}$, and the second inequality is by (8.4).

Thus, for $t = \log_2(k)$, we conclude

$$\rho_1 \geq \frac{1}{\pi \log(2)} \frac{\log(k)}{k}.$$

Now, we show that $\rho_2 = \frac{1}{2}$. First we show that for any $t \geq 1$, $f_2(\theta) \geq 1/2$. To this end, we note

$$\frac{1 - \left(1 - \frac{\theta}{\pi}\right)^t}{1 - \cos(\theta)} \geq \frac{1}{2} \Leftrightarrow \frac{1}{2} \left(1 + \cos(\theta)\right) \geq \left(1 - \frac{\theta}{\pi}\right)^t.$$

Notice that for all $\theta \in [0, \pi/2)$, if $t_1 \geq t_2 \geq 1$, then $\left(1 - \frac{\theta}{\pi}\right)^{t_1} \leq \left(1 - \frac{\theta}{\pi}\right)^{t_2}$. Hence, it suffices $\frac{1}{2} \left(1 + \cos(\theta)\right) \geq \left(1 - \frac{\theta}{\pi}\right)$. With the use of simple calculus, the latter is true and is also tight for $\theta = 0$ and $\theta = \pi/2$. It is worth observing the opposite trend of the values of ρ_1 and ρ_2 . The reason is that $\Pr[A_{i,j}]$ drops as we use more hyperplanes and, of course, $1 - \Pr[A_{i,j}]$ grows. Now, we can establish the following lower bound on the expected score of our randomized rounding procedure. Let $\theta_{i,j} = \arccos(v_i^T v_j)$.

$$\begin{aligned} \mathbb{E}[C_k] &= \sum_{e=(i,j)} \Pr[A_{i,j}] + \alpha \sum_{i < j} \left(1 - \Pr[A_{i,j}]\right) \\ &= \sum_{e=(i,j)} \left(1 - \frac{\theta_{i,j}}{\pi}\right)^t + \alpha \sum_{i < j} \left(1 - \left(1 - \frac{\theta_{i,j}}{\pi}\right)^t\right) \\ &\geq \sum_{e=(i,j)} \rho_1 \cos(\theta_{i,j}) + \alpha \sum_{i < j} \rho_2 \left(1 - \cos(\theta_{i,j})\right) \\ &\geq \min\{\rho_1, \rho_2\} \text{OPT}_{\text{SDP}} \\ &\geq \min\{\rho_1, \rho_2\} \text{OPT}_{\text{IP}} \end{aligned}$$

It suffices to set $\rho = \min\{\rho_1, \rho_2\}$. The above analysis shows that our algorithm is a ρ -approximation algorithm, where $\rho = \Omega(\frac{\log(k)}{k})$.

□

The approximation guarantees that hold for the shifted objective function have the following meaning for the original objective function. Suppose that \mathcal{P} is a ρ -approximation with respect to the shifted-objective quasi-clique partitioning problem, i.e. $\tilde{g}(\mathcal{P}) \geq \rho \tilde{g}(\mathcal{P}^*)$. Then, it holds

$$g(\mathcal{P}) \geq \rho g(\mathcal{P}^*) - (1 - \rho)\alpha \binom{n}{2}.$$

Notice that this condition is equivalent to

$$g(\mathcal{P}) - g(\mathcal{P}^*) \geq -(1 - \rho)[g(\mathcal{P}^*) + \alpha \binom{n}{2}].$$

We wish to outline that this approximation guarantee is near to a ρ -approximation if the parameter α is small enough so that the term $g(\mathcal{P}^*)$ dominates the term $\alpha \binom{n}{2}$. For example, for an input graph that consists of k cliques, we have that $g(\mathcal{P}^*) \geq \alpha \binom{n}{2}$ corresponds to $\frac{\alpha}{1-\alpha} \leq \frac{1}{k} \frac{n+k}{n+1}$, from which we conclude that it suffices that $\alpha \leq \frac{1}{k+1}$.

Alternative Roundings

One may ask whether the relaxation of Frieze and Jerrum [181], Karger, Motwani and Sudan [243] can improve significantly the approximation factor. We provide negative evidence. Before we go into further details, we notice that the main “bottleneck” in our approximation is the probability of two vertices being in the same cluster, as k grows. The probability that i, j are in the same cluster in our rounding is $p(\theta)$ where

$$p(\theta) = \left(1 - \frac{\theta}{\pi}\right)^{\log k}.$$

Suppose $\theta = \frac{\pi}{2}(1 - \epsilon)$. Then,

$$\begin{aligned} p(\theta) &= \left(1 - \frac{\theta}{\pi}\right)^{\log k} = \left(\frac{1 + \epsilon}{2}\right)^{\log k} \\ &= \frac{1}{k} + \epsilon \frac{\log k}{k} + O(\epsilon^2). \end{aligned}$$

As we see from Lemma 5 in [181], for this θ , the asymptotic expression matches ours:

$$\begin{aligned} N_k(\cos(\theta)) &\approx \frac{1}{k} + \frac{2 \log k}{k} \cos\left(\frac{\pi}{2}(1 - \epsilon)\right) \\ &= \frac{1}{k} + \pi \epsilon \frac{\log k}{k} + O(\epsilon^2). \end{aligned}$$

8.3 One-Pass Streaming Algorithm

We derive a streaming algorithm by using a greedy assignment of vertices to partitions as follows: assign each arriving vertex to a partition such that the objective function of the k graph partitioning problem, defined as a maximization problem, is increased the most. Formally, given that current vertex partition is $\mathcal{P} = (S_1, S_2, \dots, S_k)$, a vertex v is assigned to partition i such that

$$\begin{aligned} & g(S_1, \dots, S_i \cup \{v\}, \dots, S_j, \dots, S_k) \\ & \geq g(S_1, \dots, S_i, \dots, S_j \cup \{v\}, \dots, S_k), \text{ for all } j \in [k]. \end{aligned}$$

Defining $\delta g(v, S_i) = g(S_1, \dots, S_i \cup \{v\}, \dots, S_j, \dots, S_k) - g(S_1, \dots, S_i, \dots, S_j, \dots, S_k)$, the above greedy assignment of vertices corresponds to that in the following algorithm.

Greedy vertex assignment

- Assign vertex v to partition i such that $\delta g(v, S_i) \geq \delta g(v, S_j)$, for all $j \in [k]$

Special case: edge-cut and balanced vertex cardinality. This is a special case of introduced that we discussed in Section 8.2. In this case, $\delta g(v, S_i) = |N(v) \cap S_i| - \delta c(|S_i|)$, where $\delta c(x) = c(x+1) - c(x)$, for $x \in \mathbb{R}_+$, and $N(v)$ denotes the set of neighbors of vertex v . The two summation elements in the greedy index $\delta g(v, S_i)$ account for the two underlying objectives of minimizing the number of cut edges and balancing of the partition sizes. Notice that the component $|N(v) \cap S_i|$ corresponds to the number of neighbours of vertex v that are assigned to partition S_i . In other words, this corresponds to the degree of vertex v in the subgraph induced by S_i . On the other hand, the component $\delta c(|S_i|)$ can be interpreted as the marginal cost of increasing the partition i by one additional vertex.

For our special family of cost functions $c(x) = \alpha x^\gamma$, we have $\delta c(x) = \alpha \gamma x^{\gamma-1}$. For $\gamma = 1$, the greedy index rule corresponds to assigning a new vertex v to partition i with the largest number of neighbours in S_i , i.e. $|N(v) \cap S_i|$. This is one of the greedy rules considered by Stanton and Kliot [373], and is a greedy rule that may result in highly imbalanced partition sizes.

On the other hand, in case of quadratic cost $c(x) = \frac{1}{2}x^2$, the greedy index is $|N(v) \cap S_i| - |S_i|$, and the greedy assignment corresponds to assigning a new vertex v to partition i that minimizes the number of *non-neighbors* of v inside S_i , i.e. $|S_i \setminus N(v)|$. Hence, this yields the following heuristic: *place a vertex to the partition with the least number of*

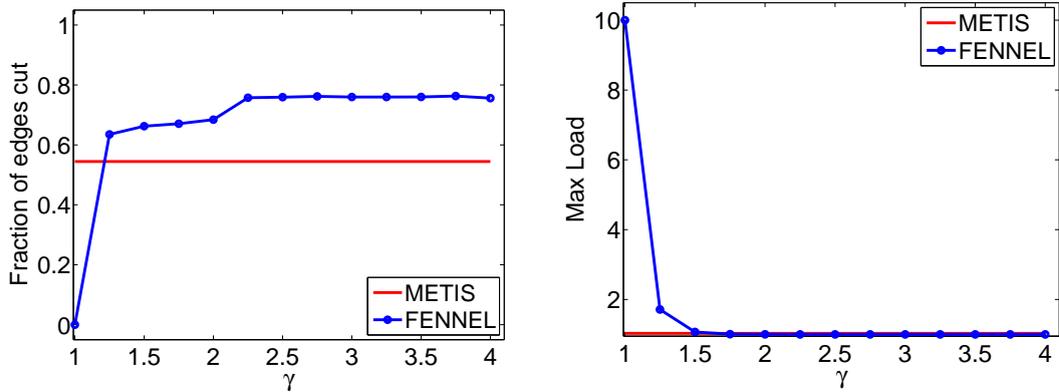


FIGURE 8.1: Fraction of edges cut λ and maximum load normalized ρ as a function of γ , ranging from 1 to 4 with a step of 0.25, over five randomly generated power law graphs with slope 2.5. The straight lines show the performance of METIS.

non-neighbors [344]. This assignment accounts for both the cost of cut edges and the balance of partition sizes.

Finally, we outline that in many applications there exist very strict constraints on the load balance. Despite the fact that we investigate the effect of the parameter γ on the load balance, one may apply the following algorithm, which enforces to consider only machines whose load is at most $\nu \times \frac{n}{k}$. This algorithm for $1 \leq \gamma \leq 2$ amounts to interpolating between the basic heuristics of [373] and [344]. The overall complexity of our algorithm is $O(n + m)$.

Greedy vertex assignment with threshold ν

- Let $I_\nu = \{i : \mu_i \leq \nu \frac{n}{k}\}$. Assign vertex v to partition $i \in I_\nu$ such that $\delta g(v, S_i) \geq \delta g(v, S_j)$, for all $j \in I_\nu$

8.4 Experimental Evaluation

In this section we present results of our experimental evaluations of the quality of graph partitions created by our method and compare with alternative methods. We first describe our experimental setup in Sections 8.4.1, and then present our findings using synthetic and real-world graphs, in Section 8.4.2 and 8.4.3, respectively.

8.4.1 Experimental Setup

The real-world graphs used in our experiments are shown in Table 9.2. Multiple edges, self loops, signs and weights were removed, if any. Furthermore, we considered the largest connected component from each graph in order to ensure that there is a non-zero

	Nodes	Edges	Description
amazon0312	400 727	2 349 869	Co-purchasing
amazon0505	410 236	2 439 437	Co-purchasing
amazon0601	403 364	2 443 311	Co-purchasing
as-735	6 474	12 572	Auton. Sys.
as-Skitter	1 694 616	11 094 209	Auton. Sys.
as-caida	26 475	53 381	Auton. Sys.
ca-AstroPh	17 903	196 972	Collab.
ca-CondMat	21 363	91 286	Collab.
ca-GrQc	4 158	13 422	Collab.
ca-HepPh	11 204	117 619	Collab.
ca-HepTh	8 638	24 806	Collab.
cit-HepPh	34 401	420 784	Citation
cit-HepTh	27 400	352 021	Citation
cit-Patents	3 764 117	16 511 740	Citation
email-Enron	33 696	180 811	Email
email-EuAll	224 832	339 925	Email
epinions	119 070	701 569	Trust
Epinions1	75 877	405 739	Trust
LiveJournal1	4 843 953	42 845 684	Social
p2p-Gnutella04	10 876	39 994	P2P
p2p-Gnutella05	8 842	31 837	P2P
p2p-Gnutella06	8 717	31 525	P2P
p2p-Gnutella08	6 299	20 776	P2P
p2p-Gnutella09	8 104	26 008	P2P
p2p-Gnutella25	22 663	54 693	P2P
p2p-Gnutella31	62 561	147 878	P2P
roadNet-CA	1 957 027	2 760 388	Road
roadNet-PA	1 087 562	1 541 514	Road
roadNet-TX	1 351 137	1 879 201	Road
Slashdot0811	77 360	469 180	Social
Slashdot0902	82 168	504 230	Social
Slashdot081106	77 258	466 661	Social
Slashdot090216	81 776	495 661	Social
Slashdot090221	82 052	498 527	Social
usroads	126 146	161 950	Road
wb-cs-stanford	8 929	2 6320	Web
web-BerkStan	654 782	6 581 871	Web
web-Google	855 802	4 291 352	Web
web-NotreDame	325 729	1 090 108	Web
web-Stanford	255 265	1 941 926	Web
wiki-Talk	2 388 953	4 656 682	Web
Wikipedia-20051105	1 596 970	18 539 720	Web
Wikipedia-20060925	2 935 762	35 046 792	Web
Twitter	41 652 230	1 468 365 182	Social

TABLE 8.2: Datasets used in our experiments.

number of edges cut. All graphs are publicly available on the Web. All algorithms have been implemented in JAVA, and all experiments were performed on a single machine, with Intel Xeon CPU at 3.6GHz, and 16GB of main memory. Wall-clock times include only the algorithm execution time, excluding the required time to load the graph into memory.

Method	BFS		Random	
	λ	ρ	λ	ρ
H	96.9%	1.01	96.9%	1.01
B [373]	97.3%	1.00	96.8%	1.00
DG [373]	0%	32	43%	1.48
LDG [373]	34%	1.01	40%	1.00
EDG [373]	39%	1.04	48%	1.01
T [373]	61%	2.11	78%	1.01
LT [373]	63%	1.23	78%	1.10
ET [373]	64%	1.05	79%	1.01
NN [344]	69%	1.00	55%	1.03
Fennel	14%	1.10	14%	1.02
METIS [245]	8%	1.00	8%	1.02

TABLE 8.3: Performance of various existing methods on amazon0312, k is set to 32.

In our synthetic experiments, we use two random graph models. The first model is the hidden partition model [122]. It is specified by four parameters: the number of vertices n , the number of clusters k , the intercluster and intracluster edge probabilities p and q , respectively. First, each vertex is assigned to one of k clusters uniformly at random. We add an edge between two vertices of the same (different) cluster(s) with probability p (q) independently of the other edges. We denote this model as $\text{HP}(n, k, p, q)$. The second model we use is a standard model for generating random power law graphs. Specifically, we first generate a power-law degree sequence with a given slope δ and use the Chung-Lu random graph model to create an instance of a power law graph [114]. The model $\text{CL}(n, \delta)$ has two parameters: the number of vertices n and the slope δ of the expected power law degree sequence.

We evaluate our algorithms by measuring two quantities from the resulting partitions. In particular, for a fixed partition \mathcal{P} we use the measures of the *fraction of edges cut* λ and the *normalized maximum load* ρ , defined as

$$\lambda = \frac{\# \text{ edges cut by } \mathcal{P}}{\# \text{ total edges}} = \frac{|\partial e(\mathcal{P})|}{m}, \text{ and}$$

$$\rho = \frac{\text{maximum load}}{\frac{n}{k}}.$$

Throughout this section, we also use the notation $\lambda_{\mathcal{M}}$ and $\mu_{\mathcal{M}}$ to indicate the partitioning method \mathcal{M} used in a particular context. In general, we omit indices whenever it is clear to which partition method we refer to. Notice that $k \geq \rho \geq 1$ since the maximum load of a cluster is at most n and there always exists at least one cluster with at least $\frac{n}{k}$ vertices.

In Section 8.4.2, we use the greedy vertex assignment without any threshold. Given that we are able to control ground truth, we are mainly interested in understanding the effect of the parameter γ on the tradeoff between the fraction of edges cut and the normalized maximum load. In Section 8.4.3, the setting of the parameters we use throughout our experiments is $\gamma = \frac{3}{2}$, $\alpha = \sqrt{k} \frac{m}{n^{3/2}}$, and $\nu = 1.1$. The choice of γ is based on our findings from Section 8.4.2 and of α based on Section 8.2. Finally, $\nu = 1.1$ is a reasonable load balancing factor for real-world settings.

As our competitors we use state-of-the-art heuristics. Specifically, in our evaluation we consider the following heuristics from [373], which we briefly describe here for completeness. Let v be the newly arrived vertex.

- **Balanced (B)**: place v to the cluster S_i with minimal size.
- **Hash partitioning (H)**: place v to a cluster chosen uniformly at random.
- **Deterministic Greedy (DG)**: place v to S_i that maximizes $|N(v) \cap S_i|$.
- **Linear Weighted Deterministic Greedy (LDG)**: place v to S_i that maximizes $|N(v) \cap S_i| \times \left(1 - \frac{|S_i|}{k}\right)$.
- **Exponentially Weighted Deterministic Greedy (EDG)**: place v to S_i that maximizes $|N(v) \cap S_i| \times \left(1 - \exp\left(|S_i| - \frac{n}{k}\right)\right)$.
- **Triangles (T)**: place v to S_i that maximizes $t_{S_i}(v)$.
- **Linear Weighted Triangles (LT)**: place v to S_i that maximizes $t_{S_i}(v) \times \left(1 - \frac{|S_i|}{k}\right)$.
- **Exponentially Weighted Triangles (ET)**: place v to S_i that maximizes $t_{S_i}(v) \times \left(1 - \exp\left(|S_i| - \frac{n}{k}\right)\right)$.
- **Non-Neighbors (NN)**: place v to S_i that minimizes $|S_i \setminus N(v)|$.

In accordance with [373], we observed that LDG is the best performing heuristic. Even if Stanton and Kliot do not compare with NN, LDG outperforms it also. Non-neighbors typically have very good load balancing properties, as LDG as well, but cut significantly more edges. Table 8.3 shows the typical performance we observe across all datasets. Specifically, it shows λ and ρ for both BFS and random order for `amazon0312`. DFS order is omitted since qualitatively it does not differ from BFS. We observe that LDG is the best competitor, Fennel outperforms all existing competitors and is inferior to METIS, but of comparable performance. In whatever follows, whenever we refer to the best competitor, unless otherwise mentioned we refer to LDG.

m	k	Fennel		METIS	
		λ	ρ	λ	ρ
7 185 314	4	62.5 %	1.04	65.2%	1.02
6 714 510	8	82.2 %	1.04	81.5%	1.02
6 483 201	16	92.9 %	1.01	92.2%	1.02
6 364 819	32	96.3%	1.00	96.2%	1.02
6 308 013	64	98.2%	1.01	97.9%	1.02
6 279 566	128	98.4 %	1.02	98.8%	1.02

TABLE 8.4: Fraction of edges cut λ and normalized maximum load ρ for Fennel and METIS [245] averaged over 5 random graphs generated according to the HP(5000,0.8,0.5) model.

8.4.2 Synthetic Datasets

Before we delve into our findings, it is worth summarizing the main findings of this section. (a) For all synthetic graphs we generated, the value $\gamma = \frac{3}{2}$ achieves the best performance *pointwise*, not in average. (b) The effect of the stream order is minimal on the results. Specifically, when $\gamma \geq \frac{3}{2}$ all orders result in the same qualitative results. When $\gamma < \frac{3}{2}$ BFS and DFS orders result in the same results which are worse with respect to load balancing –and hence better for the edge cuts– compared to the random order. (c) Fennel’s performance is comparable to METIS.

Hidden Partition: We report averages over five randomly generated graphs according to the model HP(5000, k , 0.8, 0.5) for each value of k we use. We study (a) the effect of the parameter γ , which parameterizes the function $c(x) = \alpha x^\gamma$, and (b) the effect of the number of clusters k .

We range γ from 1 to 4 with a step of 1/4, for six different values of k shown in the second column of Table 8.4. For all k , we observe, consistently, the following behavior: for $\gamma = 1$ we observe that $\lambda = 0$ and $\rho = k$. This means that one cluster receives all vertices. For any γ greater than 1, we obtain excellent load balancing with ρ ranging from 1 to 1.05, and the same fraction of edges cut with METIS up to the first decimal digit. This behavior was not expected a priori, since in general we expect λ shifting from small to large values and see ρ shifting from large to small values as γ grows. Given the insensitivity of Fennel to γ in this setting, we fix $\gamma = \frac{3}{2}$ and present in Table 8.4 our findings. For each k shown in the second column we generate five random graphs. The first column shows the average number of edges. Notice that despite the fact that we have only 5,000 vertices, we obtain graphs with several millions of edges. The four last columns show the performance of Fennel and METIS. As we see, their performance is comparable and in one case ($k=128$) Fennel clearly outperforms METIS.

Power Law: It is well known that power law graphs have no good cuts [201], but they are commonly observed in practice. We examine the effect of parameter γ for k fixed to 10.

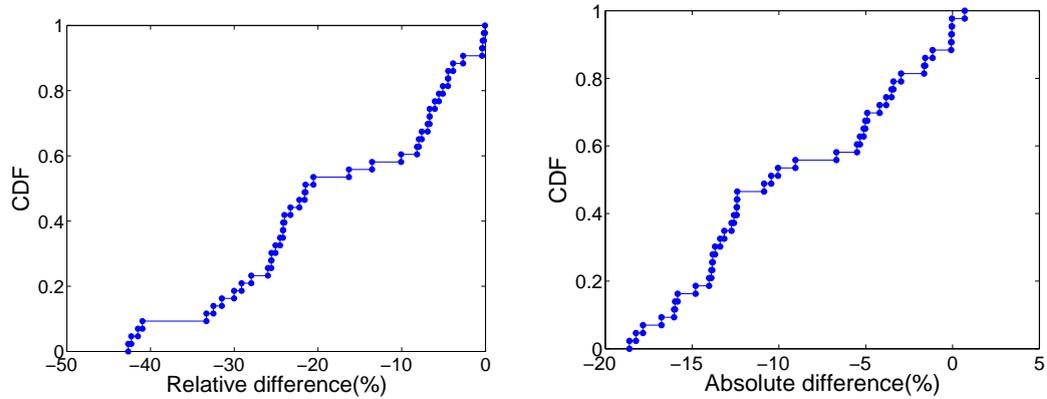


FIGURE 8.2: (Left) CDF of the relative difference $(\frac{\lambda_{\text{Fennel}} - \lambda_c}{\lambda_c}) \times 100\%$ of percentages of edges cut of our method and the best competitor. (Right) Same but for the absolute difference $(\lambda_{\text{Fennel}} - \lambda_c) \times 100\%$.

In contrast to the hidden partition experiment, we observe the expected tradeoff between λ and ρ as γ changes. We generate five random power law graphs CL(20 000,2.5), since this value matches the slope of numerous real-world networks [315]. Figure 8.1 shows the tradeoff when γ ranges from 1 to 4 with a step of 0.25 for the random stream order. The straight line shows the performance of METIS. As we see, when $\gamma < 1.5$, ρ is unacceptably large for demanding real-world applications. When $\gamma = 1.5$ we obtain essentially the same load balancing performance with METIS. Specifically, $\rho_{\text{Fennel}} = 1.02$, $\rho_{\text{METIS}} = 1.03$. The corresponding cut behavior for $\gamma = 1.5$ is $\lambda_{\text{Fennel}} = 62.58\%$, $\lambda_{\text{METIS}} = 54.46\%$. Furthermore, we experimented with the random, BFS and DFS stream orders. We observe that the only major difference between the stream orders is obtained for $\gamma = 1.25$. For all other γ values the behavior is identical. For $\gamma = 1.25$ we observe that BFS and DFS stream orders result in significantly worse load balancing properties. Specifically, $\rho_{\text{BFS}} = 3.81$, $\rho_{\text{DFS}} = 3.73$, $\rho_{\text{Random}} = 1.7130$. The corresponding fractions of edges cut are $\lambda_{\text{BFS}} = 37.83\%$, $\lambda_{\text{DFS}} = 38.85\%$, and $\lambda_{\text{Random}} = 63.51\%$.

8.4.3 Real-World Datasets

Again, before we delve into the details of the experimental results, we summarize the main points of this Section: (1) *Fennel is superior to existing streaming partitioning algorithms*. Specifically, it consistently, over a wide range of k values and over all datasets, performs better than the current state-of-the-art. *Fennel* achieves excellent load balancing with significantly smaller edge cuts. (2) For smaller values of k (less or equal than 64) the observed gain is more pronounced. (c) *Fennel is fast*. Our implementation scales well with the size of the graph. It takes about 40 minutes to partition the Twitter graph which has more than 1 billion of edges.

k	Relative Gain	$\rho_{\text{Fennel}} - \rho_c$
2	25.37%	0.47%
4	25.07%	0.36%
8	26.21%	0.18%
16	22.07%	-0.43%
32	16.59%	-0.34%
64	14.33%	-0.67%
128	13.18%	-0.17%
256	13.76%	-0.20%
512	12.88%	-0.17%
1024	11.24%	-0.44%

TABLE 8.5: The relative gain $(1 - \frac{\lambda_{\text{Fennel}}}{\lambda_c}) \times 100\%$ and load imbalance, where subindex c stands for the best competitor, averaged over all datasets in Table 9.2 as a function of k .

Twitter Graph. Twitter graph is the largest graph in our collection of graphs, with more than 1.4 billion edges. This feature makes it the most interesting graph from our collection, even if, admittedly, is a graph that can be loaded into the main memory. The results of Fennel on this graph are excellent. Specifically, Table 8.1 shows the performance of Fennel, the best competitor LDG, the baseline Hash Partition and METIS for $k = 2, 4$ and 8. All methods achieve balanced partitions, with $\rho \leq 1.1$. Fennel, is the only method that always attains this upper bound. However, this reasonable performance comes with a high gain for λ . Specifically, we see that Fennel achieves better performance of $k = 2$ than METIS. Furthermore, Fennel requires 42 minutes whereas METIS requires $8\frac{1}{2}$ hours. Most importantly, Fennel outperforms LDG consistently. Specifically, for $k = 16, 32$ and 64, Fennel achieves the following results $(\lambda, \rho) = (59\%, 1.1)$, $(67\%, 1.1)$, and $(73\%, 1.1)$, respectively. Linear weighted degrees (LDG) achieves $(76\%, 1.13)$, $(80\%, 1.15)$, and $(84\%, 1.14)$, respectively. Now we turn our attention to smaller but reasonably-sized datasets.

In Figure 8.2, we show the distribution of the difference of the fraction of edges cut of our method and that of the best competitor, conditional on that the maximum observed load is at most 1.1. This distribution is derived from the values observed by partitioning each input graph from our set averaged over a range of values of parameter k that consists of values $2, 4, \dots, 1024$. These results demonstrate that the fraction of edges cut by our method is smaller than that of the best competitor in all cases. Moreover, we observe that the median difference (relative difference) is in the excess of 20% (15%), thus providing appreciable performance gains.

Furthermore, in Table 8.5, we present the average performance gains conditional on the number of partitions k . These numerical results amount to an average relative reduction of the fraction of edges cut in the excess of 18%. Moreover, the performance gains

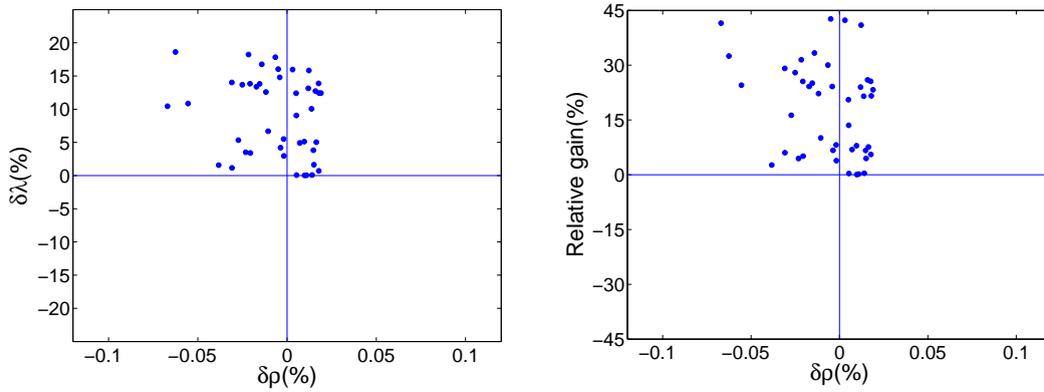


FIGURE 8.3: Absolute difference $\delta\lambda$ and Relative gain versus the maximum load imbalance $\delta\rho$.

observed are consistent across different values of parameter k , and are more pronounced for smaller values of k .

Bicriteria. In our presentation of experimental results so far, we focused on the fraction of edges cut by conditioning on the cases where the normalized maximum load was smaller than a fixed threshold. We now provide a closer look at both criteria and their relation. In Figure 8.3, we consider the difference of the fraction of edges cut vs. the difference of normalized maximum loads of the best competitor and our method. We observe that in all the cases, the differences of normalized maximum loads are well within 10% while the fraction of edges cut by our method is significantly smaller. These results confirm that the observed reduction of the fraction of edges cut by our method is not at the expense of an increased maximum load.

Speed of partitioning. We now turn our attention to the efficiency of our method with respect to the running time to partition a graph. Our graph partitioning algorithm is a one-pass streaming algorithm, which allows for fast graph partitioning. In order to support this claim, in Figure 8.4, we show the run time it took to partition each graph from our dataset vs. the graph size in terms of the number of edges. We observe that it takes in the order of minutes to partition large graphs of tens of millions of edges. As we also mentioned before, partitioning the largest graph from our dataset collection took about 40 minutes.

8.5 System Evaluation

Evaluating a partitioning algorithm is not an easy task from a systems perspective, since it depends on system characteristics. For instance, in a large-scale production data-center it is more important to balance the traffic across clusters than the traffic or the amount of computation executed within a cluster. However, in a small-scale cluster

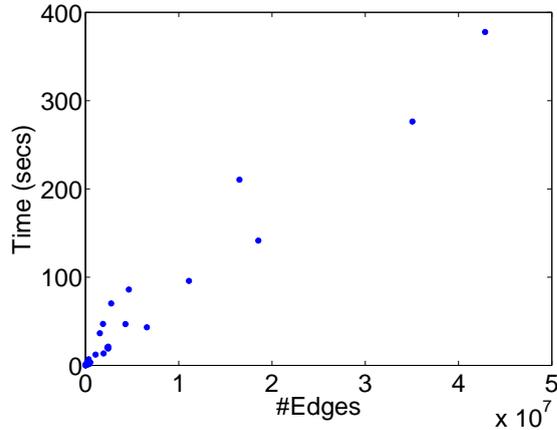


FIGURE 8.4: Fennel: time vs. number of edges.

# Clusters (k)	Run time [s]		Communication [MB]	
	Hash	Fennel	Hash	Fennel
4	32.27	25.49	321.41	196.9
8	17.26	15.14	285.35	180.02
16	10.64	9.05	222.28	148.67

TABLE 8.6: The average duration of a step and the average amount of MB exchanged per node and per step during the execution of PageRank on LiveJournal data set.

consisting of few tens of nodes, rented by a customer from a large cloud provider such as Amazon’s Elastic Map-Reduce, it is important to minimize the network traffic across the nodes but also to balance well the computational load on each node. Given this diversity of scenarios, a detailed evaluation is out of the scope of this work. Here, we perform a basic experiment to verify the superiority of our proposed method versus the de facto standard of hash partitioning with respect to speeding up a large-scale computation.

We select PageRank as the computation of interest. Notice, that an advantage of Fennel is that it gives a flexibility in choosing a suitable objective that accommodates the needs of the specific application. We demonstrate the efficiency and flexibility of Fennel with the typical Elastic Map-Reduce scenario in mind. We set up a cluster and we vary the number of nodes to 4, 8 and 16 nodes. Each node is equipped with Intel Xeon CPU at 2.27 GHz and 12 GB of main memory. On the cluster we run Giraph, a graph processing platform running on the top of Hadoop. We implemented a PageRank algorithm on Giraph and we run it on Live Journal data set¹.

Since the complexity of PageRank depends on the number of edges and not vertices, we use a version of the Fennel objective (eq. 8.1) that balances the number of edges per cluster. In particular, we choose the $c_{IN}(\mathcal{P}) = \sum_{i=1}^k e(S_i, S_i)^\gamma$ with $\gamma = 1.5$.

¹Twitter data set was too large to fit on a 16-nodes Giraph cluster.

We compare with hash partitioning, the default partitioning scheme used by Giraph. We look at two metrics. The first is the average duration of an iteration of the PageRank algorithm. This metric is directly proportional to the actual running time and incorporates both the processing and the communication time. The second metric is the average number of Megabytes transmitted by a cluster node in each iteration. This metric directly reflects the quality of the cut and is proportional to the incurred network load.

The results are shown in Table 8.6. We see that Fennel has the best run time in all cases. This is because it achieves the best balance between the computation and communication load. Hash partitioning takes 25% more time than Fennel and it also has a much higher traffic load.

8.6 Discussion

In this section we discuss some of the extensions that can be accommodated by our framework and discuss some details about distributed implementation.

Asymmetric edge costs. As discussed in Section 8.5, in some application scenarios some edges that cross partition boundaries may be more costly than other. For example, this is the case if individual graph partitions are assigned to machines in a data center and these machines are connected with an asymmetric network topology, so that the available network bandwidth varies across different pairs of machines, e.g. intra-rack vs. inter-rack machines in standard data center architectures [203]. Another example are data center network topologies where the number of hops between different pairs of machines vary substantially, e.g. torus topologies [127]. In such scenarios, it may be beneficial to partition a graph by accounting for the aforementioned asymmetries of edge-cut costs. This can be accommodated by appropriately defining the inter-partition cost function in our framework.

Distributed implementation. Our streaming algorithm requires computing marginal value indices that can be computed in a distributed fashion by maintaining local views on a global state. For concreteness, let us consider the traditional objective where the inter-partition cost is a linear function of the total number of cut edges and the intra-partition cost is a sum of convex functions of vertex cardinalities of individual partitions. In this case, computing the marginal value indices requires to compute per each vertex arrival: (1) the number of neighbors of given vertex that were already assigned to given cluster of vertices, and (2) the number of vertices that were already assigned per cluster. The former corresponds to a set-intersection query and can be efficiently implemented

by standard methods, e.g. using data structures such as minwise hashing [349]. The latter is a simple count tracking problem. Further optimizations could be made by trading accuracy for reduction of communication overhead by updating of the local views at a smaller rate than the rate of vertex arrival.

Chapter 9

PEGASUS: A System for Large-Scale Graph Processing

9.1 Introduction

In this Chapter we describe PEGASUS, an open source Peta Graph Mining library which performs typical graph mining tasks such as computing the diameter of a graph, computing the radius of each node, finding the connected components, (see also Chapter 7), and computing the importance score of nodes. The main idea behind PEGASUS is to capitalize on matrix-vector multiplication as a main primitive for the software engineer. Inspired by the work of [387] which showed that triangles can be estimated by few matrix-vector multiplications, PEGASUS introduces a set of different operators which solve a variety of graph mining tasks together with an optimized implementation of matrix-vector multiplications in MAPREDUCE. PEGASUS is a solid engineering effort which allows us to manipulate large-scale graphs. Since the introduction of PEGASUS, other large-scale graph processing systems have been introduced, among them Google's Pregel [295], LinkedIn's Giraph [3] and GraphLab [291]. It is worth mentioning that Giraph uses several algorithms and ideas from PEGASUS, including the connected components algorithm. Also, PEGASUS has been included in HADOOP for Windows Azure [1].

Outline: This Chapter is organized as follows: Section 9.2 presents the proposed method. Section 9.3 presents HADOOP implementations and Section 9.4 timings. Section 9.5 shows findings of PEGASUS in several real-world networks.

9.2 Proposed Method

Consider the following assignment $v' \leftarrow M \times v$ where $M \in \mathbb{R}^{m \times n}$, $v \in \mathbb{R}^n$. The i -th coordinate of v' is $v'_i = \sum_{j=1}^n m_{i,j} v_j$, $i = 1, \dots, m$. Typically in our applications, M is the adjacency matrix representation of a graph and therefore we are going to assume in the following that $m = n$, unless otherwise noticed.

There are three types of operations in the previous formula:

1. `combine2`: multiply $m_{i,j}$ and v_j .
2. `combineAll`: sum n multiplication results for node i .
3. `assign`: overwrite the previous value of v_i with the new result to make v'_i .

We introduce an abstraction of the basic matrix-vector multiplication, called Generalized Iterative Matrix-Vector multiplication. The corresponding programming primitive is the `GIM-V` primitive on which PEGASUS is based. The ‘Iterative’ in `GIM-V` denotes that we apply the \times_G operation until a convergence criterion is met. Specifically, let us define the operator \times_G as follows:

$$v' = M \times_G v$$

where $v'_i = \text{assign}(v_i, \text{combineAll}_i(\{x_j \mid j = 1..n, \text{ and } x_j = \text{combine2}(m_{i,j}, v_j)\}))$.

The functions `combine2()`, `combineAll()`, and `assign()` have the following interpretation, generalizing the product, sum and assignment of the traditional matrix-vector multiplication:

1. `combine2($m_{i,j}, v_j$)` : combine $m_{i,j}$ and v_j .
2. `combineAll $_i(x_1, \dots, x_n)$` : combine all the results from `combine2()` for node i .
3. `assign(v_i, v_{new})` : decide how to update v_i with v_{new} .

In the following sections we show how different choices of `combine2()`, `combineAll $_i()$` and `assign()` allow us to solve several important graph mining tasks. Before that, we want to highlight the strong connection of `GIM-V` with SQL. When `combineAll $_i()$` and `assign()` can be implemented by user defined functions, the operator \times_G can be expressed concisely in terms of SQL. This viewpoint is important when we implement `GIM-V` in large-scale parallel processing platforms, including HADOOP, if they can be customized to support several SQL primitives including JOIN and GROUP BY. Suppose we

```

SELECT E.sid, combineAllE.sid(combine2(E.val,V.val))
FROM E, V
WHERE E.did=V.id
GROUP BY E.sid

```

TABLE 9.1: GIM-V in terms of SQL.

have an **edge** table $E(\text{sid}, \text{did}, \text{val})$ and a **vector** table $V(\text{id}, \text{val})$, corresponding to a matrix and a vector, respectively. Then, \times_G corresponds to the SQL statement in Table 9.1. We assume that we have (built-in or user-defined) functions, $\text{combineAll}_i()$ and $\text{combine2}()$, and we also assume that the resulting table/vector will be fed into the $\text{assign}()$ function (omitted, for clarity).

In the following sections we show how we can customize **GIM-V**, to handle important graph mining operations including PageRank, Random Walk with Restart, diameter estimation, and connected components.

9.2.1 GIM-V and PageRank

Our first warm-up application of **GIM-V** is PageRank, a famous algorithm that was used by Google to calculate relative importance of web pages [89]. The PageRank vector p of n web pages satisfies the following eigenvector equation:

$$p = (cE^T + (1 - c)U)p$$

where c is a damping factor (usually set to 0.85), E is the row-normalized adjacency matrix (source, destination), and U is a matrix with all elements set to $1/n$.

To calculate the eigenvector p we can use the power method, which multiplies an initial vector with the matrix, several times. We initialize the current PageRank vector p^{cur} and set all its elements to $1/n$. Then the next PageRank p^{next} is calculated by $p^{next} = (cE^T + (1 - c)U)p^{cur}$. We continue to perform the multiplication until p converges.

PageRank is a direct application of **GIM-V**, i.e., $p^{next} = M \times_G p^{cur}$. Matrix M is E^T , i.e., the column-normalized version of the adjacency matrix. The three operations are defined as follows:

1. $\text{combine2}(m_{i,j}, v_j) = c \times m_{i,j} \times v_j$
2. $\text{combineAll}_i(x_1, \dots, x_n) = \frac{(1-c)}{n} + \sum_{j=1}^n x_j$
3. $\text{assign}(v_i, v_{new}) = v_{new}$

9.2.2 GIM-V and Random Walk with Restart

Random Walk with Restart (RWR) is closely related to Personalized pagerank, a popular algorithm to measure the relative proximity of vertices with respect to a given vertex [328]. In RWR, the proximity vector r_k of vertex k satisfies the equation:

$$r_k = cMr_k + (1 - c)e_k$$

where e_k is the k -th unit vector in \mathbb{R}^n , c is a restart probability parameter which is typically set to 0.85 [328] and M is as in Section 9.2.1. In GIM-V, RWR is formulated by $r_k^{next} = M \times_G r_k^{cur}$ where the three operations are defined as follows:

1. `combine2`($m_{i,j}, v_j$) = $c \times m_{i,j} \times v_j$
2. `combineAlli`(x_1, \dots, x_n) = $(1 - c)\delta_{ik} + \sum_{j=1}^n x_j$, where δ_{ik} is the *Kronecker delta*, equal to 1 if $i = k$ and 0 otherwise
3. `assign`(v_i, v_{new}) = v_{new}

9.2.3 GIM-V and Diameter Estimation

In Chapter 7 we discussed HADI, an algorithm that estimates the diameter and radius distribution of a large-scale graph. HADI can be presented within the framework of PEGASUS, since the number of neighbors reachable from vertex i within h hops is encoded in a probabilistic bitstring b_i^h which is updated as follows [168]:

$$b_i^{h+1} = b_i^h \text{ BITWISE-OR } \{b_k^h \mid (i, k) \in E\}$$

In GIM-V, the bitstring update of HADI is represented by

$$b^{h+1} = M \times_G b^h$$

where M is the adjacency matrix, b^{h+1} is a vector of length n which is updated by $b_i^{h+1} = \text{assign}(b_i^h, \text{combineAll}_i(\{x_j \mid j = 1..n, \text{ and } x_j = \text{combine2}(m_{i,j}, b_j^h)\}))$, and the three PEGASUS operations are defined as follows:

1. `combine2`($m_{i,j}, v_j$) = $m_{i,j} \times v_j$.
2. `combineAlli`(x_1, \dots, x_n) = $\text{BITWISE-OR}\{x_j \mid j = 1..n\}$
3. `assign`(v_i, v_{new}) = $\text{BITWISE-OR}(v_i, v_{new})$.

The \times_G operation is run iteratively until the bitstring of each vertex remains the same.

9.2.4 GIM-V and Connected Components

We propose HCC, a new algorithm for finding connected components in large graphs. The main idea is as follows: for each vertex i in the graph, we maintain a component identification number (id) c_i^h which is the minimum vertex id within h hops from i .

Initially, c_i^h of vertex i is set to i , i.e., $c_i^0 = i$. In each iteration, each vertex sends its current c_i^h to its neighbors. Then c_i^{h+1} is set to the minimum value among its current component id and the received component ids from its neighbors. The crucial observation is that this communication between neighbors can be formulated in GIM-V as follows:

$$c^{h+1} = M \times_G c^h$$

where M is the adjacency matrix, c^{h+1} is a vector of length n which is updated by $c_i^{h+1} = \text{assign}(c_i^h, \text{combineAll}_i(\{x_j \mid j = 1..n, \text{ and } x_j = \text{combine2}(m_{i,j}, c_j^h)\}))$, and the three PEGASUS operations are defined as follows:

1. $\text{combine2}(m_{i,j}, v_j) = m_{i,j} \times v_j$.
2. $\text{combineAll}_i(x_1, \dots, x_n) = \min\{x_j \mid j = 1..n\}$.
3. $\text{assign}(v_i, v_{new}) = \min(v_i, v_{new})$.

By repeating this process, component ids of nodes in a component are set to the minimum node id of the component. We iteratively do the multiplication until component ids converge. The upper bound of the number of iterations in HCC is d , where d is the diameter of the graph. We notice that because of the small-world phenomenon, see Section 7.1, the diameter of real graphs is small, and therefore in practice HCC completes after a small number of iterations. For a recent work with better practical performance, see [359].

9.3 Hadoop Implementation

Given the main goal of the PEGASUS project is to provide an efficient system to the user/programmer, we discuss different HADOOP implementation approaches, starting out with a naive implementation and progressing to faster methods for GIM-V. The proposed versions are evaluated in Section 9.4.

Algorithm 10 GIM-V BASE Stage 1.

Input: Matrix $M = \{(id_{src}, (id_{dst}, mval))\}$, Vector $V = \{(id, vval)\}$
Output: Partial vector $V' = \{(id_{src}, \text{combine2}(mval, vval))\}$

```

1: Stage1-Map(Key  $k$ , Value  $v$ ):
2: if  $(k, v)$  is of type V then
3:   Output( $k, v$ ); // ( $k: id, v: vval$ )
4: else if  $(k, v)$  is of type M then
5:    $(id_{dst}, mval) \leftarrow v$ ;
6:   Output( $id_{dst}, (k, mval)$ ); // ( $k: id_{src}$ )
7: end if
8:
9: Stage1-Reduce(Key  $k$ , Value  $v[1..m]$ ):
10:  $saved\_kv \leftarrow []$ ;
11:  $saved\_v \leftarrow []$ ;
12: for  $v \in v[1..m]$  do
13:   if  $(k, v)$  is of type V then
14:      $saved\_v \leftarrow v$ ;
15:     Output( $k, ("self", saved\_v)$ );
16:   else if  $(k, v)$  is of type M then
17:     Add  $v$  to  $saved\_kv$ ; // ( $v: (id_{src}, mval)$ )
18:   end if
19: end for
20: for  $(id'_{src}, mval') \in saved\_kv$  do
21:   Output( $id'_{src}, ("others", \text{combine2}(mval', saved\_v))$ );
22: end for

```

9.3.1 GIM-V BASE: Naive Multiplication

GIM-V BASE is a two-stage algorithm whose pseudo code is in Algorithm 10 and 11. The inputs are an edge file and a vector file. Each line of the edge file has the form $(id_{src}, id_{dst}, mval)$ which corresponds to a non-zero entry in the adjacency matrix. Similarly, each line of the vector file has the form $(id, vval)$ which corresponds to an element in vector v . **Stage1** performs the `combine2` operation by combining columns of matrix (id_{dst} of M) with rows of the vector (id of V). The output of **Stage1** are (key, value) pairs where the key is the source vertex id of the matrix (id_{src} of M) and the value is the partially combined result (`combine2`($mval, vval$)). This output of **Stage1** becomes the input of **Stage2**. **Stage2** combines all partial results from **Stage1** and updates the vector. The `combineAlli`() and `assign`() operations are done in line 15 of **Stage2**, where the “self” and “others” tags in line 15 and line 21 of **Stage1** are needed by **Stage2** to distinguish cases appropriately. We note that in Algorithm 10 and 11, `Output`(k, v) means to output data with the key k and the value v .

Algorithm 11 GIM-V BASE Stage 2.**Input:** Partial vector $V' = \{(id_{src}, vval')\}$ **Output:** Result Vector $V = \{(id_{src}, vval)\}$

```

1: Stage2-Map(Key k, Value v):
2:   Output( $k, v$ );
3:
4: Stage2-Reduce(Key k, Value  $v[1..m]$ ):
5:    $others\_v \leftarrow []$ ;
6:    $self\_v \leftarrow []$ ;
7:   for  $v \in v[1..m]$  do
8:      $(tag, v') \leftarrow v$ ;
9:     if  $tag = \text{"same"}$  then
10:       $self\_v \leftarrow v'$ ;
11:    else if  $tag = \text{"others"}$  then
12:      Add  $v'$  to  $others\_v$ ;
13:    end if
14:  end for
15: Output( $k, assign(self\_v, combineAll_k(others\_v))$ );

```

9.3.2 GIM-V BL: Block Multiplication

GIM-V BL is a fast algorithm for GIM-V which is based on block multiplication. The main idea is to group elements of the input matrix into blocks/submatrices of size b by b . Also, we group elements of input vectors into blocks of length b . In practice, grouping means we place all elements of a group into one line of input file. Each block contains only non-zero elements of the matrix/vector. The format of a matrix block with k nonzero elements is $(row_{block}, col_{block}, row_{elem_1}, col_{elem_1}, mval_{elem_1}, \dots, row_{elem_k}, col_{elem_k}, mval_{elem_k})$. Similarly, the format of a vector block with k nonzero elements is $(id_{block}, id_{elem_1}, vval_{elem_1}, \dots, id_{elem_k}, vval_{elem_k})$. Only blocks with at least one nonzero elements are saved to disk. This block encoding forces nearby edges in the adjacency matrix to be closely located; it is different from HADOOP's default behavior which does not guarantee co-locating them. After grouping, GIM-V is performed on blocks, not on individual elements. GIM-V BL is illustrated in Figure 9.1.

In Section 9.4, we observe that GIM-V BL is at least 5 times faster than GIM-V BASE. There are two main reasons for this speed-up.

- **Sorting Time** Block encoding decreases the number of items to be sorted in the shuffling stage of HADOOP. We observe that one of the main efficiency bottlenecks in HADOOP is its shuffling stage where network transfer, sorting, and disk I/O take place.
- **Compression** The size of the data decreases significantly by converting edges and vectors to block format. The reason is that in GIM-V BASE we need $2 \times 4 = 8$

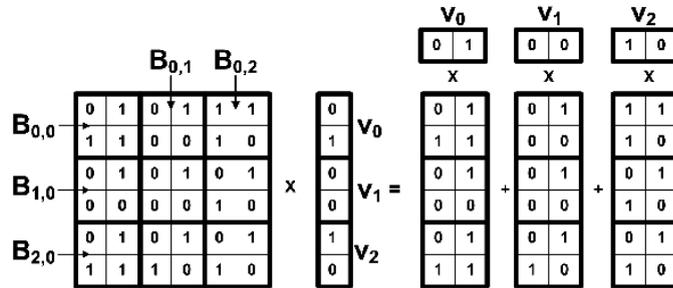


FIGURE 9.1: GIM-V BL using 2 x 2 blocks. $B_{i,j}$ represents a matrix block, and v_i represents a vector block. The matrix and vector are joined block-wise, not element-wise.

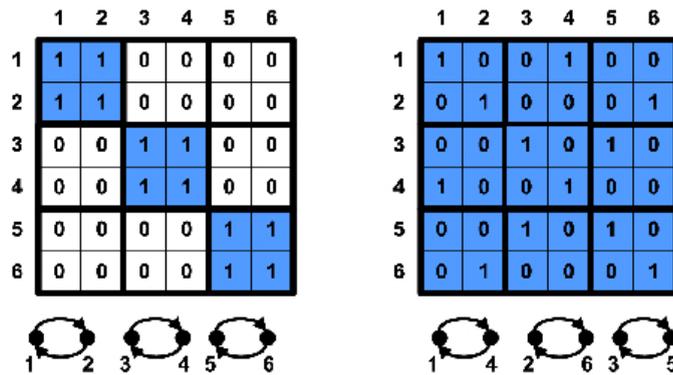


FIGURE 9.2: Clustered vs. non-clustered adjacency matrices for two isomorphic graphs. The edges are grouped into 2 by 2 blocks. The left graph uses only 3 blocks while the right graph uses 9 blocks.

bytes to save each (srcid, dstid) pair. However in GIM-V BL we can specify each *block* using a block row id and a block column id with two 4-byte Integers, and refer to elements inside the block using $2 \times \log b$ bits. This is possible because we can use $\log b$ bits to refer to a row or column inside a block. By this block method we decrease the edge file size. For instance, using block encoding we are able to decrease the size of the YahooWeb graph more than 50%.

9.3.3 GIM-V CL: Clustered Edges

We use co-clustering heuristics, see [332] as a preprocessing step to obtain a better clustering of the edge set. Figure 9.2 illustrates the concept. The preprocessing step needs to be performed only once. If the number of iterations required for the execution of an algorithm is large, then it is beneficial to perform this preprocessing step. Notice that we have two variants of GIM-V: GIM-V CL and GIM-V BL-CL, which are GIM-V BASE and GIM-V BL with clustered edges respectively.

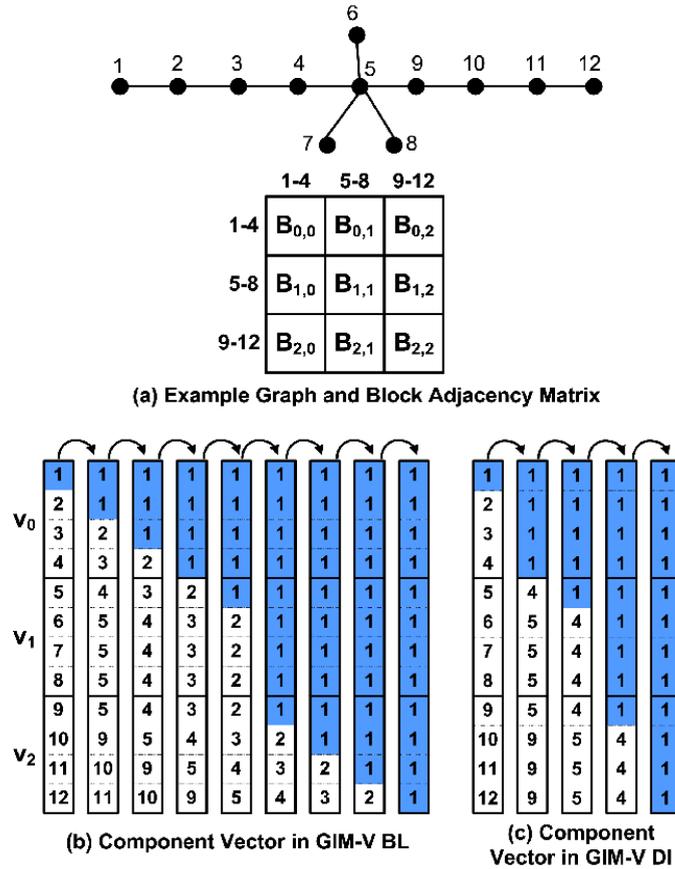


FIGURE 9.3: Propagation of component id(=1) when block width is 4. Each element in the adjacency matrix of (a) represents a 4 by 4 block; each column in (b) and (c) represents the vector after each iteration. GIM-V DL finishes in 4 iterations while GIM-V BL requires 8 iterations.

9.3.4 GIM-V DI: Diagonal Block Iteration

Reducing the number of iterations required for executing an algorithm in MAPREDUCE mitigates the computational cost a lot, since the main bottleneck of GIM-V is its shuffling and disk I/O steps. In HCC, it is possible to decrease the number of iterations when the graph has long chains. The main idea is to multiply diagonal matrix blocks and corresponding vector blocks as much as possible in one iteration. This is illustrated in Figure 9.3.

9.3.5 GIM-V NR: Node Renumbering

In HCC, the minimum vertex id is propagated to the other parts of the graph within at most d steps, where d is the diameter of the graph. If the vertex with the minimum id (which we call ‘minimum node’) is located at the center of the graph, then the number of iterations is small, close to $d/2$. However, if it is located at the boundary of the network,

Algorithm 12 Renumbering the minimum node

Input: Edge $E = \{(id_{src}, id_{dst})\}$,
current minimum vertex id $minid_{cur}$,
new minimum vertex id $minid_{new}$

Output: Renumbered Edge $V = \{(id'_{src}, id'_{dst})\}$

- 1: Renumber-Map(key k , value v):
- 2: $src \leftarrow k$;
- 3: $dst \leftarrow v$;
- 4: **if** $src = minid_{cur}$ **then**
- 5: $src \leftarrow minid_{new}$;
- 6: **else if** $src = minid_{new}$ **then**
- 7: $src \leftarrow minid_{cur}$;
- 8: **end if**
- 9: **if** $dst = minid_{cur}$ **then**
- 10: $dst \leftarrow minid_{new}$;
- 11: **else if** $dst = minid_{new}$ **then**
- 12: $dst \leftarrow minid_{cur}$;
- 13: **end if**
- 14: Output(src, dst);

then the number of iteration can be close to d . Therefore, if we preprocess the edges so that the minimum vertex id is swapped to the center vertex id, the number of iterations and the total running time of HCC would decrease.

Finding the center vertex with the minimum radius could be done with the HADI algorithm. However, the algorithm is expensive for the pre-processing step of HCC. Therefore, we instead propose the following heuristic for finding the center node: we choose the center vertex by sampling from the high-degree vertices. This heuristic is based on the fact that vertices with large degree have small radii [238].

After finding a center node, we need to renumber the edge file to swap the current minimum vertex id with the center vertex id. The MAPREDUCE algorithm for this renumbering is shown in Algorithm 12. Since the renumbering requires only filtering, it can be done with a Map-only job.

9.3.6 Analysis

Finally, we analyze the time and space complexity of GIM-V. It is not hard to observe that one iteration of GIM-V takes $O(\frac{n+m}{M} \log \frac{n+m}{M})$ time, where M stands for the number of machines. Assuming uniformity, mappers and reducers of Stage1 and Stage2 receive $O(\frac{n+m}{M})$ records per machine. The running time is dominated by the sorting time for $\frac{n+m}{M}$ records. GIM-V requires $O(V + E)$ space.

Name	Vertices	Edges	Description
YahooWeb	1,413 M	6,636 M	WWW pages in 2002
LinkedIn	7.5 M	58 M	person-person in 2006
	4.4 M	27 M	person-person in 2005
	1.6 M	6.8 M	person-person in 2004
	85 K	230 K	person-person in 2003
Wikipedia	3.5 M	42 M	doc-doc in 2007/02
	3 M	35 M	doc-doc in 2006/09
	1.6 M	18.5 M	doc-doc in 2005/11
Kronecker	177 K	1,977 M	synthetic
	120 K	1,145 M	synthetic
	59 K	282 M	synthetic
	19 K	40 M	synthetic
WWW-Barabasi	325 K	1,497 K	WWW pages in nd.edu
DBLP	471 K	112 K	document-document
flickr	404 K	2.1 M	person-person
Epinions	75 K	508 K	who trusts whom

TABLE 9.2: Order and size of networks.

9.4 Scalability

We perform experiments to answer the following questions:

- How does GIM-V scale up?
- Which of the proposed optimizations (block multiplication, clustered edges, and diagonal block iteration, vertex renumbering) gives the highest performance gains?

The graphs we use in our experiments are shown in Table 9.2. We run PEGASUS in M45 HADOOP cluster by Yahoo! and our own cluster composed of 9 machines. M45 is one of the top 50 supercomputers in the world with the total 1.5 Pb storage and 3.5 Tb memory. For the performance and scalability experiments, we used synthetic Kronecker graphs [279] since we can generate them with any size, and they are one of the most realistic graphs among synthetic graphs.

9.4.1 Results

We first show how the performance of our method changes as we add more machines. Figure 9.4 shows the running time and performance of GIM-V for PageRank with Kronecker graph of 282 million edges, and size 32 blocks if necessary.

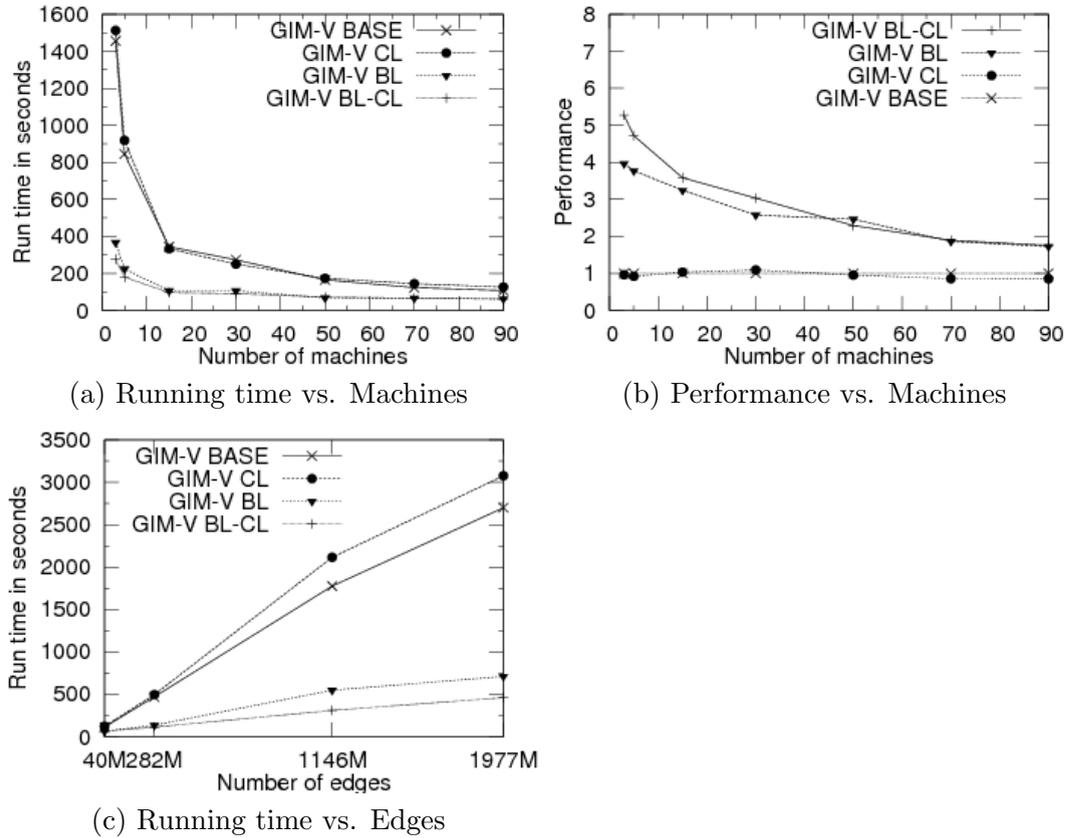


FIGURE 9.4: Scalability and Performance of GIM-V. (a) Running time decreases quickly as more machines are added. (b) The performance(=1/running time) of 'BL-CL' wins more than 5x (for n=3 machines) over the 'BASE'. (c) Every version of GIM-V shows linear scalability.

In Figure 9.4 (a), for all of the methods the running time decreases as we add more machines. Note that clustered edges(GIM-V CL) didn't help performance unless it is combined with block encoding. When it is combined, however, it showed the best performance (GIM-V BL-CL).

In Figure 9.4 (b), we see that the relative performance of each method compared to GIM-V BASE method decreases as number of machines increases. With 3 machines (minimum number of machines which HADOOP 'distributed mode' supports), the fastest method(GIM-V BL-CL) ran 5.27 times faster than GIM-V BASE. With 90 machines, GIM-V BL-CL ran 2.93 times faster than GIM-V BASE. This is expected since there are fixed component(JVM load time, disk I/O, network communication) which can not be optimized even if we add more machines.

Next we show how the performance of our methods changes as the input size grows. Figure 9.4 (c) shows the running time of GIM-V with different number of edges under 10 machines. As we can see, all of the methods scales linearly with the number of edges.

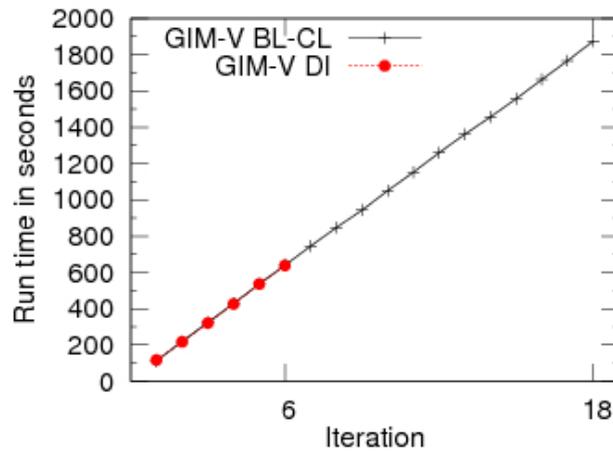


FIGURE 9.5: Comparison of GIM-V DI and GIM-V BL-CL for HCC. GIM-V DI finishes in 6 iterations while GIM-V BL-CL finishes in 18 iterations due to long chains.

Next, we compare the performance of GIM-V DI and GIM-V BL-CL for HCC in graphs with long chains. For this experiment we made a new graph whose diameter is 17, by adding a length 15 chain to the 282 million Kronecker graph which has diameter 2. As we see in Figure 9.5, GIM-V DI finished in 6 iteration while GIM-V BL-CL finished in 18 iteration. The running time of both methods for the first 6 iterations are nearly same. Therefore, the diagonal block iteration method decreases the number of iterations while not affecting the running time of each iteration much.

Finally, we compare the number of iterations with/without renumbering. Figure 9.6 shows the degree distribution of LinkedIn. Without renumbering, the minimum vertex has degree 1, which is not surprising since about 46 % of the vertices have degree 1 due to the power-law behavior of the degree distribution. We show the number of iterations after changing the minimum vertex to each of the top 5 highest-degree vertices in Figure 9.7. We see that the renumbering decreased the number of iterations to 81 % of the original. Similar results are observed for the Wikipedia graph in Figure 9.8 and 9.9. The original minimum vertex has degree 1, and the number of iterations decreased to 83 % of the original after renumbering.

9.5 Pegasus at Work

In this section we evaluate PEGASUS on real-world networks.

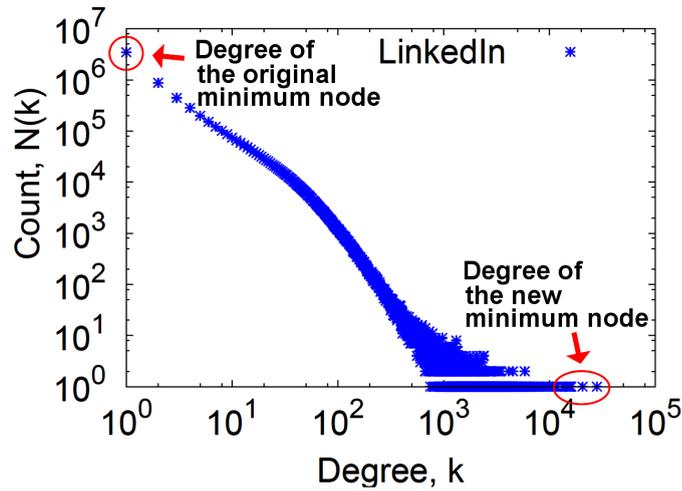


FIGURE 9.6: Degree distribution of LinkedIn. Notice that the original minimum vertex has degree 1, which is highly probable given the power-law behavior of the degree distribution. After the renumbering, the minimum vertex is replaced with a highest-degree node.

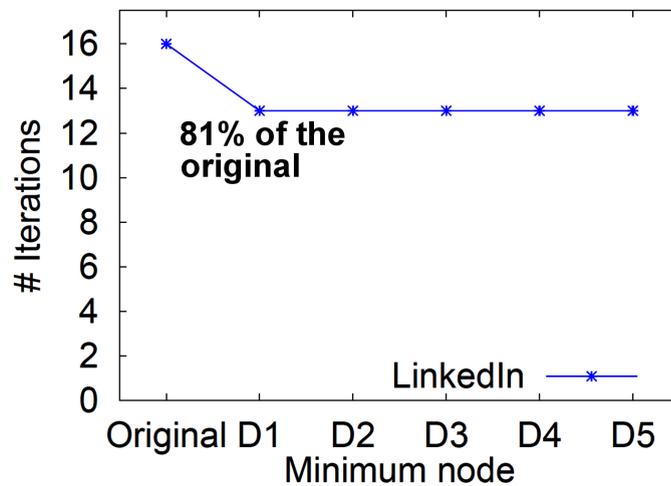


FIGURE 9.7: Number of iterations vs. the minimum vertex of LinkedIn, for connected components. D_i represents the vertex with i -th largest degree. Notice that the number of iterations decreased by 19 % after renumbering.

9.5.1 Connected Components of Real Networks

Figure 9.10 shows the evolution of connected components of LinkedIn and Wikipedia graphs. Figure 9.11 shows the distribution of connected components in the YahooWeb graph. We make the following set of observations.

Power Laws in Connected Components Distributions We observe a power law relation between the count and size of small connected components in Figure 9.10(a),(b) and Figure 9.11. This reflects that the connected components in real networks are formed by preferential attachment processes.

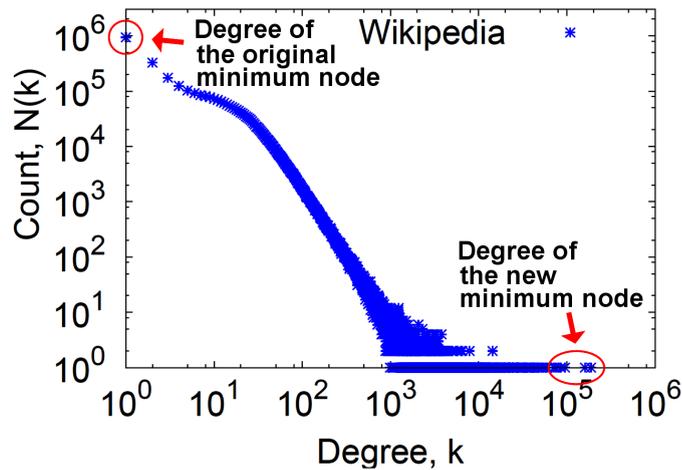


FIGURE 9.8: Degree distribution of Wikipedia. Notice that the original minimum vertex has degree 1, as in LinkedIn. After the renumbering, the minimum vertex is replaced with a highest-degree node.

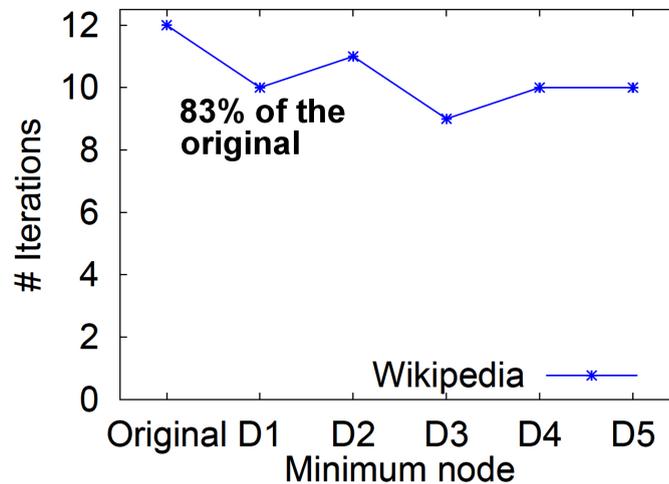
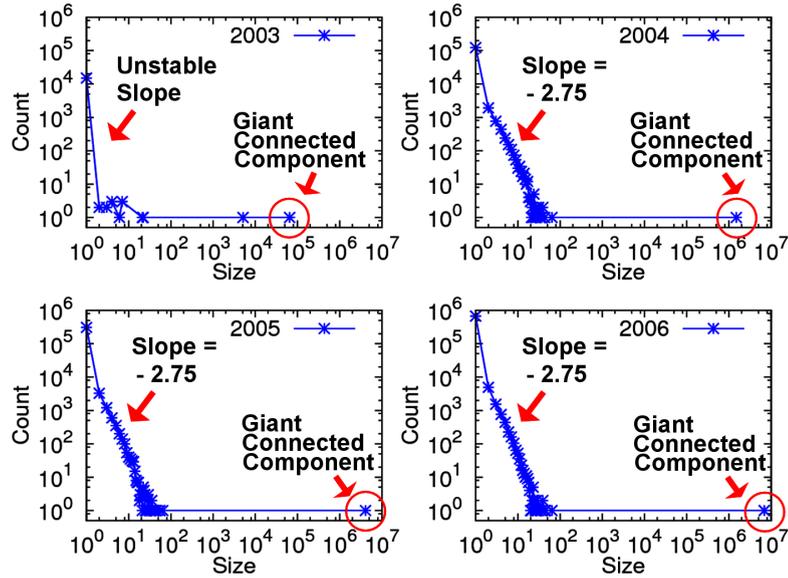


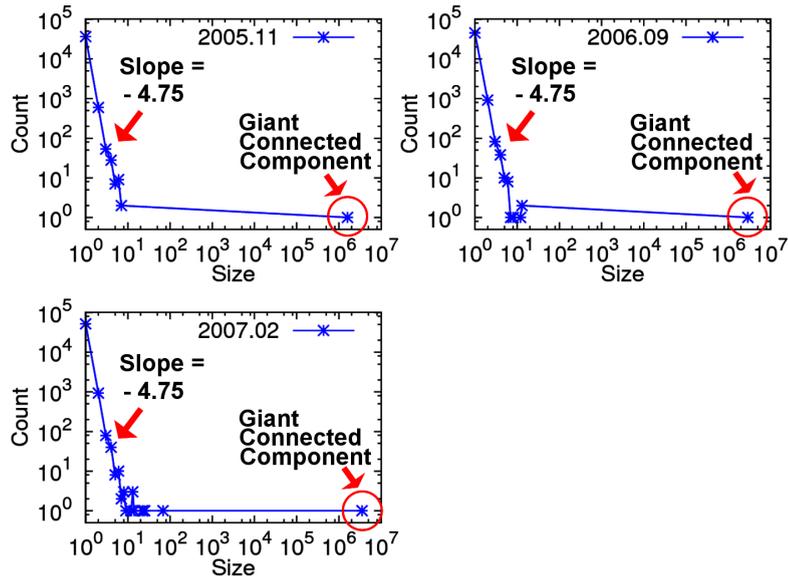
FIGURE 9.9: Number of iterations vs. the minimum vertex of Wikipedia, for connected components. D_i represents the vertex with i -th largest degree. Notice that the number of iterations decreased by 17 % after renumbering.

Absorbed Connected Components and Dunbar’s number The size of the giant component keeps growing while the second and third largest connected components do not grow beyond size 100, until they are absorbed from the giant component. This does not surprise us, since had we had two giant components it is not unlikely that some new vertex becomes connected to both.

“Anomalous” Connected Components In Figure 9.11, we see two spikes. In the first spike at size 300, more than half of the components have exactly the same structure and were made from a domain selling company where each component represents a domain to be sold. The spike happened because the company replicated sites using the same template, and injected the disconnected components into WWW network. In the second



(a) Connected Components of LinkedIn



(b) Connected Components of Wikipedia

FIGURE 9.10: The evolution of connected components. (a) The giant connected component grows for each year. However, the second largest connected component do not grow above Dunbar’s number(≈ 150) and the slope of the size distribution remains constant after the gelling point at year 2003. As in LinkedIn, notice the growth of giant connected component and the constant slope of the size distribution.

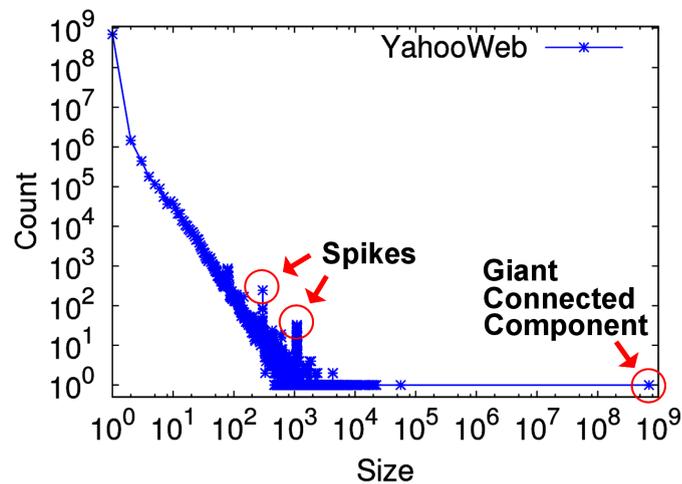


FIGURE 9.11: Connected Components of YahooWeb. Notice the two anomalous spikes which are far from the constant-slope line. Most of them are domain selling or porn sites which are replicated from templates.

spike at size 1101, more than 80 % of the components are porn sites disconnected from the giant connected component. In general, by looking carefully the distribution plot of connected components, we were able to detect interesting communities with special purposes which are disconnected from the rest of the Internet.

9.5.2 PageRank scores of Real Networks

We analyze the PageRank scores of real graphs, using PEGASUS. Figure 9.12 and 9.13 show the distribution of the PageRank scores for the Web graphs, and Figure 9.14 shows the evolution of PageRank scores for the LinkedIn and Wikipedia graphs. We observe power-law relations between the PageRank score and the number of vertices with such PageRank. The top 3 highest PageRank sites for the year 2002 are `www.careerbank.com`, `access.adobe.com`, and `top100.rambler.ru`. As expected, they have huge in-degrees (from $\approx 70\text{K}$ to $\approx 70\text{M}$).

9.5.3 Diameter of Real Networks

We analyze the diameter and radius of real networks with PEGASUS. Figure 9.15 shows the radius plot of real networks. We have following observations:

Small Diameter: For all the graphs in Figure 9.15, the average diameter is less than 6.09, verifying the six degrees of separation theory.

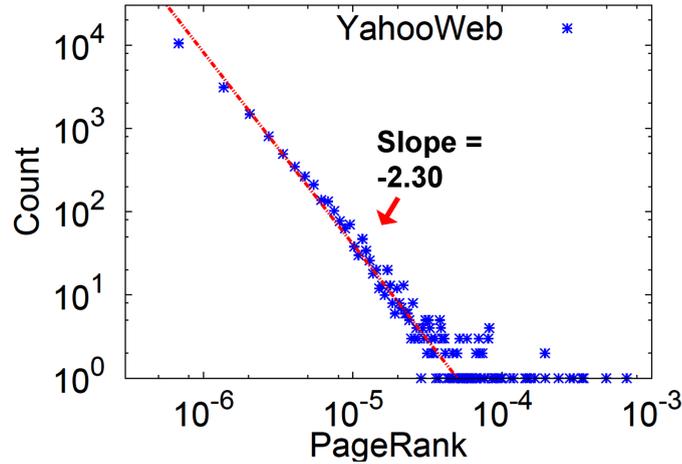


FIGURE 9.12: PageRank distribution of YahooWeb. The distribution follows a power law with an exponent 2.30.

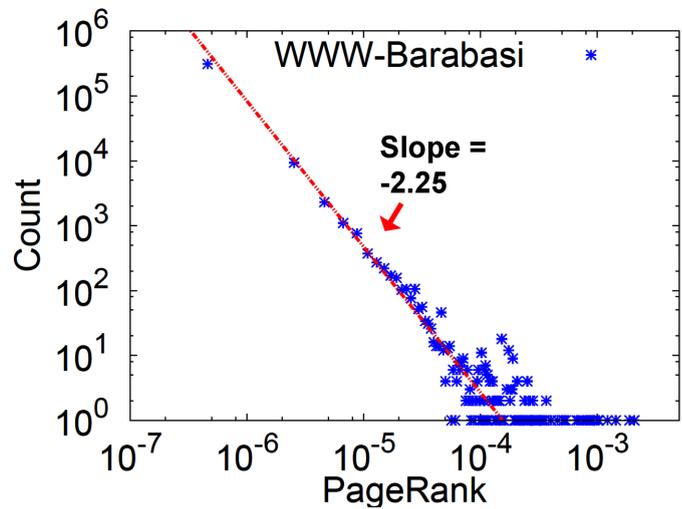
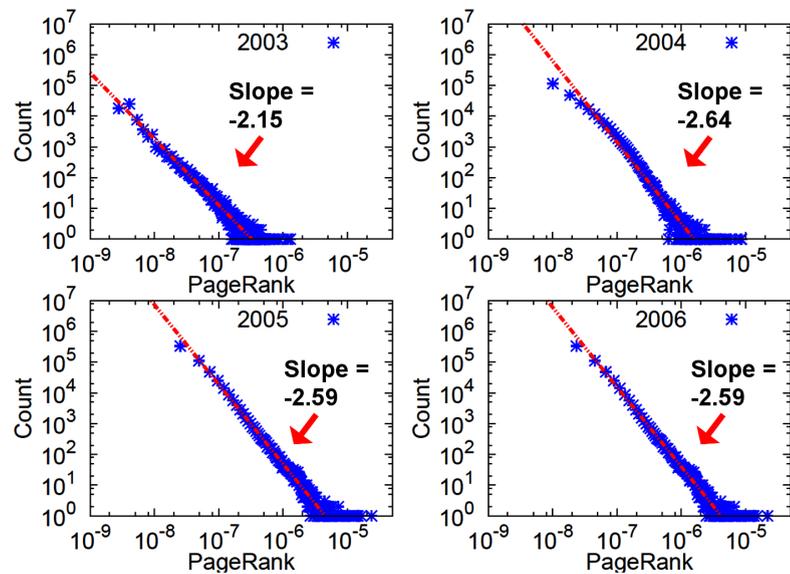


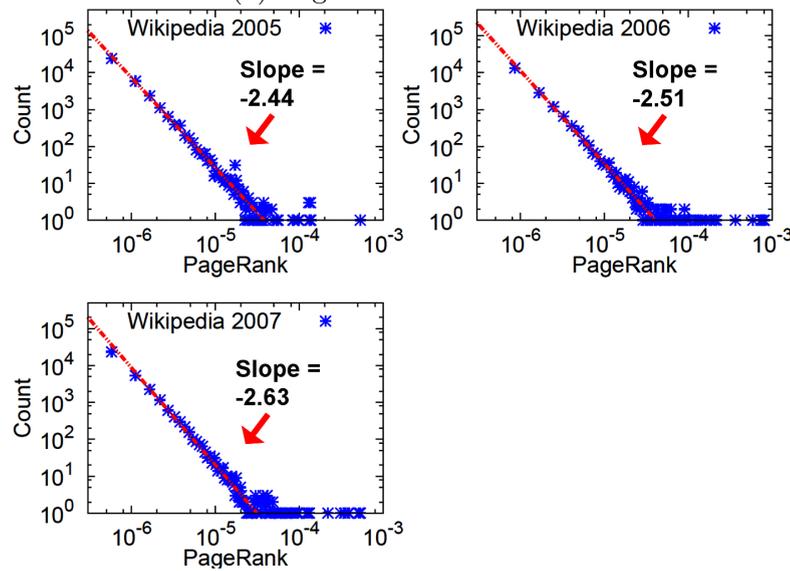
FIGURE 9.13: PageRank distribution of WWW-Barabasi. The distribution follows a power law with an exponent 2.25.

Small Changes in the Diameter over Time: For the LinkedIn graph, the average diameter remains in the range of 5.28 and 6.09 for all snapshots. For the Wikipedia graph, the average diameter remains in the range of 4.76 and 4.99 for all snapshots. Also, we do not observe a monotone pattern in the growth.

Bimodal Structure of Radius Plot For every plot, we observe a bimodal shape which reflects the structure of these real graphs. The graphs have one giant connected component where majority of vertices belongs to, and many smaller connected components whose size follows a power law. Therefore, the first mode is at radius zero which comes from single-vertex components; the second mode (e.g., at radius 6 in Epinion) comes from the giant connected component.



(a) PageRanks of LinkedIn



(b) PageRanks of Wikipedia

FIGURE 9.14: The evolution of PageRanks. (a) The distributions of PageRanks follows a power-law. However, the exponent at year 2003, which is around the gelling point, is much different from year 2004, which are after the gelling point. The exponent increases after the gelling point and becomes stable. Also notice the maximum PageRank after the gelling point is about 10 times larger than that before the gelling point due to the emergence of the giant connected component. (b) Again, the distributions of PageRanks follows a power-law. Since the gelling point is before year 2005, the three plots show similar characteristics: the maximum PageRanks and the slopes are similar.

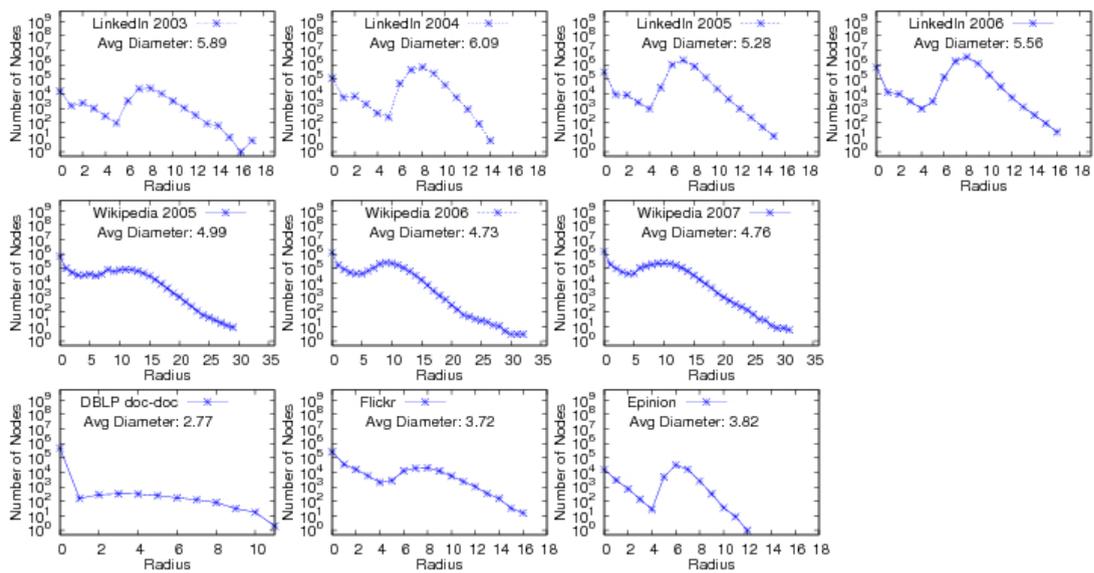


FIGURE 9.15: Radius of real graphs. X axis: radius in linear scale. Y axis: number of vertices in log scale. (Row 1) LinkedIn from 2003 to 2006. (Row 2) Wikipedia from 2005 to 2007. (Row 3) DBLP, flickr, Epinion. Notice that all the radius plots have the bimodal structure due to many smaller connected components (first mode) and the giant connected component (second mode).

II Computational Cancer Biology

Chapter 10

Approximation Algorithms for Speeding up Dynamic Programming and Denoising aCGH data

10.1 Problem & Formulation

Our approach to the aCGH denoising problem, see Section 1.2.1, is based on the well-established observation that near-by probes tend to have the same DNA copy number. We formulate the problem of denoising aCGH data as the problem of approximating a signal P with another signal F consisting of a few piecewise constant segments. Specifically, let $P = (P_1, P_2, \dots, P_n) \in \mathbb{R}^n$ be the input signal -in our setting the sequence of the noisy aCGH measurements- and let C be a constant. Our goal is to find a function $F : [n] \rightarrow \mathbb{R}$ which optimizes the following objective function:

$$\min_F \sum_{i=1}^n (P_i - F_i)^2 + C \times (|\{i : F_i \neq F_{i+1}\}| + 1). \quad (10.1)$$

The best known exact algorithm for solving the optimization problem defined by Equation 10.1 runs in $O(n^2)$ time but as our results suggest this running time is likely not to be tight. It is worth noting that existing techniques for speeding up dynamic programming [148, 149, 423] do not apply to our problem.

The main algorithmic contributions of this Chapter are two approximation algorithms for solving this recurrence. The first algorithm approximates the objective within an

additive error ϵ which we can make as small as we wish and its key idea is the reduction of the problem to halfspace range queries, a well studied computational geometric problem [14]. The second algorithm carefully decomposes the problem into a “small” (logarithmic) number of subproblems which satisfy the quadrangle inequality (Monge property).

The remainder of this Chapter is organized as follows: Section 10.2 presents the vanilla dynamic programming algorithm which runs in $O(n^2)$. Section 10.3 analyzes properties of the recurrence which will be exploited in Sections 10.4 and 10.5 where we describe the additive and multiplicative approximation algorithms respectively. In Section 10.6 we validate our model by performing an extensive biological analysis of the findings of our segmentation.

10.2 $O(n^2)$ Dynamic Programming Algorithm

In order to solve the optimization problem defined by Equation 10.1, we define the key quantity OPT_i given by the following recurrence:

$$\begin{aligned} OPT_0 &= 0 \\ OPT_i &= \min_{0 \leq j \leq i-1} [OPT_j + w(i, j)] + C, \text{ for } i > 0 \\ \text{where } w(i, j) &= \sum_{k=j+1}^i \left(P_k - \frac{\sum_{m=j+1}^i P_m}{i-j} \right)^2. \end{aligned}$$

The above recurrence has a straightforward interpretation: OPT_i is equal to the minimum cost of fitting a set of piecewise constant segments from point P_1 to P_i given that index j is a breakpoint. The cost of fitting the segment from $j+1$ to i is C . The weight function $w()$ is the minimum squared error for fitting a constant segment on points $\{P_{j+1}, \dots, P_i\}$, which is obtained for the constant segment with value $\frac{\sum_{m=j+1}^i P_m}{i-j}$, i.e., the average of the points in the segment. This recursion directly implies a simple dynamic programming algorithm. We call this algorithm CGHTRIMMER and the pseudocode is shown in Algorithm 13. The main computational bottleneck of CGHTRIMMER is the computation of the auxiliary matrix M , an upper diagonal matrix for which m_{ij} is the minimum squared error of fitting a segment from points $\{P_i, \dots, P_j\}$. To avoid a naive algorithm that would simply find the average of those points and then compute the squared error, resulting in $O(n^3)$ time, we use Claim 6.

Algorithm 13 CGHTRIMMER algorithm

Input: Signal $P = (P_1, \dots, P_n)$, Regularization parameter C
Output: Optimal Segmentation with respect to our objective (see Equation 10.1)

{ *Compute an $n \times n$ matrix M , where $M_{ji} = \sum_{k=j}^i \left(P_k - \frac{\sum_{m=j}^i P_m}{i-j+1} \right)^2$.

A is an auxiliary matrix of averages, i.e., $A_{ji} = \frac{\sum_{k=j}^i P_k}{i-j+1}$. *}

Initialize matrix $A \in \mathbb{R}^{n \times n}$, $A_{ij} = 0, i \neq j$ and $A_{ii} = P_i$.

for $i = 1$ to n **do**

for $j = i + 1$ to n **do**

$A_{i,j} \leftarrow \frac{j-i}{j-i+1} A_{i,j-1} + \frac{1}{j-i+1} P_j$

end for

end for

for $i = 1$ to n **do**

for $j = i + 1$ to n **do**

$M_{i,j} \leftarrow M_{i,j-1} + \frac{j-i}{j-i+1} (P_j - A_{i,j-1})^2$

end for

end for

{ * Solve the Recurrence.*}

for $i = 1$ to n **do**

$OPT_i \leftarrow \min_{0 \leq j \leq i-1} OPT_j + M_{j+1,i} + C$

$BREAK_i \leftarrow \arg \min_{1 \leq j \leq i} OPT_{j-1} + M_{j,i} + C$

end for

Claim 6. Let $\alpha_{(j)}$ and $m_{(j)}$ be the average and the minimum squared error of fitting a constant segment to points $\{P_1, \dots, P_j\}$ respectively. Then,

$$\alpha_{(j)} = \frac{j-1}{j} \alpha_{(j-1)} + \frac{1}{j} P_j, \quad (10.2)$$

$$m_{(j)} = m_{(j-1)} + \frac{j-1}{j} (P_j - \alpha_{(j-1)})^2. \quad (10.3)$$

The interested reader can find a proof of Claim 6 in [256]. Equations 10.2 and 10.3 provide us a way to compute means and least squared errors online. Algorithm 13 first computes matrices A and M using Equations 10.2, 10.3 and then iterates (last *for* loop) to solve the recurrence by finding the optimal breakpoint for each index i . The total running time is $O(n^2)$ (matrices A and M matrices have $O(n^2)$ entries and each requires $O(1)$ time to compute). Obviously, Algorithm 13 uses $O(n^2)$ units of space.

10.3 Analysis of The Transition Function

In the following, let $S_i = \sum_{j=1}^i P_j$. The transition function for the dynamic programming for $i > 0$ can be rewritten as:

$$OPT_i = \min_{j < i} OPT_j + \sum_{m=j+1}^i P_m^2 - \frac{(S_i - S_j)^2}{i - j} + C. \quad (10.4)$$

The transition can be viewed as a weight function $w(j, i)$ that takes the two indices j and i as parameters such that:

$$w(j, i) = \sum_{m=j+1}^i P_m^2 - \frac{(S_i - S_j)^2}{i - j} + C \quad (10.5)$$

Note that the weight function does not have the Monge property, as demonstrated by the vector $P = (P_1, \dots, P_{2k+1}) = (1, 2, 0, 2, 0, 2, 0, \dots, 2, 0, 1)$. When $C = 1$, the optimal choices of j for $i = 1, \dots, 2k$ are $j = i - 1$, i.e., we fit one segment per point. However, once we add in $P_{2k+1} = 1$ the optimal solution changes to fitting all points on a single segment. Therefore, preferring a transition to j_1 over one to j_2 at some index i does not allow us to discard j_2 from future considerations. This is one of the main difficulties in applying techniques based on the increasing order of optimal choices of j , such as the method of Eppstein, Galil and Giancarlo [148] or the method of Larmore and Schieber [268], to reduce the complexity of the $O(n^2)$ algorithm we described in Section 10.2.

We start by defining DP_i for $i = 0, 1, \dots, n$, the solution to a simpler optimization problem.

Definition 10.1. Let $DP_i, i = 0, 1, \dots, n$, satisfy the following recurrence

$$DP_i = \begin{cases} \min_{j < i} DP_j - \frac{(S_i - S_j)^2}{i - j} + C & \text{if } i > 0 \\ 0 & \text{if } i = 0 \end{cases} \quad (10.6)$$

The following observation stated as Lemma 10.2 plays a key role in Section 10.4.

Lemma 10.2. For all i , OPT_i can be written in terms of DP_i as

$$OPT_i = DP_i + \sum_{m=1}^i P_m^2.$$

Proof. We use strong induction on i . For $i = 0$ the result trivially holds. Let the result hold for all $j < i$. Then,

$$\begin{aligned}
 DP_i &= \min_{j < i} DP_j - \frac{(S_i - S_j)^2}{i - j} + C \\
 &= \min_{j < i} OPT_j - \frac{(S_i - S_j)^2}{i - j} + \sum_{m=j+1}^i P_m^2 - \sum_{m=1}^i P_m^2 + C \\
 &= OPT_i - \sum_{m=1}^i P_m^2
 \end{aligned}$$

Hence, $OPT_i = DP_i + \sum_{m=1}^i P_m^2$ for all i . □

Observe that the second order moments involved in the expression of OPT_i are absent from DP_i . Let $\tilde{w}(j, i)$ be the shifted weight function, i.e., $\tilde{w}(j, i) = -\frac{(S_i - S_j)^2}{i - j} + C$. Clearly, $w(j, i) = \tilde{w}(j, i) + \sum_{m=i}^j P_m^2$.

10.4 Additive Approximation using Halfspace Queries

In this Section, we present a novel algorithm which runs in $\tilde{O}(n^{\frac{4}{3} + \delta} \log(\frac{U}{\epsilon}))$ time and approximates the optimal objective value within additive ϵ error. We derive the algorithm gradually in the following and upon presenting the necessary theory we provide the pseudocode (see Algorithm 2) at the end of this Section. Our proposed method uses the results of [15] as stated in Corollary 2.12 to obtain a fast algorithm for the additive approximation variant of the problem. Specifically, the algorithm initializes a 4-dimensional halfspace query data structure. The algorithm then uses binary searches to compute an accurate estimate of the value DP_i for $i = 1, \dots, n$. As errors are introduced at each term, we use \tilde{DP}_i to denote the approximate value of DP_i calculated by the binary search, and \bar{DP}_i to be the optimum value of the transition function computed by examining the approximate values \tilde{DP}_j for all $j < i$. Formally,

$$\bar{DP}_i = \min_{j < i} \left[\tilde{DP}_j - \underbrace{\frac{(S_i - S_j)^2}{i - j}}_{\tilde{w}(j, i)} \right] + C.$$

Since the binary search incurs a small additive error at each step, it remains to show that these errors accumulate in a controlled way. Theorem 10.3 states that a small error at each step suffices to give an overall good approximation. We show inductively that

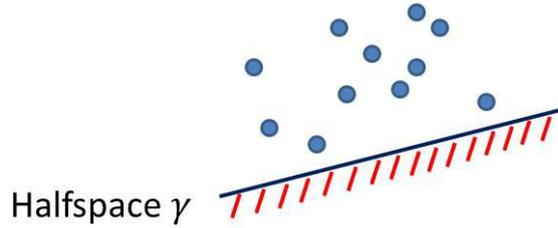


FIGURE 10.1: Answering whether or not $\bar{DP}_i \leq x + C$ reduces to answering whether the point set $\{(j, \tilde{DP}_j, 2S_j, S_j^2 + \tilde{DP}_j j) \in \mathbb{R}^4, j < i\}$ has a non-empty intersection with the halfspace $\gamma = \{y \in \mathbb{R}^4 : a_i y \leq b_i\}$ where a_i and b_i are a 4-dimensional constant vector and a constant which depend on i respectively. This type of queries can be solved efficiently, see [15].

if \tilde{DP}_i approximates \bar{DP}_i within ϵ/n , then \tilde{DP}_i is within $i\epsilon/n$ additive error from the optimal value DP_i for all i .

Theorem 10.3. *Let \tilde{DP}_i be the approximation of our algorithm to DP_i . Then, the following inequality holds:*

$$|DP_i - \tilde{DP}_i| \leq \frac{\epsilon i}{n} \tag{10.7}$$

Proof. We use induction on the number of points. Using the same notation as above, let $\bar{DP}_i = \min_{j < i} \tilde{DP}_j - w(j, i) + C$. By construction the following inequality holds:

$$|\bar{DP}_i - \tilde{DP}_i| \leq \frac{\epsilon}{n} \quad \forall i = 1, \dots, n \tag{10.8}$$

When $i = 1$ it is clear that $|DP_1 - \tilde{DP}_1| \leq \frac{\epsilon}{n}$. Our inductive hypothesis is the following:

$$|DP_j - \tilde{DP}_j| \leq \frac{j\epsilon}{n} \quad \forall j < i \tag{10.9}$$

It suffices to show that the following inequality holds:

$$|DP_i - \bar{DP}_i| \leq \frac{(i-1)\epsilon}{n} \tag{10.10}$$

since then by the triangular inequality we obtain:

$$\frac{i\epsilon}{n} \geq |DP_i - \bar{DP}_i| + |\bar{DP}_i - \tilde{DP}_i| \geq |DP_i - \tilde{DP}_i|.$$

Let j^*, \bar{j} be the optimum breakpoints for DP_i and \bar{DP}_i respectively, $j^*, \bar{j} \leq i - 1$.

$$\begin{aligned}
 DP_i &= DP_{j^*} + \tilde{w}(j^*, i) + C \\
 &\leq DP_{\bar{j}} + \tilde{w}(\bar{j}, i) + C \\
 &\leq \tilde{D}P_{\bar{j}} + \tilde{w}(\bar{j}, i) + C + \frac{\bar{j}\epsilon}{n} \text{ (by 10.9)} \\
 &= \bar{D}P_i + \frac{\bar{j}\epsilon}{n} \\
 &\leq \bar{D}P_i + \frac{(i-1)\epsilon}{n}
 \end{aligned}$$

Similarly we obtain:

$$\begin{aligned}
 \bar{D}P_i &= \tilde{D}P_{\bar{j}} + \tilde{w}(\bar{j}, i) + C \\
 &\leq \tilde{D}P_{j^*} + \tilde{w}(j^*, i) + C \\
 &\leq DP_{j^*} + \tilde{w}(j^*, i) + C + \frac{j^*\epsilon}{n} \text{ (by 10.9)} \\
 &= DP_i + \frac{j^*\epsilon}{n} \\
 &\leq DP_i + \frac{(i-1)\epsilon}{n}
 \end{aligned}$$

Combining the above two inequalities, we obtain 10.10. □

By substituting $i = n$ in Theorem 10.3 we obtain the following corollary which proves that $\tilde{D}P_n$ is within ϵ of DP_n .

Corollary 10.4. *Let $\tilde{D}P_n$ be the approximation of our algorithm to DP_n . Then,*

$$|DP_n - \tilde{D}P_n| \leq \epsilon. \tag{10.11}$$

To use the analysis above in order to come up with an efficient algorithm we need to answer two questions: (a) How many binary search queries do we need in order to obtain the desired approximation? (b) How can we answer each such query efficiently? The answer to (a) is simple: as it can easily be seen, the value of the objective function is upper bounded by U^2n , where $U = \max\{\sqrt{C}, |P_1|, \dots, |P_n|\}$. Therefore, $O(\log(\frac{U^2n}{\epsilon/n})) = \tilde{O}(\log(\frac{U}{\epsilon}))$ iterations of binary search at each index i are sufficient to obtain the desired approximation. We reduce the answer to (b) to a well studied computational geometric problem. Specifically, fix an index i , where $1 \leq i \leq n$, and consider the general form of the binary search query $\bar{D}P_i \leq x + C$, where $x + C$ is the value on which we query.

Note that we use the expression $x + C$ for convenience, i.e., so that the constant C will be simplified from both sides of the query. This query translates itself to the following decision problem, see also Figure 10.1. Does there exist an index j , such that $j < i$ and the following inequality holds:

$$x \geq \tilde{D}P_j - \frac{(S_i - S_j)^2}{i - j} \Rightarrow xi + S_i^2 \geq (x, i, S_i, -1)(j, \tilde{D}P_j, 2S_j, S_j^2 + j\tilde{D}P_j)^T?$$

Hence, the binary search query has been reduced to answering a *halfspace query*. Specifically, the decision problem for any index i becomes whether the intersection of the point set $\text{POINTS}_i = \{(j, \tilde{D}P_j, 2S_j, S_j^2 + \tilde{D}P_j j) \in \mathbb{R}^4, j < i\}$ with a hyperplane is empty. By Corollary 2.12 [15], for a point set of size n , this can be done in $\tilde{O}(n^{\frac{1}{3}+\delta})$ per query and $O(n^{\frac{1}{3}} \log n)$ amortized time per insertion of a point. Hence, the optimal value of DP_i can be found within an additive constant of ϵ/n using the binary search in $\tilde{O}(n^{\frac{1}{3}} \log(\frac{U}{\epsilon}))$ time.

Therefore, we can proceed from index 1 to n , find the approximately optimal value of OPT_i and insert a point corresponding to it into the query structure. We obtain an algorithm which runs in $\tilde{O}(n^{\frac{4}{3}+\delta} \log(\frac{U}{\epsilon}))$ time, where δ is an arbitrarily small positive constant. The pseudocode is given in Algorithm 14.

Algorithm 14 Approximation within additive ϵ using 4D halfspace queries

```

Initialize 4D halfspace query structure  $Q$ 
for  $i = 1$  to  $n$  do
   $low \leftarrow 0$ 
   $high \leftarrow nU^2$ 
  while  $high - low > \epsilon/n$  do
     $m \leftarrow (low + high)/2$ 
    { *This halfspace emptiness query is efficiently supported by Q.* }
     $flag \leftarrow (\exists j \text{ such that } xi + S_i^2 \geq (x, i, S_i, -1)(j, \tilde{D}P_j, 2S_j, S_j^2 + j\tilde{D}P_j)^T)$ 
    if  $flag$  then
       $high \leftarrow m$ 
    else
       $low \leftarrow m$ 
    end if
  end while
   $\tilde{D}P_i \leftarrow (low + high)/2$ 
  { *Point insertions are efficiently supported by Q.* }
  Insert point  $(i, \tilde{D}P_i, 2S_i, S_i^2 + \tilde{D}P_i i)$  in  $Q$ 
end for

```

TABLE 10.1: Summary of proof of Lemma 10.5.

	$w'(i_1, i_4) + w'(i_2, i_3)$	$w'(i_1, i_3) + w'(i_2, i_4)$
(S_1, S_1)	1	1
(S_2, S_2)	2	2
(S_3, S_3)	1	1
(S_1, S_2)	1	1
(S_1, S_3)	1	0
(S_2, S_3)	1	1

10.5 Multiscale Monge Decomposition

In this Section we present an algorithm which runs in $O(n \log n/\epsilon)$ time to approximate the optimal shifted objective value within a multiplicative factor of $(1+\epsilon)$. Our algorithm is based on a new technique, which we consider of independent interest. Let $w'(j, i) = \sum_{m=j+1}^i (i-j)P_m^2 - (S_i - S_j)^2$. We can rewrite the weight function w as a function of w' , namely as $w(j, i) = w'(j, i)/(i-j) + C$. The interested reader can easily check that we can express $w'(j, i)$ as a sum of non-negative terms, i.e.,

$$w'(j, i) = \sum_{j+1 \leq m_1 < m_2 \leq i} (P_{m_1} - P_{m_2})^2.$$

Recall from Section 10.2 that the weight function w is not Monge. The next lemma shows that the weight function w' is a Monge function.

Lemma 10.5. *The weight function $w'(j, i)$ is Monge (concave), i.e., for any $i_1 < i_2 < i_3 < i_4$, the following holds:*

$$w'(i_1, i_4) + w'(i_2, i_3) \geq w'(i_1, i_3) + w'(i_2, i_4).$$

Proof. Since each term in the summation is non-negative, it suffices to show that any pair of indices, (m_1, m_2) is summed as many times on the left hand side as on the right hand side. If $i_2 + 1 \leq m_1 < m_2 \leq i_3$, each term is counted twice on each side. Otherwise, each term is counted once on the left hand side since $i_1 + 1 \leq m_1 < m_2 \leq i_4$ and at most once on the right hand side since $[i_1 + 1, i_3] \cap [i_2 + 1, i_4] = [i_2 + 1, i_3]$. \square

The proof of Lemma 10.5 is summarized in Table 10.1. Specifically, let S_j be the set of indices $\{i_j + 1, \dots, i_{j+1}\}$, $j = 1, 2, 3$. Also, let (S_j, S_k) denote the set of indices (m_1, m_2) which appear in the summation such that $m_1 \in S_j, m_2 \in S_k$. The two last columns of Table 10.1 correspond to the left- and right-hand side of the Monge inequality (as in Lemma 10.5) and contain the counts of appearances of each term.

Our approach is based on the following observations:

1. Consider the weighted directed acyclic graph (DAG) on the vertex set $V = \{0, \dots, n\}$ with edge set $E = \{(j, i) : j < i\}$ and weight function $w : E \rightarrow \mathbb{R}$, i.e., edge (j, i) has weight $w(j, i)$. Solving the aCGH denoising problem reduces to finding a shortest path from vertex 0 to vertex n . If we perturb the edge weights within a factor of $(1 + \epsilon)$, as long as the weight of each edge is positive, then the optimal shortest path distance is also perturbed within a factor of at most $(1 + \epsilon)$.
2. By Lemma 10.5 we obtain that the weight function is not Monge essentially because of the $i - j$ term in the denominator.
3. Our goal is to approximate w by a Monge function w' such that $c_1 w' \leq w \leq c_2 w'$ where c_1, c_2 should be known constants.

In the following we elaborate on the latter goal. Fix an index i and note that the optimal breakpoint for that index is some index $j \in \{1, \dots, i - 1\}$. We will “bucketize” the range of index j into $m = O(\log_{1+\epsilon}(i)) = O(\log n/\epsilon)$ buckets such that the k -th bucket, $k = 1, \dots, m$, is defined by the set of indices j which satisfy

$$l_k = i - (1 + \epsilon)^k \leq j \leq i - (1 + \epsilon)^{k-1} = r_k.$$

This choice of bucketization is based on the first two observations which guide our approach. Specifically, it is easy to check that $(1 + \epsilon)^{k-1} \leq i - j \leq (1 + \epsilon)^k$. This results, for any given i , to approximating $i - j$ by a constant for each possible bucket, leading to $O(\log n/\epsilon)$ different Monge functions (one per bucket) while incurring a multiplicative error of at most $(1 + \epsilon)$. However, there exists a subtle point, as also Figure 10.2 indicates. We need to make sure that each of the Monge functions is appropriately defined so that when we consider the k -th Monge subproblem, the optimal breakpoint j_k should satisfy $j_k \in [l_k, r_k]$. Having achieved that (see Lemma 6.2) we can solve efficiently the recurrence. Specifically, OPT_i is computed as follows:

$$\begin{aligned}
 OPT_i &= \min_{j < i} \left[OPT_j + \frac{w'(i, j)}{i - j} \right] + C \\
 &= \min_k \left[\min_{j \in [l_k, r_k]} OPT_j + \frac{w'(i, j)}{i - j} \right] + C \\
 &\approx \min_k \left[\min_{j \in [l_k, r_k]} OPT_j + \frac{w'(i, j)}{(1 + \epsilon)^{k-1}} \right] + C \\
 &= \min_k \left[\min_{j \in [l_k, r_k]} OPT_j + \frac{w'(i, j)}{c_k} \right] + C.
 \end{aligned}$$

The following is one of the possible ways to define the m Monge weight functions. In what follows, λ is a sufficiently large positive constant.

$$w_k(j, i) = \begin{cases} 2^{n-i+j}\lambda & i - j < c_k = (1 + \epsilon)^{k-1} \\ 2^{i-j}\lambda & i - j > (1 + \epsilon)c_k = (1 + \epsilon)^k \\ w'(j, i)/c_k & \text{otherwise} \end{cases} \quad (10.12)$$

Lemma 10.6. *Given any vector P , it is possible to pick λ such that w_k is Monge for all $k \geq 1$. That is, for any 4-tuple $i_1 < i_2 < i_3 < i_4$, $w_k(i_1, i_4) + w_k(i_2, i_3) \geq w_k(i_1, i_3) + w_k(i_2, i_4)$.*

Proof. Since $w'(j, i) = \sum_{j+1 \leq m_1 < m_2 \leq i} (P_{m_1} - P_{m_2})^2 \leq (2K)^2 n^2$ where $K = \max_{1 \leq i \leq n} |P_i|$, we can pick λ such that $w_k(j, i) \geq w'(j, i)$. The rest of the proof is casework based on the lengths of the intervals $i_3 - i_1$, $i_4 - i_2$, and how they compare with c_k and $(1 + \epsilon)c_k$. There are 12 such cases in total. We may assume $i_3 - i_1 \leq i_4 - i_2$ without loss of generality, leaving thus 6 cases to be considered.

If $c_k \leq i_3 - i_1, i_4 - i_2 \leq (1 + \epsilon)c_k$, then:

$$\begin{aligned}
 w_k(i_1, i_3) + w_k(i_2, i_4) &= (1/c_k)(w'(i_1, i_3) + w'(i_2, i_4)) \\
 &\leq (1/c_k)(w'(i_1, i_4) + w'(i_2, i_3)) \quad (\text{by Lemma 10.5}) \\
 &\leq w_k(i_1, i_4) + w_k(i_2, i_3).
 \end{aligned}$$

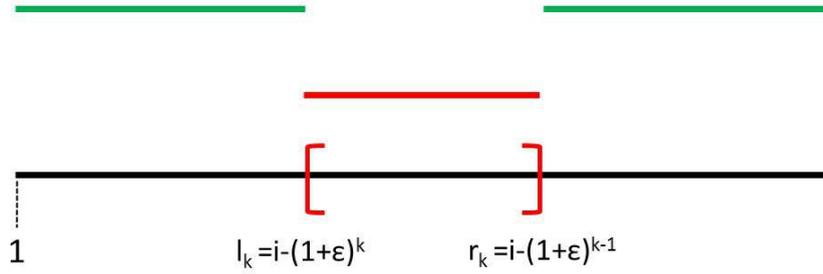


FIGURE 10.2: Solving the k -th Monge problem, $k = 1, \dots, m = O(\log_{1+\epsilon}(i))$, for a fixed index i . The k -th interval is the set of indices $\{j : l_k \leq j \leq r_k, l_k = i - (1 + \epsilon)^k, r_k = i - (1 + \epsilon)^{k-1}\}$. Ideally, we wish to define a Monge function w'_k whose maximum inside the k -th interval (red color) is smaller than the minimum outside that interval (green color). This ensures that the optimal breakpoint for the k -th Monge subproblem lies inside the red interval.

Consider the case $i_3 - i_1 < c_k$ and $i_4 - i_2 \leq (1 + \epsilon)c_k$. Then as $i_2 > i_1$, $i_3 - i_2 \leq i_3 - i_1 - 1$ and we have:

$$\begin{aligned} w_k(i_2, i_3) &= 2^{n-i_3+i_2} \lambda \\ &\geq 2 \cdot 2^{n-i_3-i_1} \lambda \\ &\geq w_k(i_1, i_3) + w_k(i_2, i_4) \end{aligned}$$

The cases of $c_k \leq i_3 - i_1$ and $c_k(1 + \epsilon) < i_4 - i_2$ can be done similarly. Note that the cases of $i_3 - i_1, i_4 - i_2 < c_k$ and $(1 + \epsilon)c_k < i_3 - i_1, i_4 - i_2$ are also covered by these.

The only case that remain is $i_3 - i_1 < c_k$ and $(1 + \epsilon)c_k < i_4 - i_2$. Since $i_3 - i_2 < i_3 - i_1 < c_k$, we have:

$$\begin{aligned} w_k(i_2, i_3) &= 2^{n-i_3+i_2} \lambda \\ &> 2^{n-i_4+i_2} \lambda \\ &= w_k(i_2, i_4) \end{aligned}$$

Similarly $w_k(i_1, i_4) = 2^{i_4-i_1} \lambda > 2^{i_3-i_1} \lambda = w_k(i_1, i_3)$. Adding them gives the desired result. \square

The pseudocode is shown in Algorithm 15. The algorithm computes the OPT values online based on the above analysis. In order to solve each Monge sub-problem our method calls the routine of Theorem 2.13. For each index i , we compute OPT_i by taking the best value over queries to all k of the Monge query structures, then we update all the structures with this value. Note that storing values of the form $2^k \lambda$ using only their exponent k suffices for comparison, so introducing $w_k(j, i)$ doesn't result in

any change in runtime. By Theorem 2.13, for each Q_k , finding $\min_{j < i} Q_k \cdot a_j + w_k(j, i)$ over all i takes $O(n)$ time. Hence, the total runtime is $O(n \log n / \epsilon)$.

Algorithm 15 Approximation within a factor of ϵ using Monge function search

Maintain $m = \log n / \log(1 + \epsilon)$ Monge function search data structures Q_1, \dots, Q_m where Q_k corresponds to the Monge function $w_k(j, i)$.

{*The weight function $w_k(j, i)$ is Monge. Specifically, $w_k(j, i) = C + \left(\sum_{m=j+1}^i (i - j) P_m^2 - \frac{(S_i - S_j)^2}{(1 + \epsilon)^k} \right)$ for all j which satisfy $(1 + \epsilon)^{k-1} \leq i - j \leq (1 + \epsilon)^k$. *}

$OPT_0 \leftarrow 0$ { * Recursion basis. * }

for $k = 1$ to m **do**

$Q_k \cdot a_0 \leftarrow 0$

end for

{ * Let $a_j^{(k)}$ denote $Q_k \cdot a_j$, i.e., the value a_j of the k -th data structure Q_k . * }

for $i = 1$ to n **do**

$OPT_i \leftarrow \infty$

for $k = 1$ to m **do**

$\text{localmin}_k \leftarrow \min_{j < i} a_j^{(k)} + w_k(j, i)$

$OPT_i \leftarrow \min\{OPT_i, \text{localmin}_k + C\}$

end for

end for

for $k = 1$ to m **do**

$Q_k \cdot a_i \leftarrow OPT_i$

end for

10.6 Validation of Our Model

In this Section we validate our model using the exact algorithm, see Section 10.2. In Section 10.6.1 we describe the datasets and the experimental setup. In Section 10.6.2 we show the findings of our method together with a detailed biological analysis.

10.6.1 Experimental Setup and Datasets

Our code is implemented in MATLAB¹. The experiments run in a 4GB RAM, 2.4GHz Intel(R) Core(TM)2 Duo CPU, Windows Vista machine. Our methods were compared to existing MATLAB implementations of the CBS algorithm, available via the Bioinformatics toolbox, and the CGHSEG algorithm [339], courteously provided to us by Franc Picard. CGHSEG was run using heteroscedastic model under the Lavielle criterion [271]. Additional tests using the homoscedastic model showed substantially worse performance

¹Code available at URL <http://www.math.cmu.edu/~ctsourak/CGHTRIMMER.zip>. Faster C code is also available, but since the competitors were implemented in MATLAB, all the results in this Section refer to our MATLAB implementation.

TABLE 10.2: Datasets, papers and the URLs where the datasets can be downloaded. \odot and \blacksquare denote which datasets are synthetic and real respectively.

	Dataset	Availability
\odot	Lai et al.	[263] http://compbio.med.harvard.edu/
\odot	Willenbrock et al.	[415] http://www.cbs.dtu.dk/~hanni/aCGH/
\blacksquare	Coriell Cell lines	[365] http://www.nature.com/ng/journal/v29/n3/
\blacksquare	Berkeley Breast Cancer	[313] http://icbp.lbl.gov/breastcancer/

and are omitted here. All methods were compared using previously developed benchmark datasets, shown in Table 10.2. Follow-up analysis of detected regions was conducted by manually searching for significant genes in the Genes-to-Systems Breast Cancer Database <http://www.itb.cnr.it/breastcancer> [404] and validating their positions with the UCSC Genome Browser <http://genome.ucsc.edu/>. The Atlas of Genetics and Cytogenetics in Oncology and Haematology <http://atlasgeneticsoncology.org/> was also used to validate the significance of reported cancer-associated genes. It is worth pointing out that since aCGH data are typically given in the log scale, we first exponentiate the points, then fit the constant segment by taking the average of the exponentiated values from the hypothesized segment, and then return to the log domain by taking the logarithm of that constant value. Observe that one can fit a constant segment by averaging the log values using Jensen’s inequality, but we favor an approach more consistent with the prior work, which typically models the data assuming i.i.d. Gaussian noise in the linear domain.

How to pick C ? The performance of our algorithm depends on the value of the parameter C , which determines how much each segment “costs.” Clearly, there is a tradeoff between larger and smaller values: excessively large C will lead the algorithm to output a single segment while excessively small C will result in each point being fit as its own segment. We pick our parameter C using data published in [415]. The data was generated by modeling real aCGH data, thus capturing their nature better than other simplified synthetic data and also making them a good training dataset for our model. We used this dataset to generate a Receiver Operating Characteristic (ROC) curve using values for C ranging from 0 to 4 with increment 0.01 using one of the four datasets in [415] (“above 20”). The resulting curve is shown in Figure 10.3. Then, we selected $C = 0.2$, which achieves high precision/specificity (0.98) and high recall/sensitivity (0.91). All subsequent results reported were obtained by setting C equal to 0.2.

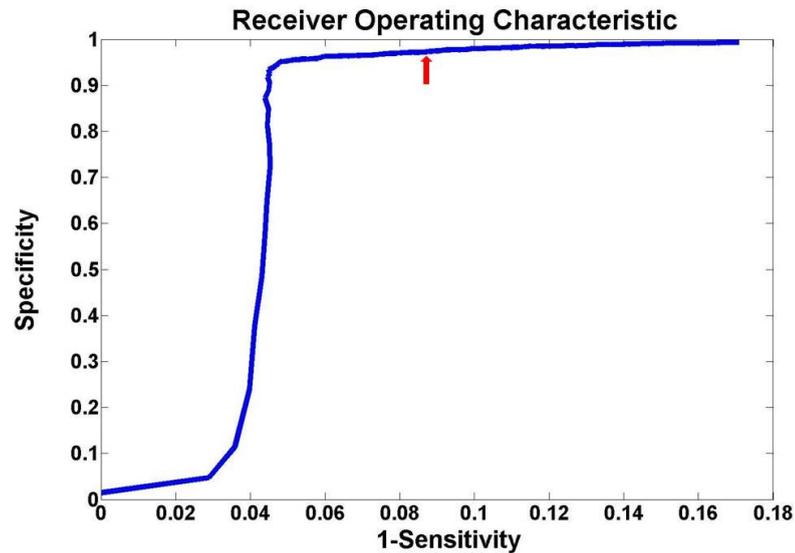


FIGURE 10.3: ROC curve of CGHTRIMMER as a function of C on data from [415]. The red arrow indicates the point (0.91 and 0.98 recall and precision respectively) corresponding to $C=0.2$, the value used in all subsequent results.

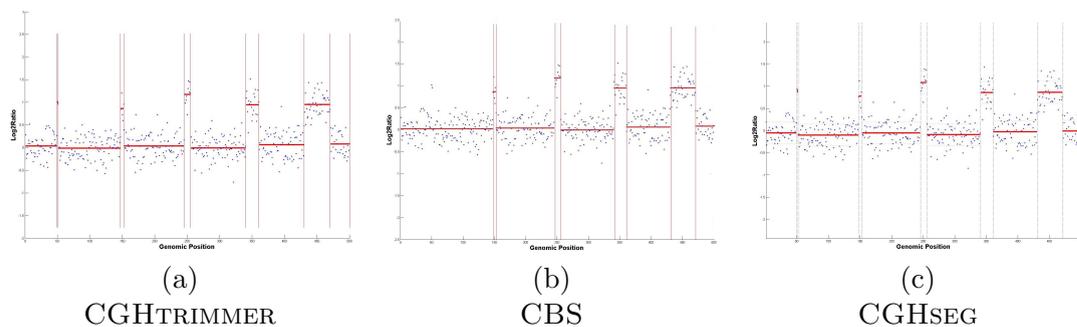


FIGURE 10.4: Performance of CGHTRIMMER, CBS, and CGHSEG on denoising synthetic aCGH data from [263]. CGHTRIMMER and CGHSEG exhibit excellent precision and recall whereas CBS misses two consecutive genomic positions with DNA copy number equal to 3.

10.6.2 Experimental Results and Biological Analysis

We show the results on synthetic data in Section 10.6.2.1, on real data where the ground truth is available to us in Section 10.6.2.2 and on breast cancer cell lines with no ground truth in Section 10.6.2.3.

10.6.2.1 Synthetic Data

We use the synthetic data published in [263]. The data consist of five aberrations of increasing widths of 2, 5, 10, 20 and 40 probes, respectively, with Gaussian noise $N(0,0.25^2)$. Figure 10.4 shows the performance of CGHTRIMMER, CBS, and CGHSEG. Both CGHTRIMMER and CGHSEG correctly detect all aberrations, while CBS misses

the first, smallest region. The running time for CGHTRIMMER is 0.007 sec, compared to 1.23 sec for CGHSEG and 60 sec for CBS.

10.6.2.2 Coriell Cell Lines

The first real dataset we use to evaluate our method is the Coriell cell line BAC array CGH data [365], which is widely considered a “gold standard” dataset. The dataset is derived from 15 fibroblast cell lines using the normalized average of \log_2 fluorescence relative to a diploid reference. To call gains or losses of inferred segments, we assign to each segment the mean intensity of its probes and then apply a simple threshold test to determine if the mean is abnormal. We follow [61] in favoring ± 0.3 out of the wide variety of thresholds that have been used [319].

Table 10.3 summarizes the performance of CGHTRIMMER, CBS and CGHSEG relative to previously annotated gains and losses in the Coriell dataset. The table shows notably better performance for CGHTRIMMER compared to either alternative method. CGHTRIMMER finds 22 of 23 expected segments with one false positive. CBS finds 20 of 23 expected segments with one false positive. CGHSEG finds 22 of 23 expected segments with seven false positives. CGHTRIMMER thus achieves the same recall as CGHSEG while outperforming it in precision and the same precision as CBS while outperforming it in recall. In cell line GM03563, CBS fails to detect a region of two points which have undergone a loss along chromosome 9, in accordance with the results obtained using the Lai et al. [263] synthetic data. In cell line GM03134, CGHSEG makes a false positive along chromosome 1 which both CGHTRIMMER and CBS avoid. In cell line GM01535, CGHSEG makes a false positive along chromosome 8 and CBS misses the aberration along chromosome 12. CGHTRIMMER, however, performs ideally on this cell line. In cell line GM02948, CGHTRIMMER makes a false positive along chromosome 7, finding a one-point segment in 7q21.3d at genomic position 97000 whose value is equal to 0.732726. All other methods also make false positive errors on this cell line. In GM7081, all three methods fail to find an annotated aberration on chromosome 15. In addition, CGHSEG finds a false positive on chromosome 11.

CGHTRIMMER also substantially outperforms the comparative methods in run time, requiring 5.78 sec for the full data set versus 8.15 min for CGHSEG (an 84.6-fold speedup) and 47.7 min for CBS (a 495-fold speedup).

TABLE 10.3: Results from applying CGHTRIMMER, CBS, and CGHSEG to 15 cell lines. Rows with listed chromosome numbers (e.g., GM03563/3) corresponded to known gains or losses and are annotated with a check mark if the expected gain or loss was detected or a “No” if it was not. Additional rows list chromosomes on which segments not annotated in the benchmark were detected; we presume these to be false positives.

Cell Line/Chromosome	CGHTRIMMER	CBS	CGHSEG
GM03563/3	✓	✓	✓
GM03563/9	✓	No	✓
GM03563/False	-	-	-
GM00143/18	✓	✓	✓
GM00143/False	-	-	-
GM05296/10	✓	✓	✓
GM05296/11	✓	✓	✓
GM05296/False	-	-	4,8
GM07408/20	✓	✓	✓
GM07408/False	-	-	-
GM01750/9	✓	✓	✓
GM01750/14	✓	✓	✓
GM01750/False	-	-	-
GM03134/8	✓	✓	✓
GM03134/False	-	-	1
GM13330/1	✓	✓	✓
GM13330/4	✓	✓	✓
GM13330/False	-	-	-
GM03576/2	✓	✓	✓
GM03576/21	✓	✓	✓
GM03576/False	-	-	-
GM01535/5	✓	✓	✓
GM01535/12	✓	No	✓
GM01535/False	-	-	8
GM07081/7	✓	✓	✓
GM07081/15	No	No	No
GM07081/False	-	-	11
GM02948/13	✓	✓	✓
GM02948/False	7	1	2
GM04435/16	✓	✓	✓
GM04435/21	✓	✓	✓
GM04435/False	-	-	8,17
GM10315/22	✓	✓	✓
GM10315/False	-	-	-
GM13031/17	✓	✓	✓
GM13031/False	-	-	-
GM01524/6	✓	✓	✓
GM01524/False	-	-	-

10.6.2.3 Breast Cancer Cell Lines

To illustrate further the performance of CGHTRIMMER and compare it to CBS and CGHSEG, we applied it to the Berkeley Breast Cancer cell line database [313]. The dataset consists of 53 breast cancer cell lines that capture most of the recurrent genomic and transcriptional characteristics of 145 primary breast cancer cases. We do not have an accepted “answer key” for this data set, but it provides a more extensive basis for detailed comparison of differences in performance of the methods on common data sets, as well as an opportunity for novel discovery. While we have applied the methods to all chromosomes in all cell lines, space limitations prevent us from presenting the full results here. The interested reader can reproduce all the results including the ones not presented here². We therefore arbitrarily selected three of the 53 cell lines and selected three chromosomes per cell line that we believed would best illustrate the comparative performance of the methods. The Genes-to-Systems Breast Cancer Database³ [404] was used to identify known breast cancer markers in regions predicted to be gained or lost by at least one of the methods. We used the UCSC Genome Browser⁴ to verify the placement of genes.

We note that CGHTRIMMER again had a substantial advantage in run time. For the full data set, CGHTRIMMER required 22.76 sec, compared to 23.3 min for CGHSEG (a 61.5-fold increase), and 4.95 hrs for CBS (a 783-fold increase).

Cell Line BT474: Figure 10.5 shows the performance of each method on the BT474 cell line. The three methods report different results for chromosome 1, as shown in Figures 10.5(a,b,c), with all three detecting amplification in the q-arm but differing in the detail of resolution. CGHTRIMMER is the only method that detects region 1q31.2-1q31.3 as aberrant. This region hosts gene NEK7, a candidate oncogene [250] and gene KIF14, a predictor of grade and outcome in breast cancer [126]. CGHTRIMMER and CBS annotate the region 1q23.3-1q24.3 as amplified. This region hosts several genes previously implicated in breast cancer [404], such as CREG1 (1q24), POU2F1 (1q22-23), RCSD1 (1q22-q24), and BLZF1 (1q24). Finally, CGHTRIMMER alone reports independent amplification of the gene CHRM3, a marker of metastasis in breast cancer patients [404].

For chromosome 5 (Figures 10.5(d,e,f)), the behavior of the three methods is almost identical. All methods report amplification of a region known to contain many breast cancer markers, including MRPL36 (5p33), ADAMTS16 (5p15.32), POLS (5p15.31), ADCY2

²<http://www.math.cmu.edu/~ctsourak/DPaCGHsupp.html>

³<http://www.itb.cnr.it/breastcancer>

⁴<http://genome.ucsc.edu/>

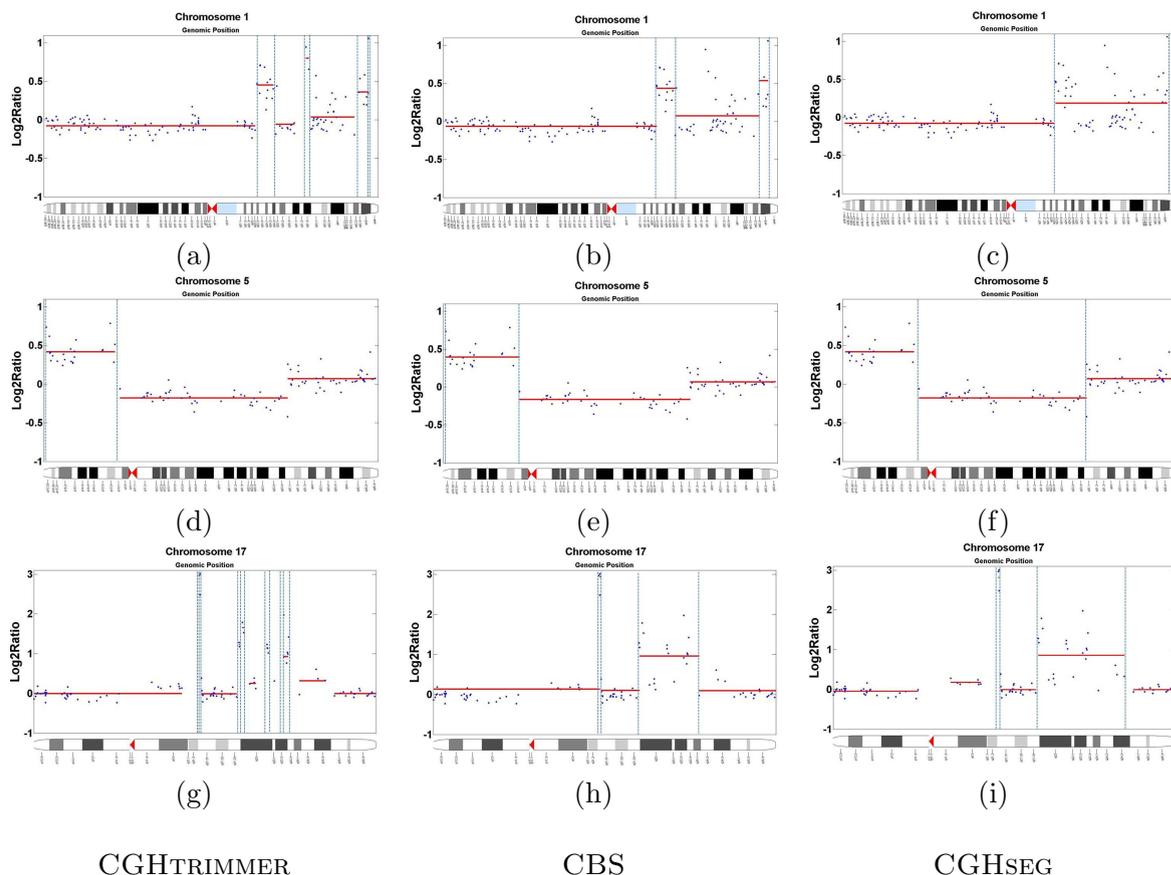


FIGURE 10.5: Visualization of the segmentation output of CGHTRIMMER, CBS, and CGHSEG for the cell line BT474 on chromosomes 1 (a,b,c), 5 (d,e,f), and 17 (g,h,i). (a,d,g) CGHTRIMMER output. (b,e,h) CBS output. (c,f,i) CGHSEG output. Segments exceeding the ± 0.3 threshold [61] are highlighted.

(5p15.31), CCT5 (5p15.2), TAS2R1 (5p15.31), ROPN1L (5p15.2), DAP (5p15.2), ANKH (5p15.2), FBXL7 (5p15.1), BASP1 (5p15.1), CDH18 (5p14.3), CDH12 (5p14.3), CDH10 (5p14.2 - 5p14.1), CDH9 (5p14.1) PDZD2 (5p13.3), GOLPH3 (5p13.3), MTMR12 (5p13.3), ADAMTS12 (5p13.3 - 5p13.2), SLC45A2 (5p13.2), TARS (5p13.3), RAD1 (5p13.2), AGXT2 (5p13.2), SKP2 (5p13.2), NIPBL (5p13.2), NUP155 (5p13.2), KRT18P31 (5p13.2), LIFR (5p13.1) and GDNF (5p13.2) [404]. The only difference in the assignments is that CBS fits one more probe to this amplified segment.

Finally, for chromosome 17 (Figures 10.5(g,h,i)), like chromosome 1, all methods detect amplification but CGHTRIMMER predicts a finer breakdown of the amplified region into independently amplified segments. All three methods detect amplification of a region which includes the major breast cancer biomarkers HER2 (17q21.1) and BRCA1 (17q21) as also the additional markers MSI2 (17q23.2) and TRIM37 (17q23.2) [404]. While the more discontinuous picture produced by CGHTRIMMER may appear to be a less parsimonious explanation of the data, a complex combination of fine-scale gains and losses in 17q is in fact well supported by the literature [323].

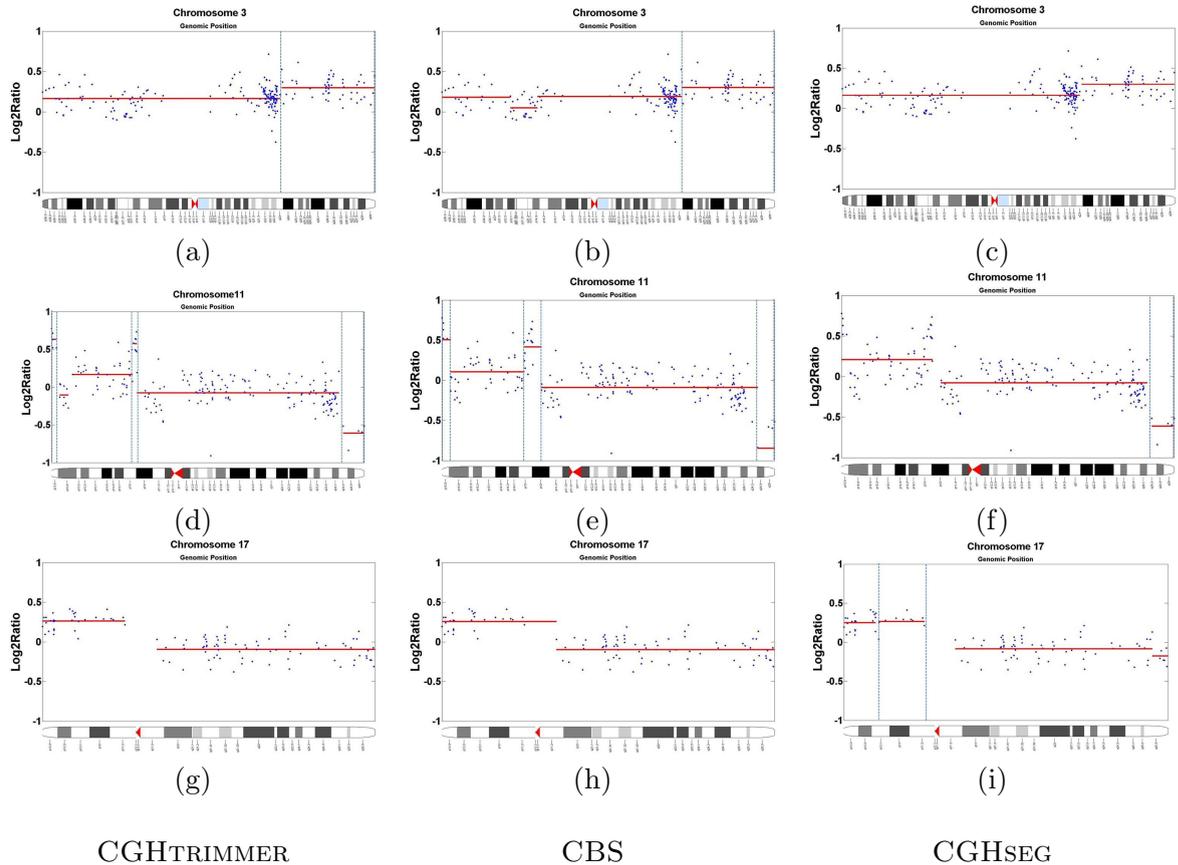


FIGURE 10.6: Visualization of the segmentation output of CGHTRIMMER, CBS, and CGHSEG for the cell line HS578T on chromosomes 3 (a,b,c), 11 (d,e,f), and 17 (g,h,i). (a,d,g) CGHTRIMMER output. (b,e,h) CBS output. (c,f,i) CGHSEG output. Segments exceeding the ± 0.3 threshold [61] are highlighted.

Cell Line HS578T: Figure 10.6 compares the methods on cell line HS578T for chromosomes 3, 11 and 17. Chromosome 3 (Figures 10.6(a,b,c)) shows identical prediction of an amplification of 3q24-3qter for all three methods. This region includes the key breast cancer markers PIK3CA (3q26.32) [273], and additional breast-cancer-associated genes TIG1 (3q25.32), MME (3q25.2), TNFSF10 (3q26), MUC4 (3q29), TFRC (3q29), DLG1 (3q29) [404]. CGHTRIMMER and CGHSEG also make identical predictions of normal copy number in the p-arm, while CBS reports an additional loss between 3p21 and 3p14.3. We are unaware of any known gain or loss in this region associated with breast cancer.

For chromosome 11 (Figures 10.6(d,e,f)), the methods again present an identical picture of loss at the q-terminus (11q24.2-11qter) but detect amplifications of the p-arm at different levels of resolution. CGHTRIMMER and CBS detect gain in the region 11p15.5, which is the site of the HRAS breast cancer metastasis marker [404]. In contrast to CBS, CGHTRIMMER detects an adjacent loss region. While we have no direct evidence this loss is a true finding, the region of predicted loss does contain EIF3F (11p15.4), identified

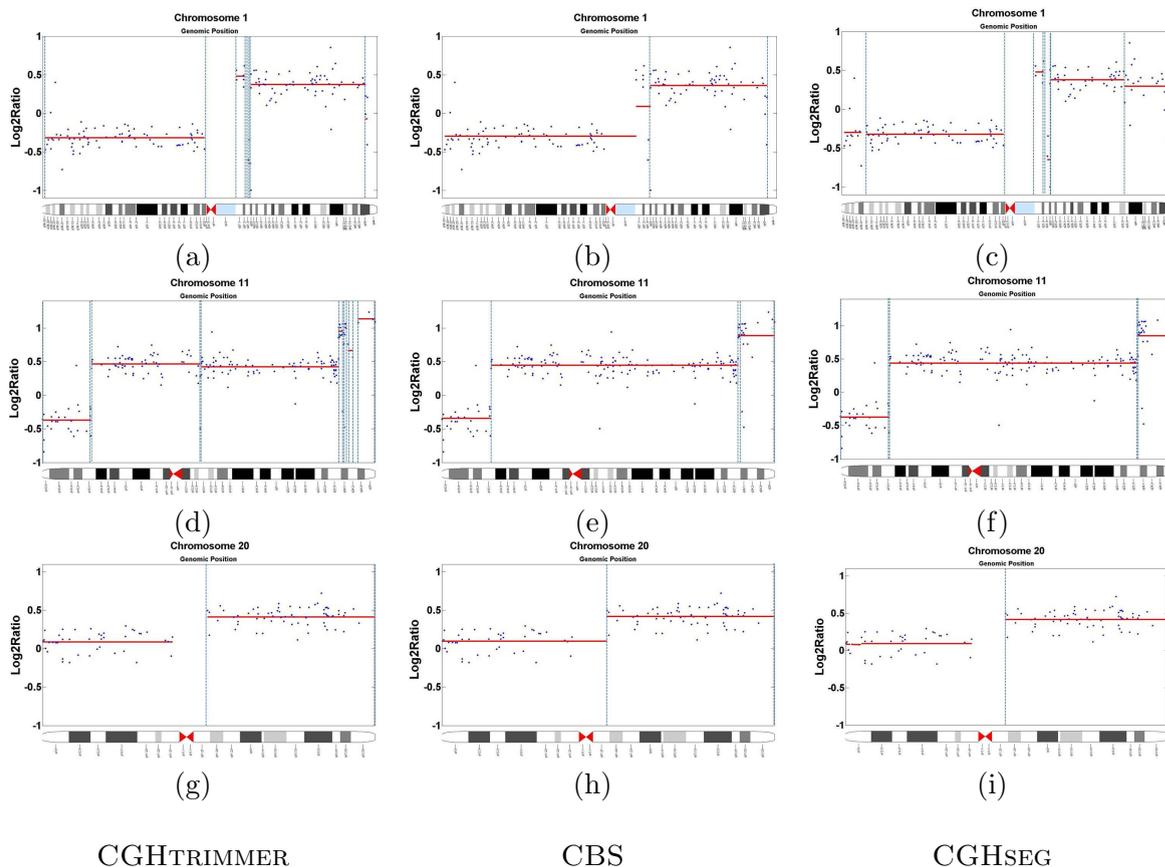


FIGURE 10.7: Visualization of the segmentation output of CGHTRIMMER, CBS, and CGHSEG for the cell line T47D on chromosomes 1 (a,b,c), 11 (d,e,f), and 20 (g,h,i). (a,d,g) CGHTRIMMER output. (b,e,h) CBS output. (c,f,i) CGHSEG output. Segments exceeding the ± 0.3 threshold [61] are highlighted.

as a possible tumor suppressor whose expression is decreased in most pancreatic cancers and melanomas [404]. Thus, we conjecture that EIF3F is a tumor suppressor in breast cancer.

On chromosome 17 (Figures 10.6(g,h,i)), the three methods behave similarly, with all three predicting amplification of the p-arm. CBS places one more marker in the amplified region causing it to cross the centromere while CGHSEG breaks the amplified region into three segments by predicting additional amplification at a single marker.

Cell Line T47D: Figure 10.7 compares the methods on chromosomes 1, 8, and 20 of the cell line T47D. On chromosome 1 (Figure 10.7(a,b,c)), all three methods detect loss of the p-arm and a predominant amplification of the q-arm. CBS infers a presumably spurious extension of the p-arm loss across the centromere into the q-arm, while the other methods do not. The main differences between the three methods appear on the q-arm of chromosome 1. CGHTRIMMER and CGHSEG detect a small region of gain proximal to the centromere at 1q21.1-1q21.2, followed by a short region of loss spanning

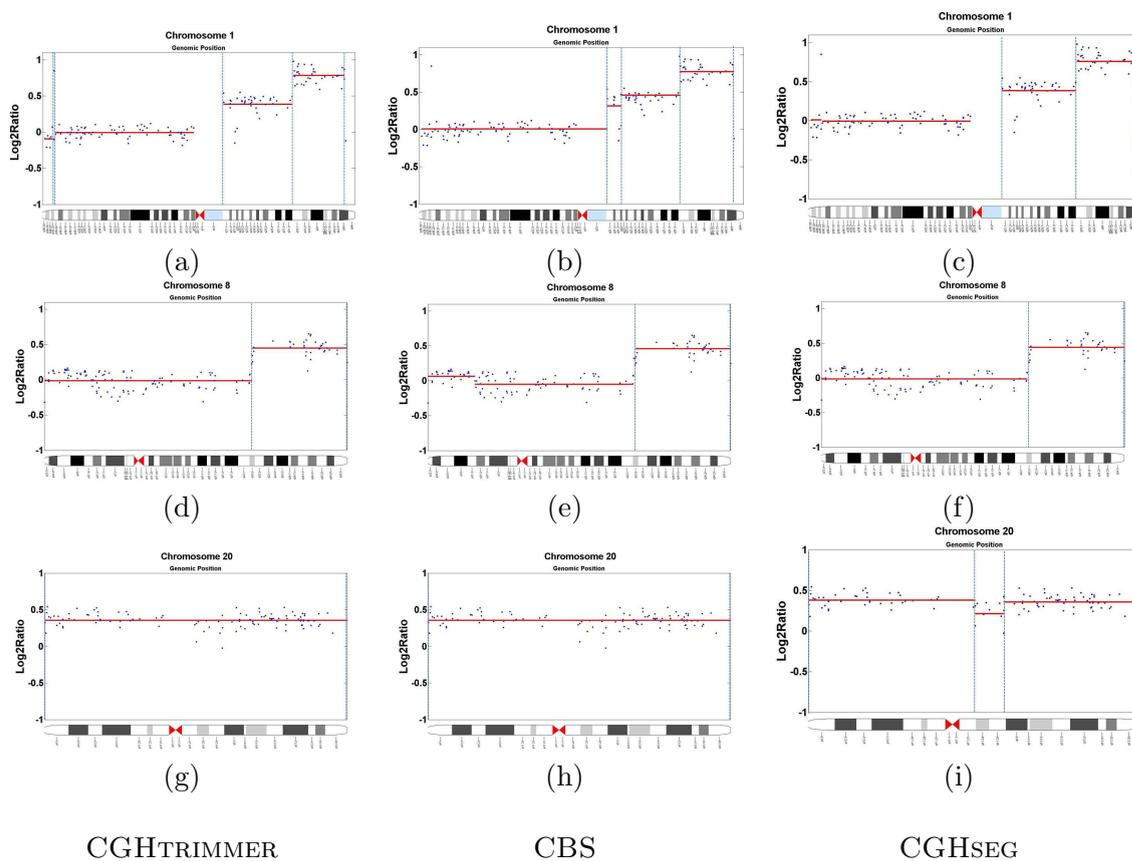


FIGURE 10.8: Visualization of the segmentation output of CGHTRIMMER, CBS, and CGHSEG for the cell line MCF10A on chromosomes 1 (a,b,c), 8 (d,e,f), and 20 (g,h,i). (a,d,g) CGHTRIMMER output. (b,e,h) CBS output. (c,f,i) CGHSEG output. Segments exceeding the ± 0.3 threshold [61] are highlighted.

1q21.3-1q22. CBS merges these into a single longer region of normal copy number. The existence of a small region of loss at this location in breast cancers is supported by prior literature [112].

The three methods provide comparable segmentations of chromosome 11 (Figure 10.7(d,e,f)). All predict loss near the p-terminus, a long segment of amplification stretching across much of the p- and q-arms, and additional amplification near the q-terminus. CGHTRIMMER, however, breaks this q-terminal amplification into several sub-segments at different levels of amplification while CBS and CGHSEG both fit a single segment to that region. We have no empirical basis to determine which segmentation is correct here. CGHTRIMMER does appear to provide a spurious break in the long amplified segment that is not predicted by the others.

Finally, along chromosome 20 (Figure 10.7(g,h,i)), the output of the methods is similar, with all three methods suggesting that the q-arm has an aberrant copy number, an observation consistent with prior studies [213]. The only exception is again that CBS

fits one point more than the other two methods along the first segment, causing a likely spurious extension of the p-arm's normal copy number into the q-arm.

Cell Line MCF10A Figure 10.8 shows the output of each of the three methods on chromosomes 1, 8, and 20 of the cell line MCF10A. On this cell line, the methods all yield similar predictions although from slightly different segmentations. All three show nearly identical behavior on chromosome 1 (Figure 10.8(a,b,c)), with normal copy number on the p-arm and at least two regions of independent amplification of the q-arm. Specifically, the regions noted as gain regions host significant genes such as PDE4DIP a gene associated with breast metastatic to bone (1q22), ECM1 (1q21.2), ARNT (1q21), MLLT11 (1q21), S100A10 (1q21.3), S100A13 (1q21.3), TPM3 (1q25) which also plays a role in breast cancer metastasis, SHC1 (1q21.3) and CKS1B (1q21.3). CBS provides a slightly different segmentation of the q-arm near the centromere, suggesting that the non-amplified region spans the centromere and that a region of lower amplification exists near the centromere. On chromosome 8 (Figure 10.8(d,e,f)) the three algorithms lead to identical copy number predictions after thresholding, although CBS inserts an additional breakpoint at 8q21.3 and a short additional segment at 8q22.2 that do not correspond to copy number changes. All three show significant amplification across chromosome 20 (Figure 10.8(g,h,i)), although in this case CGHSEG distinguishes an additional segment from 20q11.22-20q11.23 that is near the amplification threshold. It is worth mentioning that chromosome 20 hosts significant breast cancer related genes such as CYP24 and ZNF217.

Chapter 11

Robust Unmixing of Tumor States in Array Comparative Genomic Hybridization Data

11.1 Introduction

In Section 1.2.2 we discussed the importance of discovering tumor subtypes and we mentioned that the recent work by [358] showed promising results. They were, however, hampered by limitations of the hard geometric approach, particularly the sensitivity to experimental error and outlier data points caused by the simplex fitting approach. An example of simplex fitting in the plane is shown in Figure 11.1, illustrating why the strict containment model used in [102, 146, 358] is extremely sensitive to noise in the data. In the present Chapter we introduce a soft geometric unmixing model for tumor mixture separation, which relaxes the requirement for strict containment using a fitting criterion that is robust to noisy measurements. We develop a formalization of the problem and derive an efficient gradient-based optimization method. It is worth mentioning that this computational method has also been applied in the context of vertex similarity in social networks [391].

The remainder of this Chapter is organized as follows: Section 11.2 presents our proposed method and Section 11.3 we validate our method. In Section 11.4 we present a biological analysis of our findings on an aCGH data set taken from [312] and show that the method identifies state sets corresponding to known sub-types consistent with much of the analysis performed by the authors.

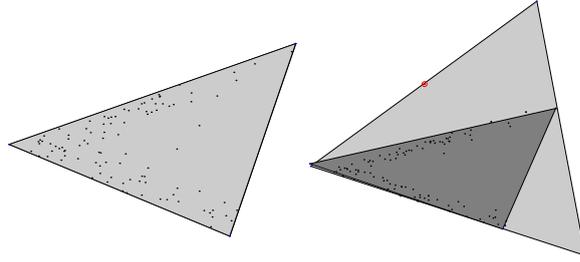


FIGURE 11.1: *Left:* The minimum area fit of a simplex containing the sample points in the plane (shown in black) using the program in § 11.2.1.1. On noiseless data, hard geometric unmixing recovers the locations of the fundamental components at the vertices. *Right:* However, the containment simplex is highly sensitive to noise and outliers in the data. A single outlier, circled above, radically changes the shape of the containment simplex fit (light gray above). In turn, this changes the estimates of basis distributions used to unmix the data. We mitigate this short coming by developing a soft geometric unmixing model (see § 11.2.1.2) that is comparatively robust to noise. The soft fit (shown dark gray) is geometrically very close to the generating sources as seen on the left.

11.2 Approach

The data are assumed to be given as g genes sampled in s tumors or tumor sections. The samples are collected in a matrix, $M \in \mathfrak{R}^{g \times s}$, in which each row corresponds to an estimate of gene copy number across the sample population obtained with aCGH. The data in M are processed as raw or baseline normalized raw input, rather than as log ratios. The “unmixing” model, described below, asserts that each sample m_i , a column of M , be well approximated by a convex combination of a fixed set of $C = [c_0 | \dots | c_k]$ of $k + 1$ unobserved basis distributions over the gene measurements. Further, the observed measurements are assumed to be perturbed by additive noise in the log domain, *i.e.*:

$$m_i = b^{\log_b(CF_i) + \eta}$$

where F_i is the vector of coefficients for the convex combination of the $(k + 1)$ basis distributions and η is additive zero mode *i.i.d.* noise.

11.2.1 Algorithms and Assumptions

Given the data model above, the inference procedure seeks to recover the $k + 1$ distributions over gene-copy number or expression that “unmix” the data. The procedure contains three primary stages:

1. Compute a reduced representation x_i for each sample m_i ,

2. Estimate the basis distributions K_{min} in the reduced coordinates and the mixture fractions F ,
3. Map the reduced coordinates K_{min} back into the “gene space” recovering C .

The second step in the method is performed by optimizing the objective in § 11.2.1.1 or the robust problem formulation in § 11.2.1.2.

Obtaining the reduced representation

We begin our calculations by projecting the data into a k dimension vector space (i.e., the intrinsic dimensionality of a $(k + 1)$ -vertex simplex). We accomplish this using principal components analysis (PCA) [335], which decomposes the input matrix M into a set of orthogonal basis vectors of maximum variance and retain only the k components of highest variance. PCA transforms the $g \times s$ measurement matrix M into a linear combination $XV + A$, where V is a matrix of the principal components of M , X provides a representation of each input sample as a linear combination of the components of V , and A is a $k \times s$ matrix in which each row contains s copies of the mean value of the corresponding row of M . The matrix X thus provides a reduced-dimension representation of M , and becomes the input to the sample mixture identification method in Stage 2. V and A are retained to allow us to later construct estimated aCGH vectors corresponding to the inferred mixture components in the original dimension g .

Assuming the generative model of the data above, PCA typically recovers a sensible reduced representation, as low magnitude log additive noise induces “shot-noise” behavior in the subspace containing the simplex with small perturbations in the orthogonal complement subspace. An illustration of this stage of our algorithm can be found in Figure 11.2.

Sample mixture identification

Stage 2 invokes either a hard geometric unmixing method that seeks the minimum volume simplex enclosing the input point set X (Program 11.1) or a soft geometric unmixing method that fits a simplex to the points balancing the desire for a compact simplex with that for containment of the input point set (Program 11.2). For this purpose, we place a prior over simplices, preferring those with small volume that fit or enclose the point set of X . This prior captures the intuition that the most plausible set of components explaining a given data set are those that can explain as much as possible of the observed data while leaving in the simplex as little empty volume, corresponding to mixtures that could be but are not observed, as possible.

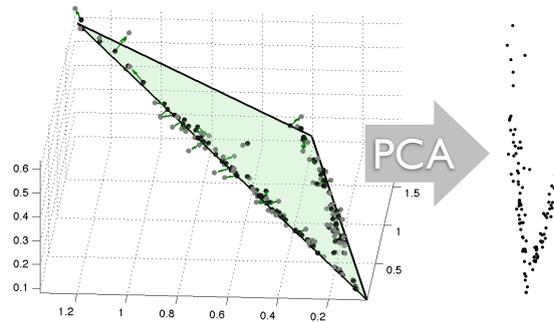


FIGURE 11.2: An illustration of the reduced coordinates under the unmixing hypothesis: points (shown in gray) sampled from the 3-simplex embedded in \mathbb{R}^3 and then perturbed by log-normal noise, producing points shown in black with sample correspondence given by the green arrows. Note that the dominant subspace remains in the planar variation induced by the simplex, and a 2D reduced representation for simplex fitting is thus sufficient.

Upon completion, Stage 2 obtains estimates of the vertex locations K_{min} , representing the inferred cell types from the aCGH data in reduced coordinates, and a set of mixture fractions describing the amount of each observed tumor sample attributed to each mixture component. The mixture fractions are encoded in a $(k+1) \times s$ matrix F , in which each column corresponds to the inferred mixture fractions of one observed tumor sample and each row corresponds to the amount of a single component attributed to all tumor samples. We define F_{ij} to be the fraction of component i assigned to tumor sample j and F_j to be vector of all mixture fractions assigned to a given tumor sample j . To ensure that the observations are modeled as convex combinations of the basis vertices, we require that $F\mathbf{1} = 1$.

Cell type identification

The reduced coordinate components from Stage 2, K_{min} , are projected up to a $g \times (k+1)$ matrix C in which each column corresponds to one of the $k+1$ inferred components and each row corresponds to the approximate copy number of a single gene in a component. We perform this transformation using the matrices V and A produced by PCA in Stage 1 with the formula $C = V^T K_{min} + A$, augmenting the average to $k+1$ columns.

Finally the complete inference procedure is summarized in the following pseudocode:

Given tumor sample matrix M , the desired number of mixture components k , and the strength of the volume prior γ :

1. Factor the sample matrix M such that $M^T = XV + A$
2. Produce the reduced k -dimensional representation by retaining the top k components in X

3. Minimize Program 11.1, obtaining an estimate of the simplex K_{min}^0
4. Minimize Program 11.2 starting at K_{min}^0 , obtaining K_{min} and F
5. Obtain the centers C in gene space as $C = A + V^T K_{min}$

11.2.1.1 Hard Geometric Unmixing

Hard geometric unmixing is equivalent to finding a minimum volume $(k + 1)$ -simplex containing a set of s points $\{X\}$ in \mathfrak{R}^k . A non-linear program for hard geometric unmixing can be written as follows:

$$\begin{aligned}
 \min_K & : \log \text{vol}(K) & (11.1) \\
 \forall i & : x_i = K F_i \\
 \forall F_i & : F_i^T \mathbf{1} = 1, F_i \succeq 0
 \end{aligned}$$

where $\log \text{vol}$ measures the volume of simplex defined by the vertices $K \doteq [v_0 | \dots | v_k]$ and $F \succeq 0$ requires that $\forall_{ij}, F_{ij} \geq 0$. Collectively, the constraints ensure that each point be expressed exactly as a unique convex combination of the vertices. Exact nonnegative matrix factorization (NNMF), see [272], can be seen as a relaxation of hard geometric unmixing. Exact NNMF retains the top two constraints while omitting the constraint that the columns of F sum to unity – thus admitting all positive combinations rather than the restriction to convex combinations as is the case for geometric unmixing.

Approximate and exponential-time exact minimizers are available for Program 11.1, in our experiments we use the approach of [102], which sacrifices some measure of accuracy for efficiency.

11.2.1.2 Soft Geometric Unmixing

Estimates of the target distributions, derived from the fundamental components (simplex vertices), produced by hard geometric unmixing are sensitive to the wide-spectrum noise and outliers characteristic of log-additive noise (*i.e.*, multiplicative noise in the linear domain). The robust formulation below tolerates noise in the sample measurements m_i and subsequently in the reduced representations x_i , improving the stability of these estimates. The sensitivity of hard geometric unmixing is illustrated in Figure 11.1. The motivation for soft geometric unmixing is to provide some tolerance to experimental error and outliers by relaxing the constraints in Program 11.1 allowing points to lie outside the boundary of the simplex fit to the data. We extend Program 11.1 to provide

a robust formulation as follows:

$$\begin{aligned} \min_K & : \sum_{i=1}^s |x_i - KF_i|_p + \gamma \log \text{vol}(K) \\ \forall F_i & : F_i^T \mathbf{1} = 1, F_i \succeq 0 \end{aligned} \quad (11.2)$$

where the term $|x_i - KF_i|_p$ penalizes the imprecise fit of the simplex to the data and γ establishes the strength of the minimum-volume prior. Optimization of Program 11.2 is seeded with an estimate produced from Program 11.1 and refined using MATLAB's *fminsearch* with analytical derivatives for the $\log \text{vol}$ term and an *LP*-step that determines mixtures components F_i and the distance to the boundary for each point outside the simplex.

We observe that when taken as whole, Program 11.2 can be interpreted as the negative log likelihood of a Bayesian model of signal formation. In the case of array CGH data, we choose $p = 1$ (*i.e.*, optimizing relative to an ℓ_1 norm), as we observe that the errors may be induced by outliers and the ℓ_1 norm would provide a relatively modest penalty for a few points far from the simplex. From the Bayesian perspective, this is equivalent to relaxing the noise model to assume *i.i.d.* heavy-tailed additive noise. To mitigate some of the more pernicious effects of log-normal noise, we also apply a total variation-like smoother to aCGH data in our experiments. Additionally, the method can be readily extended to weighted norms if an explicit outlier model is available.

11.2.1.3 Analysis & Efficiency

The hard geometric unmixing problem in § 11.2.1.1 is a non-convex objective in the present parameterization, and was shown by [324] to be NP-hard when $k + 1 \geq \log(s)$. For the special case of minimum volume tetrahedra ($k = 3$), [427] demonstrated an exact algorithm with time complexity $\Theta(s^4)$ and a $(1 + \epsilon)$ approximate method with complexity $O(s + 1/\epsilon^6)$. Below, we examine the present definition and show that Programs 11.1 and 11.2 have structural properties that may be exploited to construct efficient gradient based methods that seek local minima. Such gradient methods can be applied in lieu of or after heuristic or approximate combinatorial methods for minimizing Program 11.1, such as [102, 146] or the $(1 + \epsilon)$ method of [427] for simplexes in \mathbb{R}^3 .

We begin by studying the volume penalization term as it appears in both procedures. The volume of a convex body is well known (see [87]) to be a log concave function. In the case of a simplex, analytic partial derivatives with respect to vertex position can be used to speed the estimation of the minimum volume configuration K_{min} . The volume

of a simplex, represented by the vertex matrix $K = [v_0|\dots|v_k]$, can be calculated as:

$$\text{vol}(K) = c_k \cdot \det(\Gamma^T K K^T \Gamma)^{1/2} = c_k \cdot \det Q \quad (11.3)$$

where c_k is the volume of the unit simplex defined on $k+1$ points and Γ is a fixed vertex-edge incidence matrix such that $\Gamma^T K = [v_1 - v_0|\dots|v_k - v_0]$. The matrix Q is an inner product matrix over the vectors from the special vertex v_0 to each of the remaining k vertices. In the case where the simplex K is non-degenerate, these vectors form a linearly independent set and Q is positive definite (PD). While the determinant is log concave over PD matrices, our parameterization is linear over the matrices K , not Q . Thus it is possible to generate a degenerate simplex when interpolating between two non-degenerate simplexes K and K' . For example, let K define a triangle with two vertices on the y -axis and produce a new simplex K' by reflecting the triangle K across the y -axis. The curve $K(\alpha) = \alpha K + (1 - \alpha)K'$ linearly interpolates between the two. Clearly, when $\alpha = 1/2$, all three vertices of $K(\alpha)$ are co-linear and thus the matrix Q is not full rank and the determinant vanishes. However, in the case of small perturbations, we can expect the simplexes to remain non-degenerate.

To derive the partial derivative, we begin by substituting the determinant formulation into our volume penalization and arrive at the following calculation:

$$\begin{aligned} \log \text{vol}(K) &= \log c_k + \frac{1}{2} \log \det Q \\ &\propto \log \prod_{d=1}^k \lambda_d(Q) = \sum_{d=1}^k \lambda_d(Q) \end{aligned}$$

therefore the gradient of $\log \text{vol}(K)$ is given by

$$\frac{\partial \log \text{vol}(K)}{\partial K_{ij}} = \sum_{d=1}^k \frac{\partial}{\partial K_{ij}} \lambda_d = \sum_{d=1}^k z_d^T (\Gamma^T E_{ij} E_{ij}^T \Gamma) z_d$$

where the eigenvector z_d satisfies the equality $Q z_d = \lambda_d z_d$ and E_{ij} is the indicator matrix for the entry ij . To minimize the volume, we move the vertices along the paths specified by the negative log gradient of the current simplex volume. The Hessian is derived by an analogous computation, making Newton's method for Program 11.1, with log barriers over the equality and inequality constraints, a possible optimization strategy.

Soft geometric unmixing (Program 11.2) trades the equality constraints in Program 11.1 for a convex, but non-differentiable term, in the objective function $\sum_{i=1}^s |x_i - KF_i|_p$ for $p|1 \leq p \leq 2$. Intuitively – points inside the simplex have no impact on the cost of the fit. However, over the course of the optimization, as the shape of the simplex changes points move from the interior to the exterior, at which time they incur a cost. To determine

this cost, we solve the nonnegative least squares problem for each mixture fraction F_i , $\min_F : (KF_i - x_i)^T(KF_i - x_i)$. This step simultaneously solves for the mixture fraction, and for exterior points, the distance to the simplex is determined. The simplex is then shifted under a standard shrinkage method based on these distances.

11.3 Experimental Methods

We evaluated our methods using synthetic experiments, allowing us to assess two properties of robust unmixing 1) the fidelity with which endmembers (sub-types) are identified and 2) the relative effect of noise on hard versus robust unmixing. We then evaluate the robust method on a real world aCGH data set published by [312] in which ground truth is not available, but for which we uncover much the structure reported by the authors.

11.3.1 Methods: Synthetic Experiments

To test the algorithms given in §11.2 we simulated data using a biologically plausible model of ad-mixtures. Simulated data provides a quantitative means of evaluation as ground truth is available for both the components C and the mixture fractions F_i associated with each measurement in the synthetic design matrix M . The tests evaluate and compare hard geometric unmixing §11.2.1.1 and soft geometric unmixing §11.2.1.2 in the presence of varying levels of log-additive Gaussian noise and varying k . By applying additive Gaussian noise in the log domain we simulate the heteroscedasticity characteristic of CGH measurements (*i.e.* higher variance with larger magnitude measurements). By varying k , the dimensionality of the simplex used to fit the data, we assess the algorithmic sensitivity to this parameter as well as that to γ governing the strength of the volume prior in Program 2. The sample generation process consists of three major steps: 1) mixture fraction generation (determining the ratio of sub-types present in a sample), 2) end-member (*i.e.* sub-type) generation and 3) the sample perturbation by additive noise in the log-ratio domain.

11.3.1.1 Mixture Sampler

Samples over mixture fractions were generated in a manner analogous to Polya's Urn Process, in which previously sampled simplicial components (*e.g.*, line segments, triangles, tetrahedra) are more likely to be sampled again. This sampling mechanism produces data distributions that are similar to those we see in low dimensional projections of aCGH data when compared against purely uniform samples over mixtures. An

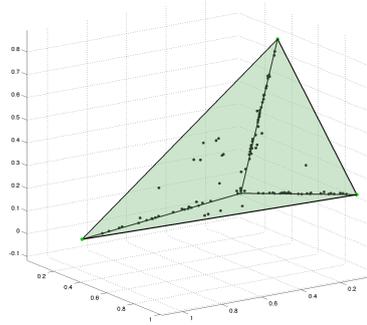


FIGURE 11.3: An example sample set generated for §11.3.1.2 shown in the “intrinsic dimensions” of the model. Note that sample points cleave to the lower dimensional substructure (edges) of the simplex.

example of a low dimensional sample set and the simplex that was used to generate the points is shown in Figure 11.3.

To generate the mixture fractions F_i for the i^{th} sample, the individual components in C^{true} are sampled without replacement from a dynamic tree model. Each node in the tree contains a dynamic distribution over the remaining components, each of which is initialized to the uniform distribution. We then sample s mixtures by choosing an initial component according to the root’s component distribution and proceed down the tree. As a tree-node is reached, its component distribution is updated to reflect the frequency with which its children are drawn. To generate the i^{th} sample, the fractional values F_i are initialized to zero. As sample generation proceeds, the currently selected component C_j updates the mixture as $F_{ij} \sim \text{uniform}[(1/2)f_p^j, 1]$ where f_p^j is the frequency of j ’s parent node. For the i^{th} mixture, this process terminates when the condition $1 \leq \sum_{j=1}^{k+1} F_{ij}$ holds. Therefore, samples generated by long paths in the tree will tend to be homogenous combinations of the components C^{true} , where as short paths will produce lower dimensional substructures. At the end of the process, the matrix of fractions F is re-normalized so that the mixtures associated with each sample sum to unity. This defines a mixture F_i^{true} for each sample – *i.e.* the convex combination over fundamental components generating the sample point.

11.3.1.2 Geometric Sampling of End-members & Noise

To determine the location of the end-members we specify an extrinsic dimension (number of genes) g , and an intrinsic dimension k (requiring $k+1$ components). We then simulate $k+1$ components by constructing a $g \times (k+1)$ matrix C^{true} of fundamental components in which each column is an end-member (*i.e.* sub-type) and each row is the copy number of one hypothetical gene, sampled from the unit Gaussian distribution and rounded to

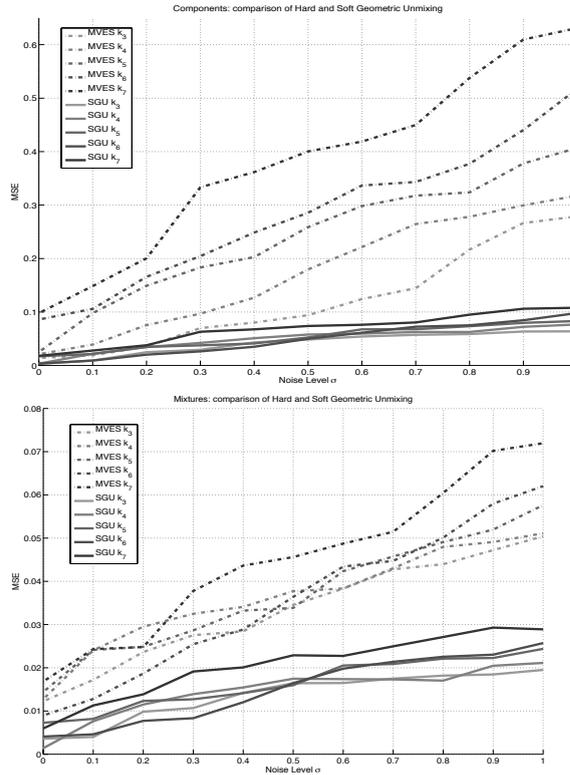


FIGURE 11.4: Left: mean squared error for the component reconstruction comparing Hard Geometric Unmixing (MVES: [102]) and Soft Geometric Unmixing (SGU) introduced in §11.2.1.2 for the experiment described in §11.3.1.2 with variable γ . The plot demonstrates that robust unmixing more accurately reconstructs the ground truth centers relative to hard unmixing in the presence of noise. Right: mean squared error for mixture reconstruction comparing MVES and SGU.

the nearest integer. Samples m_i , corresponding to the columns of the data matrix M , are then given by:

$$m_i = 2^{\log_2(C^{true} F_i^{true}) + \frac{1}{2}\sigma\eta} \quad (11.4)$$

where $\eta \sim \text{normal}(0, 1)$ and the mixture fractions F_i^{true} were obtained as in §11.3.1.1.

11.3.1.3 Evaluation

We follow Schwartz and Shackney [358] in assessing the quality of the unmixing methods by independently measuring the accuracy of inferring the components and the mixture fractions. We first match inferred mixture components to true mixture components by performing a maximum weighted bipartite matching of columns between C^{true} and the inferred components C^e , weighted by negative Euclidean distance. We will now assume that the estimates have been permuted according to this matching and continue. We then assess the quality of the mixture component identification by the root mean square

distance over all entries of all components between the matched columns of the two C matrices:

$$\mathbf{error} = \frac{1}{g(k+1)} \|C^{true} - C^e\|_F^2 \quad (11.5)$$

where $\|A\|_F = \sqrt{\sum_{ij} a_{ij}^2}$ denotes the Frobenius norm of the matrix A .

We similarly assess the quality of the mixture fractions by the root mean square distance between F^{true} and the inferred fractions F^e over all genes and samples:

$$\mathbf{error} = \frac{1}{g(k+1)} \|F^{true} - F^e\|_F^2. \quad (11.6)$$

This process was performed for $s = 100$ and $d = 10000$ to approximate a realistic tumor expression data set and evaluated for $k = 3$ to $k = 7$ and for $\sigma = \{0, 0.1, 0.2, \dots, 1.0\}$, with ten repetitions per parameter.

11.4 Results

11.4.1 Results: Synthetic Data

The results for the synthetic experiment are summarized in Figure 11.4. The figure shows the trends in MSE for hard geometric unmixing §11.2.1.1 and soft geometric unmixing §11.2.1.2 on the synthetic data described above. As hard geometric unmixing requires that each sample lie inside the fit simplex, as noise levels increase (larger σ), the fit becomes increasingly inaccurate. Further, the method MVES deteriorates to some degree as order k of the simplex increases. However, soft geometric unmixing degrades more gracefully in the presence of noise if an estimate of the noise level is available with ± 0.1 in our current model. The trend of soft unmixing exhibiting lower error and better scaling in k than hard unmixing holds for both components and mixture fractions, although components exhibit a higher average degree of variability due to the scale of the synthetic measurements when compared to the mixture fractions.

11.4.2 Array Comparative Gene Hybridization (aCGH) Data

We further illustrate the performance of our methods on a publicly available primary Ductal Breast Cancer aCGH Dataset furnished with [312]. This dataset is of interest in that each tumor sample has been sectored multiple times during biopsy which is ideal for understanding the substructure of the tumor population. The data consists of 87 aCGH profiles from 14 tumors run on a high-density ROMA platform with 83055

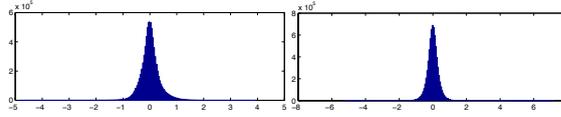


FIGURE 11.5: Empirical motivation for the $\ell_1 - \ell_1$ -total variation functional for smoothing CGH data. The left plot shows the histogram of values found in the CGH data obtained from the [312] data set. The distribution is well fit by the high kurtosis Laplacian distribution in lieu of a Gaussian. The right plot shows the distribution of differences along the probe array values. As with the values distribution, these frequencies exhibit high kurtosis.

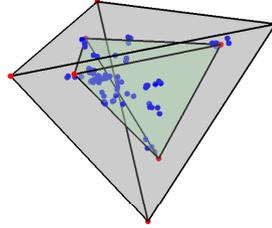


FIGURE 11.6: The simplex fit to the CGH data samples from [312] ductal data set in \mathbb{R}^3 . The gray tetrahedron was returned by the optimization of Program 11.1 and the green tetrahedron was returned by the robust unmixing routine.

probes. Profiles are derived from 4-6 sectors per tumor, with samples for tumors 5-14 sub-partitioned by cell sorting according to total DNA content, and with healthy control samples for tumors 6, 9, 12, and 13. For full details, the reader is referred to Navin et al. [312]. The processed data consists of \log_{10} ratios and which were exponentiated prior to the PCA step (Stage 1) of the method.

11.4.2.1 Preprocessing

To mitigate the effects of sensor noise on the geometric inference problem we apply a total variation (TV) functional to the raw log-domain data. The $\ell_1 - \ell_1$ -TV minimization is equivalent to a penalized projection onto the over-complete Harr basis preserving a larger degree of the signal variation when compared to discretization methods (*e.g* [205, 320]) that employ aggressive priors over the data distribution. The procedure seeks a smooth instance x of the observed signal s by optimizing the following functional:

$$\min_x : \sum_{i=1}^g |x_i - s_i|_1 + \lambda \sum_{i=1}^{g-1} |x_i - x_{i+1}|_1 \quad (11.7)$$

The functional 11.7 is convex and can be solved readily using Newton's method with log-barrier functions ([87]). The solution x can be taken as the maximum likelihood estimate of a Bayesian model of CGH data formation. That is, the above is the negative log-likelihood of a simple Bayesian model of signal formation. The measurements \hat{x}_i are

assumed to be perturbed by the *i.i.d.* Laplacian noise and the changes along the probe array are assumed to be sparse. Recall that the Laplacian distribution is defined by the probability density function $\Pr[x] = \frac{1}{z} \exp\left(\frac{-|x|}{a}\right)$. In all experiments the strength of the prior λ was set to $\lambda = 10$. The data fit this model well as illustrated in Figure 11.5. The dimension of the reduced representation k , fixing the number of fundamental components, was determined using the eigengap heuristic during the PCA computation (Stage 1). This rule ceases computing additional principal components when the difference in variances jumps above threshold.

11.4.2.2 Unmixing Analysis and Validation

The raw data was preprocessed as described above and a simplex was fit to the reduced coordinate representation using the soft geometric unmixing method (see §11.2.1.2). A three dimensional visualization of the resulting fit is shown for the [312] data set in Figure 11.6. To assess the performance with increasing dimensionality, we ran experiments for polytope dimensionality k ranging from 3 to 9. Following the eigen-gap heuristic we chose to analyze the results for $k = 6$. The γ value was picked according to the estimated noise level in the aCGH dataset and scaled relative to the unit simplex volume (here, $\gamma = 100$). The estimated 6 components/simplex vertices/pure cancer types are labeled C_1, C_2, \dots, C_6 .

Figure 11.7 shows mixture fraction assignments for the aCGH data for $k = 6$. While there is typically a non-zero amount of each component in each sample due to imprecision in assignments, the results nonetheless show distinct subsets of tumors favoring different mixture compositions and with tumor cells clearly differentiated from healthy control samples. The relative consistency within versus between tumors provides a secondary validation that soft unmixing is effective at robustly assigning mixture fractions to tumor samples despite noise inherent to the assay and that produced by subsampling cell populations. It is also consistent with observations of Navin *et al.*

It is not possible to know with certainty the true cell components or mixture fractions of the real data, but we can validate the biological plausibility of our results by examining known sites of amplification in the inferred components. We selected fourteen benchmark loci frequently amplified in breast cancers through manual literature search. Table 11.1 lists the chosen benchmarks and the components exhibiting at least 2-fold amplification of each. Figure 11.8 visualizes the results, plotting relative amplification of each component as a function of genomic coordinate and highlighting the locations of the benchmark markers. Thirteen of the fourteen benchmark loci exhibit amplification for a subset of the components, although often at minimal levels. The components also

show amplification of many other sites not in our benchmark set, but we cannot definitively determine which are true sites of amplification and which are false positives. We further tested for amplification of seven loci reported as amplified by Navin *et al.* [312] specifically in the tumors examined here and found that six of the seven are specifically amplified in one of our inferred components: PPP1R12A (C_2), KRAS (C_2), CDC6 (C_2), RARA (C_2), EFNA5 (C_2), PTPN1 (C_3), and LPXN (not detected). Our method did not infer a component corresponding to normal diploid cells as one might expect due to stromal contamination. This failure may reflect a bias introduced by the dataset, in which many samples were cell sorted to specifically select aneuploid cell fractions, or could reflect an inherent bias of the method towards more distinct components, which would tend to favor components with large amplifications.

We repeated these analyses for the hard unmixing with a higher amplification threshold due to the noise levels in the centers. It detected amplification at 11 of the 14 loci, with spurious inferences of deletion at four of the 11. For the seven sites reported in Navin *et al.*, hard unmixing identified five (failing to identify EFNA5 or LPXN) and again made spurious inferences of deletions for three of these sites, an artifact the soft unmixing eliminates. The results suggest that hard unmixing produces less precise fits of simplexes to the true data.

We can also provide a secondary analysis based on Navin *et al.*'s central result that the tumors can be partitioned into monogenomic (those appearing to show essentially a single genotype) and polygenomic (those that appear to contain multiple tumor sub-populations). We test for monogeniety in mixture fractions by finding the minimum correlation coefficient between mixture fractions of consecutive tumor sectors (ignoring normal controls) maximized over all permutations of the sectors. Those tumors with correlations above the mean over all tumors (0.69) were considered monogenomic and the remainder polygenomic. Navin *et al.* assign $\{1, 2, 6, 7, 9, 11\}$ as monogenomic and $\{3, 4, 5, 8, 10, 12, 13, 14\}$ as polygenomic. Our tests classify $\{1, 2, 5, 6, 7, 8, 11\}$ as monogenomic and $\{3, 4, 10, 12, 13, 14\}$ as polygenomic, disagreeing only in tumors 5 and 8. Our methods are thus effective at identifying true intratumor heterogeneity in almost all cases without introducing spurious heterogeneity. By contrast, hard unmixing identifies only tumors 5 and 8 as polygenomic, generally obscuring true heterogeneity in the tumors.

Our long-term goal in this work is not just to identify sub-types, but to describe the evolutionary relationships among them. We have no empirical basis for validating any such predictions at the moment but nonetheless consider the problem informally here for illustrative purposes. To explore the question of possible ancestral relationships among components, we manually examined the most pronounced regions of shared gain across

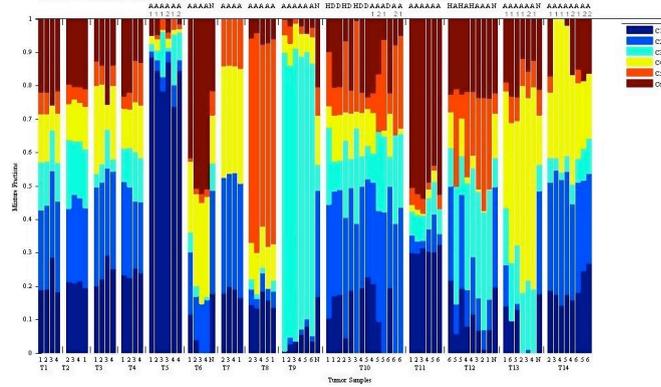


FIGURE 11.7: Inferred mixture fractions for six-component soft geometric unmixing applied to breast cancer aCGH data. Data is grouped by tumor, with multiple sectors per tumor placed side-by-side. Columns are annotated below by sector or N for normal control and above by cell sorting fraction (D for diploid, H for hypodiploid, A for aneuploid, and A1/A2 for subsets of aneuploid) where cell sorting was used.

Marker	Locus	Component
MUC1	1q21	C1,C4
PIK3CA	3q26.3	C3,C6
ESR1	6q25.1	C4
EGFR	7p12	C5
c-MYC	8q24	C1,C3,C5
PTEN	10p23	none
PGR	14q23.2	C6
CCND1	11q13	C4
BRCA2	13q12.3	C5
ESR2	14q23	C1
BRCA1, ERBB2	17q21	C5,C6
STAT5A, STAT5B	17q11.2	C5
GRB7	17q12	C6
CEA	19q13.2	C6

TABLE 11.1: Benchmark set of breast cancer markers selected for validation of real data, annotated by gene name, genomic locus, and the set of components exhibiting amplification at the given marker.

component. Figure 11.9 shows a condensed view of the six components highlighting several regions of shared amplification between components. The left half of the image shows components 3, 5, and 1, revealing a region of shared gain across all three components at 9p21 (labeled B). Components 5 and 1 share an additional amplification at 1q21 (labeled A). Components 1 and 5 have distinct but nearby amplifications on chromosome 17, with component 1 exhibiting amplification at 17q12 (labeled D) and component 5 at 17q21 (labeled C). We can interpret these images to suggest a possible evolutionary scenario: component 3 initially acquires an amplification at 9p21 (the locus

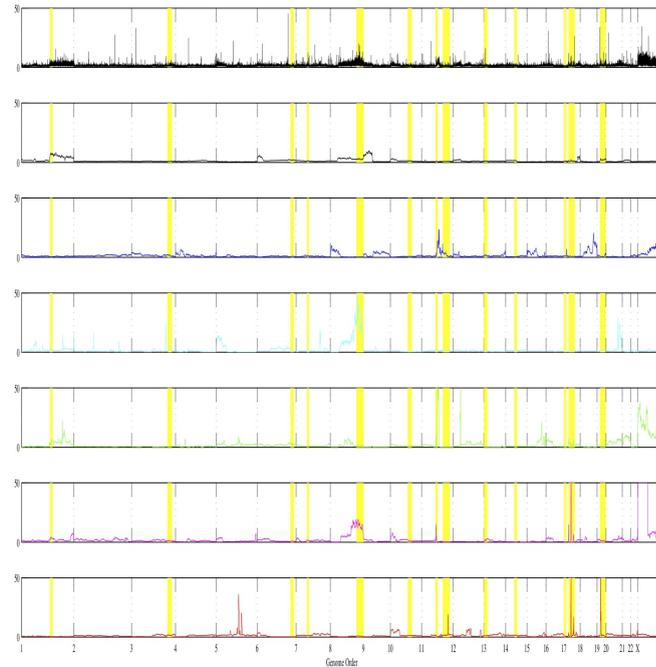


FIGURE 11.8: Copy numbers of inferred components versus genomic position. The average of all input arrays (top) is shown for comparison, with the six components below. Benchmarks loci are indicated by yellow vertical bars.

of the gene *CDKN2B/p15INK4b*), an unobserved descendent of component 3 acquires secondary amplification at 1q21 (the locus of *MUC1*), and this descendent then diverges into components 1 and 5 through acquisition of independent abnormalities at 17q12 (site of *PGAP3*) or 17q21 (site of *HER2*). The right side of the figure similarly shows some sharing of sites of amplification between components 2, 4, and 6, although the amplified regions do not lead to so simple an evolutionary interpretation. The figure is consistent with the notion that component 2 is ancestral to 4, with component 2 acquiring a mutation at 5q21 (site of *APC/MCC*) and component 4 inheriting that mutation but adding an additional one at 17q21. We would then infer that the amplification at the *HER2* locus arose independently in component 6, as well as in component 5. The figure thus suggests the possibility that the *HER2*-amplifying breast cancer sub-type may arise from multiple distinct ancestral backgrounds in different tumors. While we cannot evaluate the accuracy of these evolutionary scenarios, they provide an illustration of how the output of this method is intended to be used to make inferences of evolutionary pathways of tumor states.

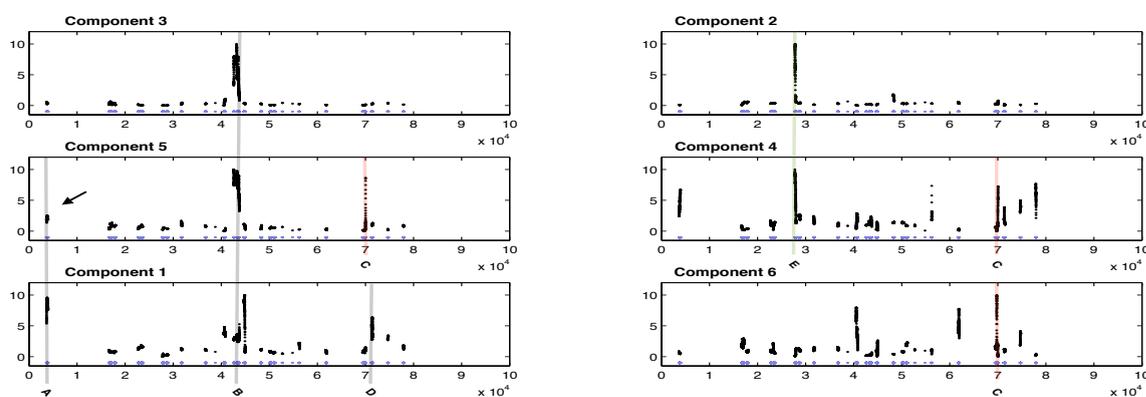


FIGURE 11.9: Plot of amplification per probe highlighting regions of shared amplification across components. The lower (blue) dots mark the location of the collected cancer benchmarks set. Bars highlight specific markers of high shared amplification for discussion in the text. *Above:* **A:** 1q21 (site of MUC1), **B:** 9p21 (site of CDKN2B), **C:** 7q21 (site of HER2), **D:** 17q12 (site of PGAP3), **E:** 5q21 (site of APC/MCC).

Chapter 12

Perfect Reconstruction of Oncogenetic Trees

12.1 Introduction

In this Chapter we answer a fundamental question regarding oncogenetic trees, see Section 1.2.3. Before we state the question, we introduce some notation first. Let $T = (V, E, r)$ be a rooted tree on V , i.e., every vertex has in-degree at most one and there are no cycles, and let $r \in V$ be the root of T . Given a finite family $\mathcal{F} = \{A_1, \dots, A_q\}$ of sets of vertices, i.e., $A_i \subseteq V(T)$ for $i = 1, \dots, q$, where each A_i is the vertex set of a r -rooted sub-tree of T , what are the necessary and sufficient conditions, if any, which allow us to uniquely reconstruct T ? In this work we treat this natural combinatorial question, namely:

“When can we reconstruct an oncogenetic tree T from a set family \mathcal{F} ?”

Despite the fact that in practice aCGH data tend to be noisy and consistent with more than one oncogenetic trees, the question is nonetheless interesting and to the best of our knowledge remains open so far [330, 350]. Theorem 12.1 provides the necessary and sufficient conditions to uniquely reconstruct an oncogenetic tree. We write $u \prec v$ ($u \not\prec v$) to denote that u is (not) a descendant of v in T .

Theorem 12.1. *Let T be an oncogenetic tree and $\mathcal{F} = \{A_1, \dots, A_q\}$ be a finite family of sets of vertices, i.e., $A_i \subseteq V(T)$ for $i = 1, \dots, q$, where each A_i is the vertex set of a r -rooted sub-tree of T . The necessary and sufficient conditions to uniquely reconstruct the tree T from the family \mathcal{F} are the following:*

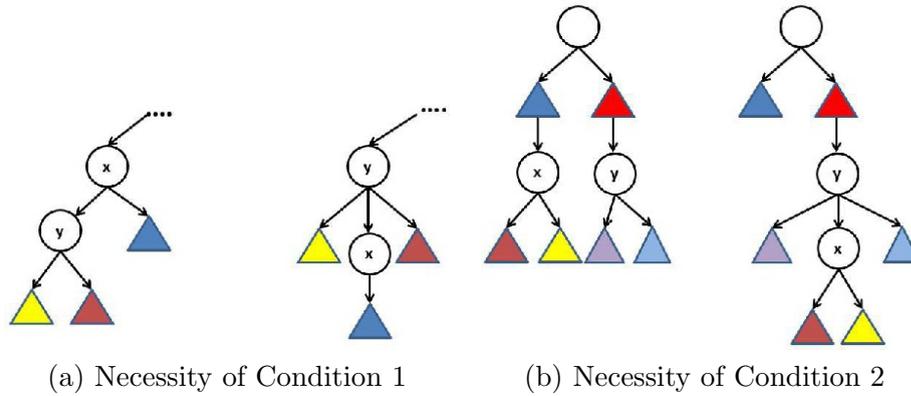


FIGURE 12.1: Illustration of necessity conditions of Theorem 12.1.

1. For any two distinct vertices $x, y \in V(T)$ such that $(x, y) \in E(T)$, there exists a set $A_i \in \mathcal{F}$ such that $x \in A_i$ and $y \notin A_i$.
2. For any two distinct vertices $x, y \in V(T)$ such that $y \not\prec x$ and $x \not\prec y$ there exist sets $A_i, A_j \in \mathcal{F}$ such that $x \in A_i, y \notin A_i$ and $x \notin A_j$ and $y \in A_j$.

In Section 12.2 we prove Theorem 12.1. It is worth noticing that our proof provides a simple procedure for the reconstruction as well.

12.2 Proof of Theorem 12.1

In the following we call a tree T consistent with the family set \mathcal{F} if all sets $A_i \in \mathcal{F}$ are vertices of rooted sub-trees of T . Notice that when two or more trees are consistent with the input dataset \mathcal{F} , then we cannot uniquely reconstruct T .

Proof. First we prove the necessity of conditions 1,2 and then their sufficiency to reconstruct T .

Necessity: For the sake of contradiction, assume that Condition 1 does not hold. Therefore, there exists two vertices $x, y \in V(T)$ such that there exists no set $A \in \mathcal{F}$ that contains one of them. Then, the two trees shown in Figure 12.1(a) are both consistent with \mathcal{F} . Therefore we cannot reconstruct T . Similarly, assume that Condition 2 does not hold. Specifically assume that for all j such that $x \in A_j$, then $y \in A_j$ too (for the symmetric case the same argument holds). Then, both trees in Figure 12.1(b) are consistent with \mathcal{F} and therefore T is not reconstructable. The symmetric case follows by the same argument.

Sufficiency: Let $x \in V(T)$ and P_x be the vertex set of the unique path from the root r to x , i.e., $P_x = \{r, \dots, x\}$. Also, define F_x to be the intersection of all sets in the family

\mathcal{F} that contain vertex x , i.e., $F_x = \bigcap_{A_i \ni x} A_i$. We prove that $F_x = P_x$. Since by the definition of an oncogenetic tree $P_x \subseteq F_x$ it suffices to show that $F_x \subseteq P_x$. Assume for the sake of contradiction that $F_x \not\subseteq P_x$. Then, there exists a vertex $v \in V(T)$ such that $v \in F_x, v \notin P_x$. We consider the following three cases.

- CASE 1 ($x \prec v$): Since by definition each set $A_i \in \mathcal{F}$ is the vertex set of a rooted sub-tree of T , $v \in P_x$ by the definition of an oncogenetic tree.
- CASE 2 ($v \prec x$): Inductively by condition 1, there exists $A_i \in \mathcal{F}$ such that $x \in A_i, v \notin A_i$. Therefore, $v \notin F_x$.
- CASE 3 ($x \not\prec v, v \not\prec x$): By condition 2, there exists $A_i \in \mathcal{F}$ such that $x \in A_i$ and $v \notin A_i$. Hence, $v \notin F_x$.

In all three cases above, we obtain a contradiction and therefore $v \in F_x \Rightarrow v \in P_x$. Therefore, $F_x \subseteq P_x$ and subsequently $F_x = P_x$. Given this fact, it is easy to reconstruct the tree T . We sketch the algorithm: compute for each x the set F_x which is the unordered set of vertices of the unique path from r to x . The ordering of the vertices which results in finding the path P_x , i.e., $(v_0 = r \rightarrow v_1 \rightarrow \dots \rightarrow v_{k-1} \rightarrow v_k = x)$ is computed using sets in \mathcal{F} which contain v_i but not v_{i+1} , $i = 0, \dots, k-1$. The existence of such sets is guaranteed by condition 1. □

III Conclusion and Future Directions

Chapter 13

Conclusion

This Chapter concludes the dissertation by posing both specific open problems and broader directions related to our work.

13.1 Open Problems

13.1.1 Rainbow Connectivity (Chapter 3)

We propose two open problems for future work.

Erdős-Rényi graphs: In Chapter 3 we give an asymptotically tight result on the rainbow connectivity of $G = G(n, p)$ at the connectivity threshold. It is reasonable to conjecture that this could be tightened. We conjecture that *whp*, $rc(G) = \max\{Z_1, \text{diameter}(G(n, p))\}$.

Random regular graphs: Our result on random regular graphs is not so tight. It is still reasonable to believe that the above conjecture also holds in this case. (Of course $Z_1 = 0$ here).

13.1.2 Random Apollonian Graphs (Chapter 4)

We propose two open problems for future work.

Diameter: We conjecture that $\lim_{t \rightarrow +\infty} \frac{d(G_t)}{\log t} \rightarrow c$, where c is a constant. Finding c is an interesting research problem.

Conductance: We conjecture that the conductance of a RAN is $\Theta\left(\frac{1}{\sqrt{t}}\right)$ *whp*. Figure 13.1 shows that $\phi(G_t) \leq \frac{1}{\sqrt{t}}$.

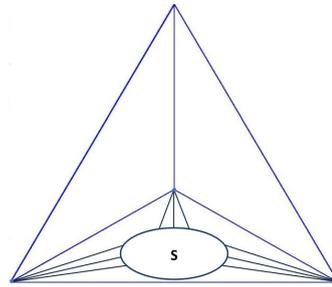


FIGURE 13.1: By the pigeonhole principle, one of the three initial faces receives $\Theta(t)$ vertices. Using Theorem 4.3 it is not hard to see that the encircled set of vertices S has conductance $\phi(S) \approx \frac{\sqrt{t}}{t} = \frac{1}{\sqrt{t}}$ whp.

13.1.3 Triangle Counting (Chapter 5)

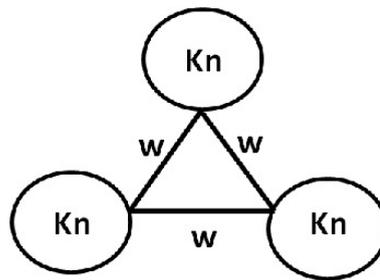


FIGURE 13.2: Weighted Graph

Certain open problems on triangle counting follow:

Weighted graphs: The sampling scheme presented in Section 5.2.1 can be adapted to weighted graphs: multiply the weight of sampled edge by $\frac{1}{p}$ and count weighted triangles. However this can be problematic as the graph shown in Figure 13.2 indicates. Specifically, for a sufficiently large w , throwing out any weighted edge results in an arbitrarily bad estimate of the count of triangles. Finding a better sampling scheme for weighted graphs is left as a problem for future work.

Better sparsification: Finding an easy-to-compute quantity which gives us an optimal way to sparsify the graph with respect to triangles is an interesting research problem.

Random projections and triangles: Can we provide some reasonable condition on G that would guarantee (5.7)?

MapReduce: Our proposed methods are easily parallelizable and developing a MAPREDUCE implementation is a natural practical direction.

Streaming model : Besides MAPREDUCE, the streaming model is particularly well tailored for large data. Can we design better algorithms for triangle counting for the

streaming model [96]? Can we estimate the average number of triangles over all vertices with a given degree d ?

13.1.4 Densest Subgraphs (Chapter 6)

Our work in Chapter 6 leaves several open problems.

Algorithm Design: Can we design approximation algorithms with better additive guarantees or a multiplicative approximation algorithm without shifting the objective? Can we characterize the worst-case behavior of LOCALSEARCHOQC? Also designing efficient randomized algorithms, e.g., [32] is an interesting research direction.

Bipartite graphs: A natural extension of our objective to a bipartite graph $G(L \cup R, E)$ is the following. For a set $A \subseteq L, B \subseteq R$ and a real parameter α , we define $f_\alpha(A, B) = e[A \cup B] - \alpha|A||B|$. Maximizing f_α is in general a hard problem. This can be easily shown using a reduction similar to the one we used in Chapter 6 where we used a planted clique in an Erdős-Rényi $G(n, 1/2)$ graph and invoking the results of [165] on detecting planted bipartite clique distributions. Understanding this objective and evaluating its performance is an interesting direction.

13.1.5 Structure of the Web Graph (Chapter 7)

Studying the Web graph is an important problem in graph mining research, given its prominence. Two problems related to our work follow.

Improving Hadi: How well do recent advances in moment estimation [234] perform in practice? Can we use them to improve the performance of HADI?

Analyzing the Web graph: Using our algorithmic tool HADI we provided insights into certain fundamental properties of the Web graph. Can we provide further insights by studying other statistics, e.g., [398]? Can we use these structural properties to come up with good models of the Web graph?

13.1.6 FENNEL: Streaming Graph Partitioning for Massive Scale Graphs (Chapter 8)

How to interpolate? In Chapter 8 we showed that by selecting an exponent γ between 1 and 2 allows us to interpolate between the two heuristics most frequently used for balanced graph partitioning. Can we find an easy-to-compute graph statistic, e.g., power law exponent, that allows us to make a good selection for γ for a given graph?

Random graphs with hidden partition: Assume that we have a dense graph with k clusters where the probability of having an edge insider and between two distinct clusters is p and q respectively, $p > q = \Theta(1)$. Can we predict how different heuristics derived from our framework will perform?

More cost functions: It is interesting to test other choices of cost functions and tackle the asymmetric edge cost issue that occurs in standard data center architectures.

Communication cost: A future goal is to minimize further the communication cost using minwise hashing [349].

13.1.7 PEGASUS: A System for Large-Scale Graph Processing (Chapter 9)

Given that PEGASUS is much more an engineering rather than a theoretical contribution, we provide a broad research direction, related to MAPREDUCE, the framework on which PEGASUS relies.

Understanding the limits of MapReduce: Afrati, Das Sarma, Salihoglu and Ullman [12] provide an interesting framework for studying MAPREDUCE algorithms. Let the replication rate of any MAPREDUCE algorithm be the average number of reducers each input is sent to. Understanding trade-offs between the replication rate and the efficiency of MAPREDUCE algorithms is an interesting research direction.

13.1.8 Approximation Algorithms for Speeding up Dynamic Programming and Denoising aCGH data (Chapter 10)

In Chapter 10, we present a new formulation for the problem of denoising aCGH data. Our formulation has already proved to be valuable in numerous settings [136].

Improving algorithmic performance: Our results strongly indicate that the $O(n^2)$ algorithm, which is — to the best of our knowledge — the fastest exact algorithm, is not tight. There is inherent structure in the optimization problem. Lemma 13.1 is such an example.

Lemma 13.1. *If $|P_{i_1} - P_{i_2}| > 2\sqrt{2C}$, then in the optimal solution of the dynamic programming using L_2 norm, i_1 and i_2 are in different segments.*

Proof. The proof is by contradiction, see also Figure 13.3. Suppose the optimal solution has a segment $[i, j]$ where $i \leq i_1 < i_2 \leq j$, and its optimal x value is x^* . Then consider splitting it into 5 intervals $[i, i_1 - 1]$, $[i_1, i_1]$, $[i_1 + 1, i_2 - 1]$, $[i_2, i_2]$, $[i_2 + 1, r]$. We let

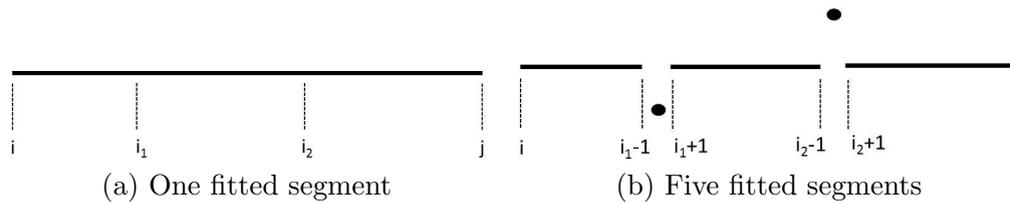


FIGURE 13.3: Lemma 13.1: if $|P_{i_1} - P_{i_2}| > 2\sqrt{2C}$ then Segmentation (b) which consists of five segments (two of which are single points) is superior to Segmentation (a) which consists of a single segment.

$x = x^*$ be the fitted value to the three intervals not containing i_1 and i_2 . Also, as $|P_{i_1} - P_{i_2}| > 2\sqrt{2C}$, $(P_{i_1} - x)^2 + (P_{i_2} - x)^2 > 2\sqrt{2C}^2 = 4C$. So by letting $x = P_{i_1}$ in $[i_1, i_1]$ and $x = P_{i_2}$ in $[i_2, i_2]$, the total decreases by more than $4C$. This is more than the added penalty of having 4 more segments, a contradiction with the optimality of the segmentation. \square

Uncovering and taking advantage of the inherent structure in a principled way should result in a faster exact algorithm. This is an interesting research direction which we leave as future work.

More applications: Another research direction is to find more applications (e.g., histogram construction [204]) to which our methods are applicable.

aCGH Denoising: Our model does not perform well in very noisy data. Developing novel machine learning techniques for detecting signal breakpoints is a practical research direction.

13.1.9 Robust Unmixing of Tumor States in Array Comparative Genomic Hybridization Data (Chapter 11)

In Chapter 11 we introduce a geometric framework which provides a novel perspective on detecting cancer subtypes. From a computational perspective, we aim to fit to a cloud of points a minimum volume enclosing simplex.

Complex geometric shapes: Can we fit instead of a simplex, a simplicial complex?

13.1.10 Perfect Reconstruction of Oncogenetic Trees (Chapter 12)

We propose the following open problems on oncogenetic trees.

Enumerative properties of oncogenetic trees: How does one calculate the size of the smallest family \mathcal{F} that is uniquely consistent with a given tree T ? Which trees are

extremal in the aforementioned sense? A broad research direction is to obtain a deeper understanding of tumorigenesis through the understanding of combinatorial properties of oncogenetic trees.

Bibliography

- [1] <http://www.zdnet.com/microsofts-big-data-plans-acknowledge-embrace-integrate-700>
- [2] A map of the interactome network of the metazoan *c. elegans*. *Science (New York, N.Y.)*, 303(5657):540–543, January 2004. URL <http://dx.doi.org/10.1126/science.1091403>.
- [3] <http://incubator.apache.org/giraph/>.
- [4] <http://wiki.apache.org/hadoop/PoweredBy>.
- [5] Snap stanford network analysis library. <http://snap.stanford.edu/>.
- [6] <http://engineering.twitter.com/2011/08/storm-is-coming-more-details-and-plans.html>.
- [7] Microsoft big data. <http://www.microsoft.com/en-us/news/features/2013/feb13/02-11BigData.aspx>.
- [8] James Abello, Mauricio G. C. Resende, and Sandra Sudarsky. Massive quasi-clique detection. In *LATIN*, 2002.
- [9] Bernardo M Ábrego, Silvia Fernández-Merchant, Michael G Neubauer, and William Watkins. Sum of squares of degrees in a graph. *Journal of Inequalities in Pure and Applied Mathematics*, 10(3):1–34, 2009.
- [10] Dimitris Achlioptas, 2010. Random Graphs and Large Networks, available at <http://www.youtube.com/watch?v=arDYNhP95JI>.
- [11] Dimitris Achlioptas and Frank McSherry. Fast computation of low rank matrix approximations. In *Proceedings of the thirty-third annual ACM symposium on Theory of computing*, pages 611–618. ACM, 2001.
- [12] Foto N Afrati, Anish Das Sarma, Semih Salihoglu, and Jeffrey D Ullman. Vision paper: Towards an understanding of the limits of map-reduce computation. *arXiv preprint arXiv:1204.1754*, 2012.

- [13] G. Agarwal and D. Kempe. Modularity-maximizing graph communities via mathematical programming. *The European Physical Journal B-Condensed Matter and Complex Systems*, 66(3):409–418, 2008.
- [14] P. K. Agarwal and J. Erickson. *Geometric Range Searching and Its Relatives*. 1999.
- [15] P. K. Agarwal, D. Eppstein, and J. Matousek. Dynamic half-space reporting, geometric optimization, and minimum spanning trees. In *Proceedings of the 33rd Annual Symposium on Foundations of Computer Science*, pages 80–89, Washington, DC, USA, 1992. IEEE Computer Society. ISBN 0-8186-2900-2. doi: 10.1109/SFCS.1992.267816. URL <http://portal.acm.org/citation.cfm?id=1398516.1398901>.
- [16] A. Aggarwal, M. Klawe, S. Moran, P. Shor, and R. Wilber. Geometric applications of a matrix searching algorithm. In *Proceedings of the second annual symposium on Computational geometry*, SCG '86, pages 285–292, New York, NY, USA, 1986. ACM. ISBN 0-89791-194-6. doi: <http://doi.acm.org/10.1145/10515.10546>. URL <http://doi.acm.org/10.1145/10515.10546>.
- [17] A. Aggarwal, M. Hansen, and T. Leighton. Solving query-retrieval problems by compacting voronoi diagrams. In *Proceedings of the twenty-second annual ACM symposium on Theory of computing*, STOC '90, pages 331–340, New York, NY, USA, 1990. ACM. ISBN 0-89791-361-2. doi: <http://doi.acm.org/10.1145/100216.100260>. URL <http://doi.acm.org/10.1145/100216.100260>.
- [18] Rudolf Ahlswede and Gyula OH Katona. Graphs with maximal number of adjacent pairs of edges. *Acta Mathematica Hungarica*, 32(1):97–120, 1978.
- [19] William Aiello, Fan Chung, and Linyuan Lu. A random graph model for massive graphs. In *Proceedings of the thirty-second annual ACM symposium on Theory of computing*, pages 171–180. Acm, 2000.
- [20] William Aiello, Fan Chung, and Linyuan Lu. A random graph model for power law graphs. *Experimental Mathematics*, 10(1):53–66, 2001.
- [21] Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1):47, 2002.
- [22] Reka Albert, Hawoong Jeong, and Albert-Laszlo Barabasi. Diameter of the world wide web. *Nature*, (401):130–131, 1999.
- [23] D Albertson and D Pinkel. Genomic microarrays in human genetic disease and cancer. *Human Molecular Genetics*, 12:145–152, 2003.

- [24] Farid Alizadeh. Interior point methods in semidefinite programming with applications to combinatorial optimization. *SIAM Journal on Optimization*, 5, 1993.
- [25] N. Alon, R. Yuster, and U Zwick. Finding and counting given length cycles (extended abstract). In *Proceedings of the Second Annual European Symposium on Algorithms*, ESA '94, pages 354–364, London, UK, 1994. Springer-Verlag. ISBN 3-540-58434-X. URL <http://portal.acm.org/citation.cfm?id=647904.739463>.
- [26] Noga Alon and Joel Spencer. *The Probabilistic Method*. John Wiley, 2008.
- [27] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, STOC '96, pages 20–29, New York, NY, USA, 1996. ACM. ISBN 0-89791-785-5. doi: 10.1145/237814.237823. URL <http://doi.acm.org/10.1145/237814.237823>.
- [28] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, pages 20–29. ACM, 1996.
- [29] Noga Alon, Raphael Yuster, and Uri Zwick. Finding and counting given length cycles. *Algorithmica*, 17(3):209–223, 1997.
- [30] J. Ignacio Alvarez-hamelin, Alain Barrat, and Alessandro Vespignani. Large scale networks fingerprinting and visualization using the k-core decomposition. In *NIPS*, 2006.
- [31] Reid Andersen and Kumar Chellapilla. Finding dense subgraphs with size bounds. In *WAW*, 2009.
- [32] Reid Andersen and Yuval Peres. Finding sparse cuts locally using evolving sets. In *STOC '09*, pages 235–244, 2009.
- [33] Reid Andersen, Fan Chung, and Kevin Lang. Local graph partitioning using pagerank vectors. In *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*, pages 475–486. IEEE, 2006.
- [34] RFS Andrade and JGV Miranda. Spectral properties of the apollonian network. *Physica A: Statistical Mechanics and its Applications*, 356(1):1–5, 2005.
- [35] José S Andrade Jr, Hans J Herrmann, Roberto FS Andrade, and Luciano R da Silva. Apollonian networks: Simultaneously scale-free, small world, euclidean, space filling, and with matching graphs. *Physical Review Letters*, 94(1):018702, 2005.

- [36] A. Angel, N. Sarkas, N. Koudas, and D. Srivastava. Dense subgraph maintenance under streaming edge weight updates for real-time story identification. *Proc. VLDB Endow.*, 5(6):574–585, February 2012.
- [37] S. Arora, R. Satish, and U. Vazirani. Expander flows, geometric embeddings and graph partitioning. In *STOC*, pages 222–231, 2004.
- [38] Sanjeev Arora, David Karger, and Marek Karpinski. Polynomial time approximation schemes for dense instances of NP-hard problems. In *STOC*, 1995.
- [39] Sanjeev Arora, Rong Ge, Sushant Sachdeva, and Grant Schoenebeck. Finding overlapping communities in social networks: toward a rigorous approach. In *Proceedings of the 13th ACM Conference on Electronic Commerce*, pages 37–54. ACM, 2012.
- [40] Yuichi Asahiro, Kazuo Iwama, Hisao Tamaki, and Takeshi Tokuyama. Greedily finding a dense subgraph. *J. Algorithms*, 34(2), 2000.
- [41] Yuichi Asahiro, Refael Hassin, and Kazuo Iwama. Complexity of finding dense subgraphs. *Discr. Ap. Math.*, 121(1-3), 2002.
- [42] J. H. Atkins and L. J. Gershell. From the analyst’s couch: Selective anticancer drugs. *Nat Rev Cancer*, 2:645–646, 2002.
- [43] Camille Stephan-Otto Attolini and Franziska Michor. Evolutionary theory of cancer. *Annals of the New York Academy of Sciences*, 1168(1):23–51, 2009.
- [44] H. Avron. Counting triangles in large graphs using randomized matrix trace estimation, 2010.
- [45] Lars Backstrom, Paolo Boldi, Marco Rosa, Johan Ugander, and Sebastiano Vigna. Four degrees of separation. In *WebSci*, pages 33–42, 2012.
- [46] David A. Bader and Kamesh Madduri. A graph-theoretic analysis of the human protein-interaction network using multicore parallel algorithms. *Parallel Comput.*, 2008.
- [47] Z. Bar-Yossef, R. Kumar, and D. Sivakumar. Reductions in streaming algorithms, with an application to counting triangles in graphs. In *Proceedings of the 13th annual ACM-SIAM symposium on Discrete algorithms (SODA)*, pages 623–632, 2002.
- [48] Ziv Bar-Yossef, T. S. Jayram, Ravi Kumar, D. Sivakumar, and Luca Trevisan. Counting distinct elements in a data stream. In *Proceedings of the 6th International Workshop on Randomization and Approximation Techniques*, RANDOM

- '02, pages 1–10, London, UK, UK, 2002. Springer-Verlag. ISBN 3-540-44147-6. URL <http://dl.acm.org/citation.cfm?id=646978.711822>.
- [49] Ziv Bar-Yossef, Ravi Kumar, and D. Sivakumar. Reductions in streaming algorithms, with an application to counting triangles in graphs. In *Proceedings of the thirteenth annual ACM-SIAM symposium on Discrete algorithms*, SODA '02, pages 623–632, Philadelphia, PA, USA, 2002. Society for Industrial and Applied Mathematics. ISBN 0-89871-513-X. URL <http://dl.acm.org/citation.cfm?id=545381.545464>.
- [50] Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.
- [51] D. Barry and J.A. Hartigan. A bayesian analysis for change point problems. *Journal of the American Statistical Association*, 88(421):309–319, 1993. ISSN 01621459. URL <http://www.jstor.org/stable/2290726>.
- [52] Marc Barthélemy. Spatial networks. *Physics Reports*, 499(1):1–101, 2011.
- [53] Vladimir Batagelj and Matjaz Zaversnik. An $O(m)$ algorithm for cores decomposition of networks. *Arxiv*, arXiv.cs/0310049, 2003.
- [54] Vladimir Batagelj and Matjaž Zaveršnik. Short cycle connectivity. *Discrete Mathematics*, 307(3):310–318, 2007.
- [55] Luca Becchetti, Paolo Boldi, Carlos Castillo, and Aristides Gionis. Efficient semi-streaming algorithms for local triangle counting in massive graphs. In *KDD '08: Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 16–24, New York, NY, USA, 2008. ACM. ISBN 978-1-60558-193-4. doi: <http://doi.acm.org/10.1145/1401890.1401898>.
- [56] Niko Beerenwinkel and Seth Sullivant. Markov models for accumulating mutations. *Biometrika*, 96(3):645–661, 2009.
- [57] Niko Beerenwinkel, Jörg Rahnenführer, Martin Däumer, Daniel Hoffmann, Rolf Kaiser, Joachim Selbig, and Thomas Lengauer. Learning multiple evolutionary pathways from cross-sectional data. In *Proceedings of the eighth annual international conference on Resaerch in computational molecular biology*, pages 36–44. ACM, 2004.
- [58] Niko Beerenwinkel, Jörg Rahnenführer, Rolf Kaiser, Daniel Hoffmann, Joachim Selbig, and Thomas Lengauer. Mtreemix: a software package for learning and using mixture models of mutagenetic trees. *Bioinformatics*, 21(9):2106–2107, 2005.

- [59] Niko Beerenwinkel, Nicholas Eriksson, and Bernd Sturmfels. Conjunctive bayesian networks. *Bernoulli*, 13(4):893–909, 2007.
- [60] W. Bein, M.J. Golin, L. Larmore, and Y. Zhang. The knuth-yao quadrangle-inequality speedup is a consequence of total monotonicity. *ACM Trans. Algorithms*, 6:1–22, December 2009. ISSN 1549-6325. doi: <http://doi.acm.org/10.1145/1644015.1644032>. URL <http://doi.acm.org/10.1145/1644015.1644032>.
- [61] B.A. Bejjani, R. Saleki, B.C. Ballif, E.A. Rorem, K. Sundin, A. Theisen, C.D. Kashork, and L.G. Shaffer. Use of targeted array-based cgh for the clinical diagnosis of chromosomal imbalance: is less more? *American Journal of Medical Genetics A*, pages 259–67, 2005.
- [62] R.E. Bellman. *Dynamic Programming*. Dover Publications, 2003.
- [63] András A. Benczúr and David R. Karger. Randomized approximation schemes for cuts and flows in capacitated graphs. *CoRR*, cs.DS/0207078, 2002.
- [64] András A. Benczúr and David R. Karger. Approximating s-t minimum cuts in $\tilde{O}(n^2)$ time. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, STOC '96, pages 47–55, New York, NY, USA, 1996. ACM. ISBN 0-89791-785-5. doi: <http://doi.acm.org/10.1145/237814.237827>. URL <http://doi.acm.org/10.1145/237814.237827>.
- [65] Jonathan W Berry, Bruce Hendrickson, Randall A LaViolette, and Cynthia A Phillips. Tolerating the community detection resolution limit with edge weighting. *Physical Review E*, 83(5):056119, 2011.
- [66] D. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 2000.
- [67] Kevin Beyer, Peter J. Haas, Berthold Reinwald, Yannis Sismanis, and Rainer Gemulla. On synopses for distinct-value estimation under multiset operations. In *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*, SIGMOD '07, pages 199–210, New York, NY, USA, 2007. ACM. ISBN 978-1-59593-686-8. doi: 10.1145/1247480.1247504. URL <http://doi.acm.org/10.1145/1247480.1247504>.
- [68] S. Bhamidi, G. Bresler, and A. Sly. Mixing time of exponential random graphs. In *Proceedings of the 49th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 803–812, 2008.
- [69] Krishna Bharat, Bay-Wei Chang, Monika Rauch Henzinger, and Matthias Ruhl. Who links to whom: Mining linkage between web sites. In *Proceedings of the*

- 2001 *IEEE International Conference on Data Mining*, ICDM '01, pages 51–58, Washington, DC, USA, 2001. IEEE Computer Society. ISBN 0-7695-1119-8. URL <http://dl.acm.org/citation.cfm?id=645496.757722>.
- [70] G. Bignell, J. Huang, J. Greshock, S. Watt, A. Butler, S. West, M. Grigorova, K. Jones, W. Wei, M. Stratton, A. Futreal, B. Weber, M. Shapero, and R. Wooster. High-resolution analysis of dna copy number using oligonucleotide microarrays. *Genome Research*, pages 287–295, 2004.
- [71] A. H. Bild, A. Potti, and J. R. Nevins. Opinion: Linking oncogenic pathways with therapeutic opportunities. *Nat Rev Cancer*, 6:735–741, 2006.
- [72] Olivier Bodini, Alexis Darrasse, and Michele Soria. Distances in random apollonian network structures. *arXiv preprint arXiv:0712.2129*, 2007.
- [73] Paolo Boldi, Massimo Santini, and Sebastiano Vigna. Permuting web and social graphs. *Internet Mathematics*, 6(3):257–283, 2009.
- [74] Paolo Boldi, Marco Rosa, and Sebastiano Vigna. Hyperanf: approximating the neighbourhood function of very large graphs on a budget. In *Proceedings of the 20th international conference on World wide web*, WWW '11, pages 625–634, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0632-4. doi: 10.1145/1963405.1963493. URL <http://doi.acm.org/10.1145/1963405.1963493>.
- [75] Paolo Boldi, Marco Rosa, and Sebastiano Vigna. Hyperanf: approximating the neighbourhood function of very large graphs on a budget. In *Proceedings of the 20th international conference on World wide web*, WWW '11, pages 625–634, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0632-4. doi: 10.1145/1963405.1963493. URL <http://doi.acm.org/10.1145/1963405.1963493>.
- [76] Béla Bollobás. *Modern graph theory*, volume 184. Springer Verlag, 1998.
- [77] Béla Bollobás. *Random graphs*, volume 73. Cambridge university press, 2001.
- [78] Béla Bollobás and Oliver Riordan. Robustness and vulnerability of scale-free random graphs. *Internet Mathematics*, 1(1):1–35, 2004.
- [79] Béla Bollobás and Oliver Riordan. The diameter of a scale-free random graph. *Combinatorica*, 24(1):5–34, 2004.
- [80] Béla Bollobás, Oliver Riordan, Joel Spencer, Gábor Tusnády, et al. The degree sequence of a scale-free random graph process. *Random Structures & Algorithms*, 18(3):279–290, 2001.
- [81] A. Bonato. *A course on the Web graph*. American Mathematical Society, 2008.

- [82] Anthony Bonato. A survey of models of the web graph. In *Combinatorial and algorithmic aspects of networking*, pages 95–95. Springer, 2005.
- [83] John Adrian Bondy and Uppaluri Siva Ramachandra Murty. *Graph theory with applications*, volume 290. Macmillan London, 1976.
- [84] Christian Borgs, Jennifer Chayes, Constantinos Daskalakis, and Sebastien Roch. First to market is not everything: an analysis of preferential attachment with fitness. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 135–144. ACM, 2007.
- [85] Christian Borgs, Jennifer Chayes, Jian Ding, and Brendan Lucier. The hitchhiker’s guide to affiliation networks: A game-theoretic approach. *arXiv preprint arXiv:1008.1516*, 2010.
- [86] David W Boyd. The sequence of radii of the apollonian packing. *mathematics of computation*, 39(159):249–254, 1982.
- [87] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, New York, NY, USA, 2004. ISBN 0521833787.
- [88] U. Brandes, D. Delling, M. Gaertler, R. Görke, M. Hoefer, Z. Nikoloski, and D. Wagner. On finding graph clusterings with maximum modularity. In *Graph-Theoretic Concepts in Computer Science*, pages 121–132. Springer, 2007.
- [89] S. Brin and L. Page. The anatomy of a large-scale hypertextual (web) search engine. In *Proc. 7th International World Wide Web Conference (WWW7)/Computer Networks*, pages 107–117, 1998. Published as Proc. 7th International World Wide Web Conference (WWW7)/Computer Networks, volume 30, number 1-7.
- [90] Andrei Broder, Ravi Kumar, Farzin Maghoul, Prabhakar Raghavan, Sridhar Rajagopalan, Raymie Stata, Andrew Tomkins, and Janet Wiener. Graph structure in the web. *Computer networks*, 33(1):309–320, 2000.
- [91] Andrei Z. Broder, Moses Charikar, Alan M. Frieze, and Michael Mitzenmacher. Min-wise independent permutations. *Journal of Computer and System Sciences*, pages 630–659, 2000.
- [92] Coen Bron and Joep Kerbosch. Algorithm 457: finding all cliques of an undirected graph. *CACM*, 16(9), 1973.
- [93] Nicolas Broutin and Luc Devroye. Large deviations for the weighted height of an extended class of trees. *Algorithmica*, 46(3-4):271–297, 2006.

- [94] Mauro Brunato, Holger H. Hoos, and Roberto Battiti. On effectively finding maximal quasi-cliques in graphs. In *Learning and Intelligent Optimization*. 2008.
- [95] Gregory Buehrer and Kumar Chellapilla. A scalable pattern mining approach to web graph compression with communities. In *WSDM*, 2008.
- [96] Luciana S Buriol, Gereon Frahling, Stefano Leonardi, Alberto Marchetti-Spaccamela, and Christian Sohler. Counting triangles in data streams. In *Proceedings of the twenty-fifth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 253–262. ACM, 2006.
- [97] R. E. Burkard, B. Klinz, and R. Rüdiger. Perspectives of monge properties in optimization. *Discrete Appl. Math.*, 70:95–161, September 1996. ISSN 0166-218X. doi: 10.1016/0166-218X(95)00103-X. URL <http://portal.acm.org/citation.cfm?id=240858.240860>.
- [98] M.G. Butler, F.J. Meaney, and C.G. Palmer. Clinical and cytogenetic survey of 39 individuals with prader-labhart-willi syndrome. *American Journal of Medical Genetics*, 23:793–809, 1986.
- [99] Alessio Cardillo, Salvatore Scellato, Vito Latora, and Sergio Porta. Structural properties of planar graphs of urban street patterns. *Physical Review E*, 73(6):066107, 2006.
- [100] Yair Caro, Arie Lev, Yehuda Roditty, Zsolt Tuza, and Raphael Yuster. On rainbow connection. *Electron. J. Combin*, 15(1):R57, 2008.
- [101] Sourav Chakraborty, Eldar Fischer, Arie Matsliah, and Raphael Yuster. Hardness and algorithms for rainbow connection. *Journal of Combinatorial Optimization*, 21(3):330–347, 2011.
- [102] T.H. Chan, C.Y. Chi, Y.M. Huang, and W.K. Ma. A convex analysis based minimum-volume enclosing simplex algorithm for hyperspectral unmixing. *IEEE Trans. Signal Processing*, 57(11):4418–4432, 2009.
- [103] L Sunil Chandran, Anita Das, Deepak Rajendraprasad, and Nithin M Varma. Rainbow connection number and connected dominating sets. *Electronic Notes in Discrete Mathematics*, 38:239–244, 2011.
- [104] Moses Charikar. Greedy approximation algorithms for finding dense components in a graph. In *APPROX*, 2000.
- [105] Gary Chartrand, Garry L Johns, Kathleen A McKeon, and Ping Zhang. Rainbow connection in graphs. *Mathematica Bohemica*, 133(1):85–98, 2008.

- [106] B. Chazelle and F.P. Preparata. Halfspace range search: an algorithmic application of k-sets. In *Proceedings of the first annual symposium on Computational geometry*, SCG '85, pages 107–115, New York, NY, USA, 1985. ACM. ISBN 0-89791-163-6. doi: <http://doi.acm.org/10.1145/323233.323248>. URL <http://doi.acm.org/10.1145/323233.323248>.
- [107] B. Chazelle, L.J. Guibas, and D. T. Lee. The power of geometric duality. *BIT*, 25: 76–90, June 1985. ISSN 0006-3835. doi: <http://dx.doi.org/10.1007/BF01934990>. URL <http://dx.doi.org/10.1007/BF01934990>.
- [108] Norishige Chiba and Takao Nishizeki. Arboricity and subgraph listing algorithms. *SIAM Journal on Computing*, 14(1):210–223, 1985.
- [109] Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, Michael Mitzenmacher, Alessandro Panconesi, and Prabhakar Raghavan. On compressing social networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 219–228. ACM, 2009.
- [110] Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, Alessandro Panconesi, and Prabhakar Raghavan. Models for the compressible web. *SIAM Journal on Computing*, 42(5):1777–1802, 2013.
- [111] Karen Stromme Christensen. *Cities and complexity*. Sage Publications Thousand Oaks, CA, 1999.
- [112] N. Chunder, S. Mandal, D. Basu, A. Roy, S. Roychoudhury, and C. K. Panda. Deletion mapping of chromosome 1 in early onset and late onset breast tumors—a comparative study in eastern india. *Pathol Res Pract*, 199:313–321, 2003.
- [113] Fan Chung, Linyuan Lu, and Van Vu. Spectra of random graphs with given expected degrees. *Proceedings of the National Academy of Sciences*, 100(11):6313–6318, 2003.
- [114] Fan R. K. Chung and Linyuan Lu. The average distance in a random graph with given expected degrees. *Internet Mathematics*, 1(1), 2003.
- [115] Fan RK Chung. Spectral graph theory (cbms regional conference series in mathematics, no. 92). 1997.
- [116] Colin Clark. Urban population densities. *Journal of the Royal Statistical Society. Series A (General)*, 114(4):490–496, 1951.

- [117] K. L. Clarkson and P. W. Shor. Applications of random sampling in computational geometry, ii. *Discrete Comput. Geom.*, 4:387–421, September 1989. ISSN 0179-5376. doi: <http://dx.doi.org/10.1007/BF02187740>. URL <http://dx.doi.org/10.1007/BF02187740>.
- [118] Aaron Clauset, Cosma Rohilla Shalizi, and Mark EJ Newman. Power-law distributions in empirical data. *SIAM review*, 51(4):661–703, 2009.
- [119] Edith Cohen. Size-estimation framework with applications to transitive closure and reachability. *J. Comput. Syst. Sci.*, 55(3):441–453, December 1997. ISSN 0022-0000. doi: 10.1006/jcss.1997.1534. URL <http://dx.doi.org/10.1006/jcss.1997.1534>.
- [120] Jonathan Cohen. Graph twiddling in a mapreduce world. *Computing in Science & Engineering*, 11(4):29–41, 2009.
- [121] P. Comon. Independent component analysis. *Signal Processing*, 36:287–314, 1994.
- [122] A. Condon and R. M. Karp. Algorithms for graph partitioning on the planted partition model. In *RANDOM-APPROX*, pages 221–232, 1999.
- [123] Colin Cooper and Alan Frieze. A general model of web graphs. *Random Structures & Algorithms*, 22(3):311–335, 2003.
- [124] Colin Cooper and Ryuhei Uehara. Scale free properties of random k-trees. *Mathematics in Computer Science*, 3(4):489–496, 2010.
- [125] T. Cormen, C. Leiserson, and R. Rivest. *Introduction to Algorithms*. The MIT Press, 1990.
- [126] T. Corson, A. Huang, M.S. Tsao, and B.L. Gallie. Kif14 is a candidate oncogene in the 1q minimal region of genomic gain in multiple cancers. *Oncogene*, 24:47414753, 2005.
- [127] P. Costa, A. Donnelly, A. Rowstron, and G. O’Shea. Camdoop: exploiting in-network aggregation for big data applications. In *NSDI’12*, pages 3–3, 2012.
- [128] Alexis Darrasse and Michele Soria. Degree distribution of random apollonian network structures and boltzmann sampling. *DMTCS Proceedings*, (1), 2008.
- [129] Alexis Darrasse, Hsien-Kuei Hwang, Olivier Bodini, and Michele Soria. The connectivity-profile of random increasing k-trees. *arXiv preprint arXiv:0910.3639*, 2009.

- [130] J. Dean and S. Ghemawat. Mapreduce: simplified data processing on large clusters. *Commun. ACM*, 51(1):107–113, January 2008. ISSN 0001-0782. doi: 10.1145/1327452.1327492. URL <http://doi.acm.org/10.1145/1327452.1327492>.
- [131] R. Desper, F. Jiang, O. P. Kallioniemi, H. Moch, C. H. Papadimitriou, and A. A. Schaffer. Inferring tree models for oncogenesis from comparative genome hybridization data. *J Comp Biol*, 6:37–51, 1999.
- [132] Richard Desper, Feng Jiang, Olli-P. Kallioniemi, Holger Moch, Christos H. Papadimitriou, and Alejandro A. Schäffer. Inferring tree models for oncogenesis from comparative genome hybridization data. *Journal of Computational Biology*, 6(1):37–52, 1999.
- [133] Richard Desper, Feng Jiang, Olli-P. Kallioniemi, Holger Moch, Christos H. Papadimitriou, and Alejandro A. Schäffer. Inferring tree models for oncogenesis from comparative genome hybridization data. *Journal of Computational Biology*, 6(1):37–51, 1999.
- [134] Richard Desper, Feng Jiang, Olli-P. Kallioniemi, Holger Moch, Christos H. Papadimitriou, and Alejandro A. Schäffer. Distance-based reconstruction of tree models for oncogenesis. *Journal of Computational Biology*, 7(6):789–803, 2000.
- [135] Richard Desper, Feng Jiang, Olli-P. Kallioniemi, Holger Moch, Christos H. Papadimitriou, and Alejandro A. Schäffer. Distance-based reconstruction of tree models for oncogenesis. *Journal of Computational Biology*, 7(6):789–803, 2000.
- [136] Jiarui Ding and Sohrab P. Shah. Robust hidden semi-markov modeling of array cgh data. In *Bioinformatics and Biomedicine (BIBM), 2010 IEEE International Conference on*, pages 603–608. IEEE, 2010.
- [137] Debora Donato, Luigi Laura, Stefano Leonardi, and Stefano Millozzi. The web as a graph: How far we are. *ACM Trans. Internet Technol.*, 7(1), February 2007. ISSN 1533-5399. doi: 10.1145/1189740.1189744. URL <http://doi.acm.org/10.1145/1189740.1189744>.
- [138] B. Dorow, D. Widdows, L. Katarina, J.-P. Eckmann, S. Danilo, and E. Moses. Using curvature and markov clustering in graphs for lexical acquisition and word sense discrimination. In *In Workshop MEANING-2005*, 2004.
- [139] Jonathan PK Doye and Claire P. Massen. Self-similar disk packings as model spatial scale-free networks. *arXiv preprint cond-mat/0407779*, 2004.
- [140] P. Drineas, Alan Frieze, Ravi Kannan, Santosh Vempala, and V. Vinay. Clustering in large graphs and matrices. In *Proceedings of the tenth annual ACM-SIAM*

- symposium on Discrete algorithms*, SODA '99, pages 291–299, Philadelphia, PA, USA, 1999. Society for Industrial and Applied Mathematics. ISBN 0-89871-434-6. URL <http://portal.acm.org/citation.cfm?id=314500.314576>.
- [141] Xiaoxi Du, Ruoming Jin, Liang Ding, Victor E. Lee, and John H. Thornton, Jr. Migration motif: a spatial - temporal pattern mining approach for financial markets. In *KDD*, 2009.
- [142] Marianne Durand and Philippe Flajolet. Loglog counting of large cardinalities (extended abstract). In *ESA*, pages 605–617, 2003.
- [143] Richard Durrett. *Random graph dynamics*, volume 20. Cambridge university press, 2007.
- [144] David Easley and Jon Kleinberg. *Networks, crowds, and markets*. Cambridge Univ Press, 2010.
- [145] E. Eckmann, J.-P. and Moses. Curvature of co-links uncovers hidden thematic layers in the world wide web. *PNAS*, 99(9):5825–5829, April 2002. doi: 10.1073/pnas.032093399. URL <http://dx.doi.org/10.1073/pnas.032093399>.
- [146] R. Ehrlich and W. Full. Sorting out geology — unmixing mixtures. In *Use and Abuse of Statistical Methods in the Earth Sciences*, pages 33–46. Oxford University Press, 1987.
- [147] P. H. Eilers and R. X. de Menezes. Quantile smoothing of array cgh data. *Bioinformatics*, 21(7):1146–1153, April 2005. ISSN 1367-4803. URL <http://view.ncbi.nlm.nih.gov/pubmed/15572474>.
- [148] D. Eppstein, Z. Galil, and R. Giancarlo. Speeding up dynamic programming. In *Proceedings of the 29th Annual Symposium on Foundations of Computer Science*, pages 488–496, Washington, DC, USA, 1988. IEEE Computer Society. ISBN 0-8186-0877-3. doi: 10.1109/SFCS.1988.21965. URL <http://portal.acm.org/citation.cfm?id=1398513.1398584>.
- [149] D. Eppstein, Z. Galil, R. Giancarlo, and G.F. Italiano. Sparse dynamic programming i: linear cost functions. *J. ACM*, 39:519–545, July 1992. ISSN 0004-5411. doi: <http://doi.acm.org/10.1145/146637.146650>. URL <http://doi.acm.org/10.1145/146637.146650>.
- [150] D. Eppstein, Z. Galil, R. Giancarlo, and G.F. Italiano. Sparse dynamic programming ii: convex and concave cost functions. *J. ACM*, 39:546–567, July 1992. ISSN 0004-5411. doi: <http://doi.acm.org/10.1145/146637.146656>. URL <http://doi.acm.org/10.1145/146637.146656>.

- [151] David Eppstein, Maarten Löffler, and Darren Strash. Listing all maximal cliques in sparse graphs in near-optimal time. In *ISAAC*, 2010.
- [152] Paul Erdős and A Rényi. On the evolution of random graphs. *Publ. Math. Inst. Hungar. Acad. Sci.*, 5:17–61, 1960.
- [153] P. Erdős and A. Rényi. On random graphs i. *Publicationes Mathematicae*, 6: 290297, 1959.
- [154] AB Ericksen. A matter of security. *Graduating Engineer & Computer Careers*, pages 24–28, 2007.
- [155] R. Etzioni, S. Hawley, D. Billheimer, L. D. True, and B. Knudsen. Analyzing patterns of staining in immunohistochemical studies: application to a study of prostate cancer recurrence. *Cancer Epidem Biomark*, 14:1040–1046, 2005.
- [156] Logan Everett, Li-San Wang, and Sridhar Hannenhalli. Dense subgraph computation via stochastic search: application to detect transcriptional modules. In *ISMB*, 2006.
- [157] Alex Fabrikant, Elias Koutsoupias, and Christos H Papadimitriou. Heuristically optimized trade-offs: A new paradigm for power laws in the internet. In *Automata, languages and programming*, pages 110–122. Springer, 2002.
- [158] Michalis Faloutsos, Petros Faloutsos, and Christos Faloutsos. On power-law relationships of the internet topology. In *ACM SIGCOMM Computer Communication Review*, volume 29, pages 251–262. ACM, 1999.
- [159] Eric R Fearon, Bert Vogelstein, et al. A genetic model for colorectal tumorigenesis. *Cell*, 61(5):759–767, 1990.
- [160] U. Feige and R. Krauthgamer. A polylogarithmic approximation of the minimum bisection. *SIAM J. Comput.*, 31(4):1090–1118, April 2002. ISSN 0097-5397.
- [161] Uriel Feige. Approximating maximum clique by removing subgraphs. *SIAM Journal of Discrete Mathematics*, 18(2), 2005.
- [162] Uriel Feige and Michael Langberg. Approximation algorithms for maximization problems arising in graph partitioning. *J. Algorithms*, 41(2), 2001.
- [163] Uriel Feige, Guy Kortsarz, and David Peleg. The dense k-subgraph problem. *Algorithmica*, 29(3), 2001.
- [164] Joan Feigenbaum, Sampath Kannan, Andrew McGregor, Siddharth Suri, and Jian Zhang. On graph problems in a semi-streaming model. *Theoretical Computer Science*, 348(2):207–216, 2005.

- [165] Vitaly Feldman, Elena Grigorescu, Lev Reyzin, Santosh Vempala, and Ying Xiao. Statistical algorithms and a lower bound for detecting planted cliques. In *Proceedings of the 45th annual ACM symposium on Symposium on theory of computing*, pages 655–664. ACM, 2013.
- [166] J.-A. Ferrez, K. Fukuda, and Th.M. Liebling. Parallel computation of the diameter of a graph. In *HPCSA*, 1998.
- [167] C. M. Fiduccia and R. M. Mattheyses. A linear-time heuristic for improving network partitions. In *DAC '82*, pages 175–181, 1982.
- [168] Philippe Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. *Journal of Computer and System Sciences*, 1985.
- [169] Philippe Flajolet, ric Fusy, Olivier Gandouet, and et al. Hyperloglog: The analysis of a near-optimal cardinality estimation algorithm. In *IN AOFA 07: PROCEEDINGS OF THE 2007 INTERNATIONAL CONFERENCE ON ANALYSIS OF ALGORITHMS*, 2007.
- [170] Gary William Flake, Steve Lawrence, and C. Lee Giles. Efficient identification of web communities. In *KDD*, 2000.
- [171] Abraham Flaxman, Alan Frieze, and Trevor Fenner. High degree vertices and eigenvalues in the preferential attachment graph. *Internet Mathematics*, 2(1):1–19, 2005.
- [172] Abraham D Flaxman, Alan M Frieze, and Juan Vera. A geometric preferential attachment model of networks. *Internet Mathematics*, 3(2):187–205, 2006.
- [173] Abraham D Flaxman, Alan M Frieze, and Juan Vera. A geometric preferential attachment model of networks ii. *Internet Mathematics*, 4(1):87–111, 2007.
- [174] Santo Fortunato. Community detection in graphs. *Physics Reports*, 486(3-5), 2010.
- [175] R. M. Foster. The average impedance of an electrical network. *Contributions to Applied Mechanics (Reissner Anniversary Volume)*, page 333340, 1949.
- [176] Ove Frank and David Strauss. Markov graphs. *Journal of the American Statistical Association*, 81(395):832–842, 1986. URL <http://www.jstor.org/stable/2289017>.
- [177] Eugene Fratkin, Brian T. Naughton, Douglas L. Brutlag, and Serafim Batzoglou. MotifCut: regulatory motifs finding with maximum density subgraphs. In *ISMB*, 2006.

- [178] Jane Fridlyand, Antoine M. Snijders, Dan Pinkel, Donna G. Albertson, and Ajay N. Jain. Hidden Markov models approach to the analysis of array CGH data. *J. Multivar. Anal.*, 90(1):132–153, 2004. ISSN 0047-259X. doi: <http://dx.doi.org/10.1016/j.jmva.2004.02.008>.
- [179] A. Frieze and C.E Tsourakakis. Some properties of random apollonian networks. *Internet Mathematics*, 2013.
- [180] A. Frieze, A. Gionis, and C.E Tsourakakis. Algorithmic techniques for modeling and mining large graphs (amazing). In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining (to appear)*. ACM, 2013.
- [181] A. M. Frieze and M. Jerrum. Improved approximation algorithms for max k-cut and max bisection. *Algorithmica*, 18(1):67–81, 1997.
- [182] Alan M. Frieze and Charalampos E. Tsourakakis. Rainbow connectivity of sparse random graphs. In *APPROX-RANDOM*, pages 541–552, 2012.
- [183] Alan M. Frieze and Charalampos E. Tsourakakis. On certain properties of random apollonian networks. In *WAW*, pages 93–112, 2012.
- [184] Alan M. Frieze and Charalampos E. Tsourakakis. Rainbow connection of sparse random graphs. *The Electronic Journal of Combinatorics*, 19, 2012.
- [185] Ioannis Fudos and Christoph M. Hoffmann. A graph-constructive approach to solving systems of geometric constraints. *ACM Trans. Graph.*, 16(2):179–216, 1997.
- [186] G. Gallo, M. D. Grigoriadis, and R. E. Tarjan. A fast parametric maximum flow algorithm and applications. *Journal of Computing*, 18(1), 1989.
- [187] Y Gao and C Hobson. Random k-tree as a model for complex networks. In *Workshop on Algorithms and Models for the Web-Graph, WAW*, 2006.
- [188] Yong Gao. The degree distribution of random k-trees. *Theoretical Computer Science*, 410(8):688–695, 2009.
- [189] M. R. Garey, D. S. Johnson, and L. Stockmeyer. Some simplified np-complete problems. In *STOC '74*, pages 47–63, 1974.
- [190] Alan F Gates, Olga Natkovich, Shubham Chopra, Pradeep Kamath, Shravan M Narayanamurthy, Christopher Olston, Benjamin Reed, Santhosh Srinivasan, and Utkarsh Srivastava. Building a high-level dataflow system on top of map-reduce: the pig experience. *Proceedings of the VLDB Endowment*, 2(2):1414–1425, 2009.

- [191] Moritz Gerstung, Michael Baudis, Holger Moch, and Niko Beerenwinkel. Quantifying cancer progression with conjunctive Bayesian networks. *Bioinformatics*, 25(21):2809–2815, 2009.
- [192] Moritz Gerstung, Michael Baudis, Holger Moch, and Niko Beerenwinkel. Quantifying cancer progression with conjunctive bayesian networks. *Bioinformatics*, 25(21):2809–2815, 2009.
- [193] David Gibson, Ravi Kumar, and Andrew Tomkins. Discovering large dense subgraphs in massive graphs. In *VLDB*, 2005.
- [194] Edgar N Gilbert. Random graphs. *The Annals of Mathematical Statistics*, 30(4):1141–1144, 1959.
- [195] E.N. Gilbert and E.F. Moore. Variable-length binary encodings. *Bell System Tech.*, 38:933–966, July 1959.
- [196] M. Girvan and M. E. J. Newman. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences*, 99(12):7821–7826, 2002.
- [197] David F. Gleich and C. Seshadhri. Vertex neighborhoods, low conductance cuts, and good seeds for local community methods. In *KDD2012*, pages 597–605, 2012.
- [198] Michel X. Goemans and David P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *JACM*, 42(6), 1995.
- [199] A. V. Goldberg. Finding a maximum density subgraph. Technical report, University of California at Berkeley, 1984.
- [200] T. R. Golub, D. K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. P. Mesirov, H. Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, C. D. Bloomfield, and E. S. Lander. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science*, 286:531–537, 1999.
- [201] J. Gonzalez, Y. Low, H. Gu, D. Bickson, and C. Guestrin. Powergraph: Distributed graph-parallel computation on natural graphs. In *OSDI*, 2012.
- [202] Ronald L Graham, Jeffrey C Lagarias, Colin L Mallows, Allan R Wilks, and Catherine H Yan. Apollonian circle packings: number theory. *Journal of Number Theory*, 100(1):1–45, 2003.
- [203] A. Greenberg and et al. V12: a scalable and flexible data center network. In *SIGCOMM '09*, 2009.

- [204] S. Guha, N. Koudas, and K. Shim. Approximation and streaming algorithms for histogram construction problems. *ACM Trans. Database Syst.*, 31:396–438, March 2006. ISSN 0362-5915. doi: <http://doi.acm.org/10.1145/1132863.1132873>. URL <http://doi.acm.org/10.1145/1132863.1132873>.
- [205] S. Guha, Y. Li, and D. Neuberg. Bayesian hidden Markov modeling of array CGH data. *Harvard University, Paper 24*, 2006. URL <http://www.bepress.com/harvardbiostat/paper24/>.
- [206] S. Guha, Y. Li, and D. Neuberg. Bayesian hidden markov modeling of array cgh data. *Harvard University, Paper 24*, 2006. URL <http://www.bepress.com/harvardbiostat/paper24/>.
- [207] Katrin Hainke, Jörg Rahnenführer, and Roland Fried. Disease progression models: A review and comparison. Technical report, Dortmund University, Technical Report, 2011.
- [208] András Hajnal and Endre Szemerédi. Proof of a conjecture of erdős. *Combinatorial theory and its applications*, 2:601–623, 1970.
- [209] Jing He and Hongyu Liang. On rainbow- k -connectivity of random graphs. *Information Processing Letters*, 112(10):406–410, 2012.
- [210] Fritz Heider. Attitudes and cognitive organization. *The Journal of psychology*, 21(1):107–112, 1946.
- [211] D.S. Hirschberg. A linear space algorithm for computing maximal common subsequences. *Commun. ACM*, 18:341–343, June 1975. ISSN 0001-0782. doi: <http://doi.acm.org/10.1145/360825.360861>. URL <http://doi.acm.org/10.1145/360825.360861>.
- [212] J. Håstad. Clique is hard to approximate within $n^{1-\epsilon}$. *Acta Mathematica*, 182(1), 1999.
- [213] JG Hodgson, K Chin, C Collins, and JW Gray. Genome amplification of chromosome 20 in breast cancer. *Breast Cancer Res Treat.*, 2003.
- [214] L. Hsu, S. Self, D. Grove, K. Wang, J. Delrow, L. Loo, and P. Porter. Denoising array based comparative genomic hybridization data using wavelets. *Biostatistics*, 6:211–226, 2005.
- [215] Qiang Huang, Guo Pei Yu, Steven A McCormick, Juan Mo, Bhakti Datta, Manoj Mahimkar, Philip Lazarus, Alejandro A Schäffer, Richard Desper, and Stimson P Schantz. Genetic differences detected by comparative genomic hybridization in

- head and neck squamous cell carcinomas from different tumor sites: construction of oncogenetic trees for tumor progression. *Genes, Chromosomes and Cancer*, 34(2):224–233, 2002.
- [216] Philippe Hupé, Nicolas Stransky, Jean-Paul Thiery, François Radvanyi, and Emmanuel Barillot. Analysis of array cgh data: from signal ratio to gain and loss of dna regions. *Bioinformatics*, 20(18):3413–3422, 2004. ISSN 1367-4803. doi: <http://dx.doi.org/10.1093/bioinformatics/bth418>.
- [217] Soon hyung Yook, Zoltn N. Oltvai, and Albert Iszl Barabasi. Functional and topological characterization of protein interaction networks. *Proteomics*, 4:928–942, 2004.
- [218] M Hglund, D Gisselsson, N Mandah, B Johansson, F Mertens, F Mitelman, and T Sll. Multivariate analyses of genomic imbalances in solid tumors reveal distinct and converging pathways of karyotypic evolution. *Genes, Chromosomes and Cancer*, 2001.
- [219] Piotr Indyk and David P. Woodruff. Tight lower bounds for the distinct elements problem. In *FOCS*, pages 283–288, 2003.
- [220] Alon Itai and Michael Rodeh. Finding a minimum circuit in a graph. *SIAM Journal on Computing*, 7(4):413–423, 1978.
- [221] Matthew O Jackson. *Social and economic networks*. Princeton University Press, 2010.
- [222] H. V. Jagadish, N. Koudas, S. Muthukrishnan, V. Poosala, K. C. Sevcik, and S. Torsten. Optimal histograms with quality guarantees. In *Proceedings of the 24rd International Conference on Very Large Data Bases, VLDB '98*, pages 275–286, San Francisco, CA, USA, 1998. Morgan Kaufmann Publishers Inc. ISBN 1-55860-566-5. URL <http://portal.acm.org/citation.cfm?id=645924.671191>.
- [223] Svante Janson, Tomasz Łuczak, and Andrzej Ruciński. Random graphs. 2000. *Wiley–Intersci. Ser. Discrete Math. Optim*, 2000.
- [224] Mark R Jerrum, Leslie G Valiant, and Vijay V Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theoretical Computer Science*, 43:169–188, 1986.
- [225] Bin Jiang. A topological pattern of urban street networks: universality and peculiarity. *Physica A: Statistical Mechanics and its Applications*, 384(2):647–655, 2007.

- [226] Ruoming Jin, Yang Xiang, Ning Ruan, and David Fuhry. 3-hop: a high-compression indexing scheme for reachability query. In *SIGMOD*, 2009.
- [227] William B Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a hilbert space. *Contemporary mathematics*, 26(189-206):1, 1984.
- [228] Kees Jong, Elena Marchiori, Gerrit Meijer, Aad van der Vaart, and Bauke Ylstra. Breakpoint identification and smoothing of array comparative genomic hybridization data. *Bioinformatics*, 20(18):3636–3637, 2004.
- [229] JáJ'a Joseph. *An Introduction to Parallel Algorithms*. Addison-Wesley, 1992.
- [230] Hossein Jowhari and Mohammad Ghodsi. New streaming algorithms for counting triangles in graphs. In *Computing and Combinatorics*, pages 710–716. Springer, 2005.
- [231] A. Kallioniemi, O.P. Kallioniemi, D. Sudar, D. Rutovitz, J.W. Gray, F. Waldman, and D. Pinkel. Comparative genomic hybridization for molecular cytogenetic analysis of solid tumors. *Science*, 258(5083):818–821, 1992. URL <http://www.sciencemag.org/cgi/content/abstract/258/5083/818>.
- [232] Gabriela Kalna and Desmond J. Higham. A clustering coefficient for weighted networks, with application to gene expression data. *AI Commun.*, 20:263–271, December 2007. ISSN 0921-7126. URL <http://portal.acm.org/citation.cfm?id=1365534.1365536>.
- [233] A. Kamb, S. Wee, and C. Lengauer. Why is cancer drug discovery so difficult? *Nat Rev Drug Discov*, 6:115–120, 2007.
- [234] Daniel M. Kane, Jelani Nelson, and David P. Woodruff. An optimal algorithm for the distinct elements problem. In *Proceedings of the twenty-ninth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, PODS '10, pages 41–52, New York, NY, USA, 2010. ACM. ISBN 978-1-4503-0033-9. doi: 10.1145/1807085.1807094. URL <http://doi.acm.org/10.1145/1807085.1807094>.
- [235] Daniel M Kane, Kurt Mehlhorn, Thomas Sauerwald, and He Sun. Counting arbitrary subgraphs in data streams. In *Automata, Languages, and Programming*, pages 598–609. Springer, 2012.
- [236] U. Kang, C. E. T., and C. Faloutsos. Pegasus: A peta-scale graph mining system. In *ICDM*, pages 229–238, 2009.
- [237] U. Kang, Charalampos E. Tsourakakis, Ana Paula Appel, Christos Faloutsos, and Jure Leskovec. Radius plots for mining tera-byte scale graphs: Algorithms, patterns, and observations. In *SDM*, pages 548–558, 2010.

- [238] U. Kang, Charalampos E. Tsourakakis, Ana Paula Appel, Christos Faloutsos, and Jure Leskovec. Radius plots for mining tera-byte scale graphs: Algorithms, patterns, and observations. In *SDM*, pages 548–558, 2010.
- [239] U. Kang, C. E. Tsourakakis, A. P. Appel, C. Faloutsos, and J. Leskovec. Hadi: Mining radii of large graphs. *ACM Trans. Knowl. Discov. Data*, 5(2):8:1–8:24, February 2011. ISSN 1556-4681. doi: 10.1145/1921632.1921634. URL <http://doi.acm.org/10.1145/1921632.1921634>.
- [240] U. Kang, Charalampos E. Tsourakakis, and Christos Faloutsos. Pegasus: mining peta-scale graphs. *Knowl. Inf. Syst.*, 27(2):303–325, 2011.
- [241] R. Kannan and V. Vinay. Analyzing the structure of large graphs. *Manuscript*, 1999.
- [242] T. Karagiannis, C. Gkantsidis, D. Narayanan, and A. Rowstron. Hermes: Clustering users in large-scale e-mail services. In *ACM Symposium on Cloud Computing 2010*, 2010.
- [243] D. Karger, R. Motwani, and M. Sudan. Approximate graph coloring by semidefinite programming. *J. ACM*, 45(2):246–265, 1998.
- [244] G. Karypis and V. Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Sci. Comput.*, 20(1):359–392, December 1998. ISSN 1064-8275. doi: 10.1137/S1064827595287997. URL <http://dx.doi.org/10.1137/S1064827595287997>.
- [245] George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *JSC*, 20(1), 1998.
- [246] B. W. Kernighan and S. Lin. An efficient heuristic procedure for partitioning graphs. *Bell Syst. Tech. J.*, 49(2):291–307, February 1970.
- [247] Subhash Khot. Ruling out ptas for graph min-bisection, dense k -subgraph, and bipartite clique. *Journal of Computing*, 36(4), 2006.
- [248] Samir Khuller and Barna Saha. On finding dense subgraphs. In *ICALP*, 2009.
- [249] Jeong Han Kim and Van H Vu. Concentration of multivariate polynomials and its applications. *Combinatorica*, 20(3):417–434, 2000.
- [250] M. Kimura and Y. Okano. Identification and assignment of the human nima-related protein kinase 7 gene (*nek7*) to human chromosome 1q31.3. *Cytogenet. Cell Genet.*, 2001.

- [251] M. Klawe and D. Kleitman. An almost linear time algorithm for generalized matrix searching. *SIAM J. Discret. Math.*, 3:81–97, January 1990. ISSN 0895-4801. doi: 10.1137/0403009. URL <http://portal.acm.org/citation.cfm?id=84499.84509>.
- [252] Jon Kleinberg. The small-world phenomenon: an algorithm perspective. In *Proceedings of the thirty-second annual ACM symposium on Theory of computing*, pages 163–170. ACM, 2000.
- [253] Jon M Kleinberg, Ravi Kumar, Prabhakar Raghavan, Sridhar Rajagopalan, and Andrew S Tomkins. The web as a graph: Measurements, models, and methods. In *Computing and combinatorics*, pages 1–17. Springer, 1999.
- [254] Ton Kloks. *Treewidth: computations and approximations*, volume 842. Springer-Verlag New York Incorporated, 1994.
- [255] D. E. Knuth. Optimum binary search trees. *Acta Inf.*, 1:14–25, 1971.
- [256] Donald E. Knuth. *Seminumerical Algorithms*. Addison-Wesley, 1981.
- [257] Mihail N. Kolountzakis, Gary L. Miller, Richard Peng, and Charalampos E. Tsourakakis. Efficient triangle counting in large graphs via degree-based vertex partitioning. *Internet Mathematics*, 8(1-2), 2012.
- [258] A. Konstantin and H. Räcke. Balanced graph partitioning. In *SPAA '04*, pages 120–124, 2004. ISBN 1-58113-840-7.
- [259] Ioannis Koutis, Gary L Miller, and Richard Peng. Approaching optimality for solving sdd linear systems. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pages 235–244. IEEE, 2010.
- [260] R. Krauthgamer, J. (S.) Naor, and R. Schwartz. Partitioning graphs into balanced components. In *SODA '09*, pages 942–949, 2009.
- [261] Michael Krivelevich and Raphael Yuster. The rainbow connection of a graph is (at most) reciprocal to its minimum degree. *Journal of Graph Theory*, 63(3):185–191, 2010.
- [262] Ravi Kumar, Prabhakar Raghavan, Sridhar Rajagopalan, D Sivakumar, Andrew Tomkins, and Eli Upfal. Stochastic models for the web graph. In *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium on*, pages 57–65. IEEE, 2000.
- [263] Weil R. Lai, Mark D. Johnson, Raju Kucherlapati, and Peter J. Park. Comparative analysis of algorithms for identifying amplifications and deletions in array cgh data.

- Bioinformatics*, 21(19):3763–3770, 2005. ISSN 1367-4803. doi: <http://dx.doi.org/10.1093/bioinformatics/bti611>.
- [264] Anukool Lakhina, John W Byers, Mark Crovella, and Peng Xie. Sampling biases in ip topology measurements. In *INFOCOM 2003. Twenty-Second Annual Joint Conference of the IEEE Computer and Communications. IEEE Societies*, volume 1, pages 332–341. IEEE, 2003.
- [265] Stefan Lämmer, Björn Gehlsen, and Dirk Helbing. Scaling laws in the spatial structure of urban road networks. *Physica A: Statistical Mechanics and its Applications*, 363(1):89–95, 2006.
- [266] P. Lamy, C. L. Anderson, L. Dyrskjot, N. Topping, and C. Wiuf. A hidden Markov model to estimate population mixture and allelic copy-numbers in cancers using Affymetrix SNP arrays. *BMC Bioinformatics*, 8:434, 2007.
- [267] Michael A. Langston and et al. A combinatorial approach to the analysis of differential gene expression data: The use of graph algorithms for disease prediction and screening. In *Methods of Microarray Data Analysis IV*. 2005.
- [268] L. L. Larmore and B. Schieber. On-line dynamic programming with applications to the prediction of rna secondary structure. *J. Algorithms*, 12(3):490–515, 1991.
- [269] Matthieu Latapy. Main-memory triangle computations for very large (sparse (power-law)) graphs. *Theoretical Computer Science*, 407(1):458–473, 2008.
- [270] Silvio Lattanzi and D Sivakumar. Affiliation networks. In *Proceedings of the 41st annual ACM symposium on Theory of computing*, pages 427–434. ACM, 2009.
- [271] Marc Lavielle. Using penalized contrasts for the change-point problem. *Signal Processing*, 85(8):1501–1510, 2005.
- [272] D.D. Lee and H.S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788791, 1999.
- [273] Jong Woo Lee, Young Hwa Soung, Su Young Kim, Hae Woo Lee, Won Sang Park, Suk Woo Nam, Sang Ho Kim, Jung Young Lee, Nam Jin Yoo, , and Sug Hyung Lee. Pik3ca gene is frequently mutated in breast carcinomas and hepatocellular carcinomas. *Oncogene*, 2005.
- [274] Victor E. Lee, Ning Ruan, Ruoming Jin, and Charu C. Aggarwal. A survey of algorithms for dense subgraph discovery. In *Managing and Mining Graph Data*. 2010.

- [275] J. Lejeune, M. Gautier, and R. Turpin. Etude des chromosomes somatiques de neuf enfants mongoliens. *Comptes Rendus Hebdomadaires des Seances de l'Academie des Sciences (Paris)*, 248(11):17211722, 1959.
- [276] J. Lejeune, J. Lafourcade, R. Berger, J. Vialatta, M. Boeswillwald, P. Seringe, and R. Turpin. Trois ca de deletion partielle du bras court d'un chromosome 5. *Compte Rendus de l'Academie des Sciences (Paris)*, 257:3098, 1963.
- [277] J Leskovec, K. J. Lang, A. Dasgupta, and M. W. Mahoney. Statistical properties of community structure in large social and information networks. In *WWW '08*, 2008.
- [278] Jure Leskovec and Christos Faloutsos. Scalable modeling of real graphs using kronecker multiplication. In *Proceedings of the 24th international conference on Machine learning*, pages 497–504. ACM, 2007.
- [279] Jure Leskovec, Deepayan Chakrabarti, Jon M. Kleinberg, and Christos Faloutsos. Realistic, mathematically tractable graph generation and evolution, using kronecker multiplication. *PKDD*, pages 133–145, 2005.
- [280] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graph evolution: Densification and shrinking diameters. *ACM Trans. Knowl. Discov. Data*, 1(1):2, 2007. ISSN 1556-4681. doi: <http://doi.acm.org/10.1145/1217299.1217301>.
- [281] Jure Leskovec, Lars Backstrom, Ravi Kumar, and Andrew Tomkins. Microscopic evolution of social networks. In *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '08, pages 462–470, New York, NY, USA, 2008. ACM. ISBN 978-1-60558-193-4. doi: <http://doi.acm.org/10.1145/1401890.1401948>. URL <http://doi.acm.org/10.1145/1401890.1401948>.
- [282] Jure Leskovec, Deepayan Chakrabarti, Jon Kleinberg, Christos Faloutsos, and Zoubin Ghahramani. Kronecker graphs: An approach to modeling networks. *The Journal of Machine Learning Research*, 11:985–1042, 2010.
- [283] Jure Leskovec, Daniel Huttenlocher, and Jon Kleinberg. Predicting positive and negative links in online social networks. In *Proceedings of the 19th international conference on World wide web*, WWW '10, pages 641–650, New York, NY, USA, 2010. ACM. ISBN 978-1-60558-799-8. doi: <http://doi.acm.org/10.1145/1772690.1772756>. URL <http://doi.acm.org/10.1145/1772690.1772756>.
- [284] Jurij Leskovec, Deepayan Chakrabarti, Jon Kleinberg, and Christos Faloutsos. Realistic, mathematically tractable graph generation and evolution, using kronecker

- multiplication. In *Knowledge Discovery in Databases: PKDD 2005*, pages 133–145. Springer, 2005.
- [285] Xueliang Li and Yuefang Sun. *Rainbow connections of graphs*. Springer, 2012.
- [286] Doron Lipson, Yonatan Aumann, Amir Ben-Dor, Nathan Linial, and Zohar Yakhini. Efficient calculation of interval scores for dna copy number data analysis. In *RECOMB*, pages 83–100, 2005.
- [287] Richard J Lipton and Robert Endre Tarjan. A separator theorem for planar graphs. *SIAM Journal on Applied Mathematics*, 36(2):177–189, 1979.
- [288] Jun Liu, Jaaved Mohammed, James Carter, Sanjay Ranka, Tamer Kahveci, and Michael Baudis. Distance-based clustering of CGH data. *Bioinformatics*, 22(16):1971–1978, 2006.
- [289] Thomas Longerich, Michael Martin Mueller, Kai Breuhahn, Peter Schirmacher, Axel Benner, and Christiane Heiss. Oncogenetic tree modeling of human hepatocarcinogenesis. *International Journal of Cancer*, 130(3):575–583, 2012.
- [290] László Lovász and Miklós Simonovits. The mixing rate of markov chains, an isoperimetric inequality, and computing the volume. In *Foundations of Computer Science, 1990. Proceedings., 31st Annual Symposium on*, pages 346–354. IEEE, 1990.
- [291] Y. Low, J. Gonzalez, A. Kyrola, D. Bickson, C. Guestrin, and J. M. Hellerstein. Graphlab: A new framework for parallel machine learning. In *UAI*, pages 340–349, 2010.
- [292] Jun Ma and Shaohan Ma. Efficient parallel algorithms for some graph theory problems. *JCST*, 1993.
- [293] Avner Magen and Anastasios Zouzias. Near optimal dimensionality reductions that preserve volumes. In *Approximation, Randomization and Combinatorial Optimization. Algorithms and Techniques*, pages 523–534. Springer, 2008.
- [294] Mohammad Mahdian and Ying Xu. Stochastic kronecker graphs. In *Algorithms and models for the web-graph*, pages 179–186. Springer, 2007.
- [295] G. Malewicz, M. H. Austern, A.J.C Bik, J.C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: a system for large-scale graph processing. In *SIGMOD '10*, pages 135–146, 2010.
- [296] Benoit B Mandelbrot. *Fractals: form, chance, and dimension*. WH Freeman San Francisco, 1977.

- [297] Jean-Francois Marckert and Marie Albenque. Some families of increasing planar maps. *Electronic Journal of Probability*, 13:1624–1671, 2008.
- [298] Guglielmo Marconi. Wireless telegraphic communication. *Nobel Price Lecture*, 1909.
- [299] J. Matousek. Reporting points in halfspaces. In *FOCS*, pages 207–215, 1991.
- [300] Colin McDiarmid, Angelika Steger, and Dominic JA Welsh. Random planar graphs. *Journal of Combinatorial Theory, Series B*, 93(2):187–205, 2005.
- [301] Mary McGlohon, Leman Akoglu, and Christos Faloutsos. Weighted graphs and disconnected components: patterns and a generator. *KDD*, 2008.
- [302] Milena Mihail and Christos Papadimitriou. On the eigenvalue power law. In *Randomization and approximation techniques in computer science*, pages 254–262. Springer, 2002.
- [303] Stanley Milgram. The small world problem. *Psychology*, 2(1):60–67, 1967.
- [304] Gary L. Miller, Richard Peng, Russell Schwartz, and Charalampos E. Tsourakakis. Approximate dynamic programming using halfspace queries and multiscale monge decomposition. In *SODA*, pages 1675–1682, 2011.
- [305] Michael Mitzenmacher. A brief history of generative models for power law and lognormal distributions. *Internet mathematics*, 1(2):226–251, 2004.
- [306] Robert J. Mokken. Cliques, clubs and clans. *Quality & Quantity*, 13(2), 1979.
- [307] G. Monge. Memoire sue la theorie des deblais et de remblais. *Histoire de l'Academie Royale des Sciences de Paris*, pages 666–704, 1781.
- [308] Gordon E Moore et al. Cramming more components onto integrated circuits. *Proceedings of the IEEE*, 86(1):82–85, 1998.
- [309] Robert Morris. Counting large numbers of events in small registers. *Communications of the ACM*, 21(10):840–842, 1978.
- [310] C. L. Myers, M. J. Dunham, S. Y. Kung, and O. G. Troyanskaya. Accurate detection of aneuploidies in array cgh and gene expression microarray data. *Bioinformatics*, 20(18):3533–3543, December 2004. ISSN 1367-4803. URL <http://view.ncbi.nlm.nih.gov/pubmed/15284100>.
- [311] Kolountzakis M. N., G. L. Miller, R. Peng, and Tsourakakis C. E. Efficient triangle counting in large graphs via degree-based vertex partitioning. pages 15–24, 2010.

- [312] Nicholas Navin, Alexander Krasnitz, Linda Rodgers, Kerry Cook, Jennifer Meth, Jude Kendall, Michael Riggs, Yvonne Eberling, Jennifer Troge, Vladimir Grubor, Dan Levy, Par Lundin, Susanne Maner, Anders Zetterberg, James Hicks, and Michael Wigler. Inferring tumor progression from genomic heterogeneity. *Genome Research*, 20:68–80, March 2010.
- [313] R. Neve, K. Chin, J. Fridlyand, J. Yeh, F. Baehner, T. Fevr, N. Clark, N. Bayani, J. Coppe, and F. Tong. A collection of breast cancer cell lines for the study of functionally distinct cancer subtypes. *Cancer Cell*, 10:515–527, 2006.
- [314] M. E. J. Newman. Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23):8577–8582, 2006.
- [315] M. E. J. Newman. The structure and function of complex networks. *SIAM review*, 45(2):167–256, 2003.
- [316] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical review E*, 69(2):026113, 2004.
- [317] M. E. J. Newman, D. J. Watts, and S. H. Strogatz. Random graph models of social networks. *Proceedings of the National Academy of Sciences of the United States of America*, 99(Suppl 1):2566–2572, February 2002. ISSN 0027-8424. doi: 10.1073/pnas.012582999. URL <http://dx.doi.org/10.1073/pnas.012582999>.
- [318] Andrew Y. Ng, Michael I. Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *NIPS*, 2001.
- [319] G. Ng, J. Huang, I. Roberts, and N. Coleman. Defining ploidy-specific thresholds in array comparative genomic hybridization to improve the sensitivity of detection of single copy alterations in cell lines. *Journal of Molecular Diagnostics*, 2006.
- [320] A. B. Olshen, E. S. Venkatraman, R. Lucito, and M. Wigler. Circular binary segmentation for the analysis of array-based dna copy number data. *Biostatistics*, 5(4):557–572, October 2004. ISSN 1465-4644. URL <http://view.ncbi.nlm.nih.gov/pubmed/15475419>.
- [321] A. B. Olshen, E.S. Venkatraman, R. Lucito, and M. Wigler. Circular binary segmentation for the analysis of array-based dna copy number data. *Biostatistics*, 5(4):557–572, October 2004. URL <http://view.ncbi.nlm.nih.gov/pubmed/15475419>.
- [322] Lorenzo Orecchia. *Fast Approximation Algorithms for Graph Partitioning using Spectral and Semidefinite-Programming Techniques*. PhD thesis, EECS Department, University of California, Berkeley, May 2011. URL <http://www.eecs.berkeley.edu/Pubs/TechRpts/2011/EECS-2011-56.html>.

- [323] B Orsetti, M Nugoli, N Cervera, L Lasorsa, P Chuchana, L Ursule, C Nguyen, R Redon, S du Manoir, C Rodriguez, and Charles Theillet. Genomic and expression profiling of chromosome 17 in breast cancer reveals complex patterns of alterations and novel candidate genes. *Cancer Research*, 2004.
- [324] A. Packer. NP-hardness of largest contained and smallest containing simplices for v - and h -polytopes. *Discrete Computational Geometry*, 28:349–377, 2002.
- [325] Rasmus Pagh and Charalampos E. Tsourakakis. Colorful triangle counting and a mapreduce implementation. *IPL*, 112(7), 2012.
- [326] S. Paik, C. Kim, and N. Wolmark. Her2 status and benefit from adjuvant trastuzumab in breast cancer. *N Engl J Med*, 358:1409–1411, 2008.
- [327] Christopher R. Palmer, Phillip B. Gibbons, and Christos Faloutsos. Anf: a fast and scalable tool for data mining in massive graphs. In *KDD*, pages 81–90, 2002.
- [328] Jia-Yu Pan, Hyung-Jeong Yang, Christos Faloutsos, and Pinar Duygulu. Automatic multimedia cross-modal correlation discovery. *ACM SIGKDD*, August 2004.
- [329] Alois Panholzer and Georg Seitz. Ordered increasing k-trees: Introduction and analysis of a preferential attachment network model. *arXiv preprint arXiv:1003.0320*, 2010.
- [330] C. Papadimitriou. Personal communication.
- [331] Christos H. Papadimitriou and Mihalis Yannakakis. Clique problem for planar graphs. *INFO. PROC. LETT.*, 13(4):131–133, 1981.
- [332] Spiros Papadimitriou and Jimeng Sun. Disco: Distributed co-clustering with map-reduce. *ICDM*, 2008.
- [333] Jeffrey Pattillo, Alexander Veremyev, Sergiy Butenko, and Vladimir Boginski. On the maximum quasi-clique problem. *Discr. Ap. Math.*, 161(1–2), 2012.
- [334] Andrew Pavlo, Erik Paulson, Alexander Rasin, Daniel J. Abadi, David J. Dewitt, Samuel Madden, and Michael Stonebraker. A comparison of approaches to large-scale data analysis. *SIGMOD*, June 2009. URL <http://database.cs.brown.edu/sigmod09/benchmarks-sigmod09.pdf>.
- [335] K. Pearson. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine*, 2:559–572, 1901.
- [336] M. D. Pegram, G. Konecny, and D. J. Slamon. The molecular and cellular biology of her2/neu gene amplification/overexpression and the clinical development of

- herceptin (trastuzumab) therapy for breast cancer. *Cancer Treat Res*, 103:57–75, 2000.
- [337] G. Pennington, C. E. Smith, S. Shackney, and R. Schwartz. Reconstructing tumor phylogenies from single-cell data. *Journal of Bioinformatics and Computational Biology*, 5:407–427, 2007.
- [338] C. M. Perou, T. Sorlie, M. B. Eisen, M. van der Rijn, S. S. Jeffrey and C. A. Rees, J. R. Pollack, D. T. Ross, H. Johnsen, L. A. Akslen, O. Fluge, A. Pergamenschikov, C. Williams, S. X. Zhu, P. E. Lønning, A. L. Børresen-Dale, P. O. Brown, and D. Botstein. Molecular portraits of human breast tumors. *Nature*, 406:747–752, 2000.
- [339] F. Picard, S. Robin, M. Lavielle, C. Vaisse, and J. J. Daudin. A statistical approach for array CGH data analysis. 6, 2005.
- [340] F. Picard, S. Robin, E. Lebarbier, and J.J. Daudin. A segmentation/clustering model for the analysis of array cgh data. *Biometrics*, 63, 2007.
- [341] D Pinkel, R Segreaves, D Sudar, S Clark, I Poole, D Kowbel, C Collins, W-L Kuo, C Chen, and Y Zha. High resolution analysis of DNA copy number variation using comparative genomic hybridization to microarrays. *Nature Genetics*, 25:207 – 211, 1998.
- [342] J. R. Pollack, C. M. Perou, A. A. Alizadeh, M. B. Eisen, A. Pergamenschikov, C. F. Williams, S. S. Jeffrey, D. Botstein, and P. O. Brown. Genome-wide analysis of dna copy-number changes using cdna microarrays. *Nature Genetics*, 23:41–46, 1999.
- [343] Pascal Pons and Matthieu Latapy. Computing communities in large networks using random walks. In *Computer and Information Sciences-ISCIS 2005*, pages 284–293. Springer, 2005.
- [344] V. Prabhakaran and et al. Managing large graphs on multi-cores with graph awareness. In *USENIX ATC'12*, 2012.
- [345] Pawel Pralat and Nicholas Wormald. Growing protean graphs. *Internet Mathematics*, 4(1):1–16, 2007.
- [346] J. Pujol and et al. The little engine(s) that could: Scaling online social networks. In *ACM SIGCOMM 2010*, 2010.
- [347] Garry Robins, Pip Pattison, Yuval Kalish, and Dean Lusher. An introduction to exponential random graph (p^*) models for social networks. *Social Networks*, 29 (2):173–191, May 2007.

- [348] Jacques Rougemont and Pascal Hingamp. Dna microarray data and contextual analysis of correlation graphs. *BMC Bioinformatics*, 4:15, 2003.
- [349] V. Satuluri, S. Parthasarathy, and Y. Ruan. Local graph sparsification for scalable clustering. In *SIGMOD '11*, pages 721–732, 2011. ISBN 978-1-4503-0661-4.
- [350] A.A. Schäffer. Personal communication.
- [351] T. Schank and D. Wagner. Approximating clustering-coefficient and transitivity. *Journal of Graph Algorithms and Applications*, 9(2):265–275, 2004. URL <http://www.ubka.uni-karlsruhe.de/vvv/ira/2004/9/9.pdf>.
- [352] Thomas Schank and Dorothea Wagner. Finding, counting and listing all triangles in large graphs, an experimental study. In *Experimental and Efficient Algorithms*, pages 606–609. Springer, 2005.
- [353] K. Schloegel, G. Karypis, and V. Kumar. Parallel multilevel algorithms for multi-constraint graph partitioning (distinguished paper). In *Proceedings from the 6th International Euro-Par Conference on Parallel Processing*, Euro-Par '00, pages 296–310, 2000.
- [354] K. Schloegel, G. Karypis, and V. Kumar. Parallel static and dynamic multi-constraint graph partitioning. *Concurrency and Computation: Practice and Experience*, 14(3):219 – 240, 2002.
- [355] B. Schölkopf and A. J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- [356] B. Schölkopf, A. J. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10:1299–1319, 1998.
- [357] Alexander Schrijver. *Combinatorial Optimization : Polyhedra and Efficiency (Algorithms and Combinatorics)*. Springer, 2004.
- [358] R. Schwartz and S. Shackney. Applying unmixing to gene expression data for tumor phylogeny inference. *BMC Bioinformatics*, 11:42, 2010.
- [359] Thomas Seidl, Brigitte Boden, and Sergej Fries. Cc-mr-finding connected components in huge graphs with mapreduce. *Machine Learning and Knowledge Discovery in Databases*, pages 458–473, 2012.
- [360] A. Sen and M. Srivastava. On tests for detecting change in mean. *Annals of Statistics*, 3:98–108, 1975.

- [361] S. E. Shackney, C. A. Smith, A. Pollice, K. Brown, R. Day, T. Julian, and J. F. Silverman. Intracellular patterns of Her-2/neu, ras, and ploidy abnormalities in primary human breast cancers predict postoperative clinical disease-free survival. *Clin Cancer Res*, 10:3042–3052, 2004.
- [362] S. P. Shah, X. Xuan, R. J. Deleeuw, M. Khojasteh, W. L. Lam, R. Ng, and K. P. Murphy. Integrating copy number polymorphisms into array cgh analysis using a robust hmm. *Bioinformatics*, 22(14), July 2006. ISSN 1460-2059. URL <http://view.ncbi.nlm.nih.gov/pubmed/16873504>.
- [363] Yanxin Shi, Fan Guo, Wei Wu, and Eric P. Xing. Gimsan: A new statistical method for analyzing whole-genome array cgh data. In *RECOMB*, pages 151–165, 2007.
- [364] B P Sinha, B B Bhattacharya, S Ghose, and P K Srimani. A parallel algorithm to compute the shortest paths and diameter of a graph and its vlsi implementation. *IEEE Trans. Comput.*, 1986.
- [365] AM Snijders, N Nowak, R Segraves, S Blackwood, N Brown, J Conroy, G Hamilton, AK Hindle, B Huey, K Kimura, S Law, K Myambo, J Palmer, B Ylstra, JP Yue, JW Gray, AN Jain, D Pinkel, and DG Albertson. Assembly of microarrays for genome-wide measurement of DNA copy number. *Nature Genetics*, 29: 263–264, 2001.
- [366] Therese Sorlie and et al. Repeated observation of breast tumor subtypes in independent gene expression data sets. *PNAS*, 100(14), 2003.
- [367] C. Sotiriou, S. Y. Neo, L. M. McShane, E. L. Korn, P. M. Long, A. Jazaeri, P. Martiat, S. B. Fox, A. L. Harris, and E. T. Liu. Breast cancer classification and prognosis based on gene expression profiles from a population-based study. *Proc Natl Acad Sci USA*, 100:10393–10398, 2003.
- [368] Mauro Sozio and Aristides Gionis. The community-search problem and how to plan a successful cocktail party. In *KDD*, 2010.
- [369] Daniel Spielman, 2010. Graphs and Networks, Yale www.cs.yale.edu/homes/spielman/462/2010.
- [370] Daniel A. Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. In *STOC*, pages 563–568, 2008.
- [371] Daniel A Spielman and Shang-Hua Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *Proceedings of the*

- thirty-sixth annual ACM symposium on Theory of computing*, pages 81–90. ACM, 2004.
- [372] Daniel A Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs and its application to nearly-linear time graph partitioning. *arXiv preprint arXiv:0809.3232*, 2008.
- [373] I. Stanton and G. Kliot. Streaming graph partitioning for large distributed graphs. In *ACM KDD*, pages 1222–1230, 2012.
- [374] Isabelle Stanton. Streaming balanced graph partitioning algorithms for random graphs. In *Arxiv*, available at <http://arxiv.org/abs/1212.1121>, 2012.
- [375] Gilbert Strang. *Introduction to linear algebra*. Wellesley Cambridge Pr, 2003.
- [376] Michael R Stratton, Peter J Campbell, and P Andrew Futreal. The cancer genome. *Nature*, 458(7239):719–724, 2009.
- [377] Siddharth Suri and Sergei Vassilvitskii. Counting triangles and the curse of the last reducer. In *Proceedings of the 20th international conference on World wide web*, pages 607–614. ACM, 2011.
- [378] Terence Tao and Van H Vu. *Additive combinatorics*, volume 105. Cambridge University Press, 2006.
- [379] Robert Tibshirani and Pei Wang. Spatial smoothing and hot spot detection for CGH data using the fused lasso. *Biostatistics (Oxford, England)*, 9(1):18–29, January 2008. URL <http://dx.doi.org/10.1093/biostatistics/kxm013>.
- [380] Ali Tofigh. *Using trees to capture reticulate evolution: lateral gene transfers and cancer progression*. PhD thesis, KTH, 2009.
- [381] K.C. Toh, M.J Todd, and R.H. Tutuncu. Sdpt3 – a matlab software package for semidefinite programming. *Optimization Methods and Software*, 11:545–581, 1999.
- [382] David Tolliver, Charalampos E. Tsourakakis, Ayshwarya Subramanian, Stanley Shackney, and Russell Schwartz. Robust unmixing of tumor states in array comparative genomic hybridization data. *Bioinformatics [ISMB]*, 26(12):106–114, 2010.
- [383] C. E. Tsourakakis. Towards quantifying vertex similarity in networks. *Internet Mathematics*, 2014.
- [384] C. E. Tsourakakis, C. Gkantsidis, B. Radunovic, and M. Vojnovic. Fennel: Streaming graph partitioning for massive distributed graphs. November 2012. Microsoft Technical Report MSR-TR-2012-213.

- [385] C.E. Tsourakakis, F. Bonchi, A. Gionis, F. Gullo, and M. Tsiarli. Denser than the densest subgraph: Extracting optimal quasi-cliques with quality guarantees. *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 104–112, 2013.
- [386] Charalampos E. Tsourakakis. Perfect reconstruction of oncogenetic trees. *CoRR*, abs/1201.4618, 2012.
- [387] Charalampos E. Tsourakakis. Fast counting of triangles in large real networks without counting: Algorithms and laws. In *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining, ICDM '08*, pages 608–617, Washington, DC, USA, 2008. IEEE Computer Society. ISBN 978-0-7695-3502-9. doi: 10.1109/ICDM.2008.72. URL <http://dx.doi.org/10.1109/ICDM.2008.72>.
- [388] Charalampos E Tsourakakis. Counting triangles in real-world networks using projections. *Knowledge and Information Systems*, 26(3):501–520, 2011.
- [389] Charalampos E Tsourakakis. Mach: Fast randomized tensor decompositions. *SIAM Data Mining (SDM)*, 2010.
- [390] Charalampos E. Tsourakakis. Fast counting of triangles in large real networks without counting: Algorithms and laws. In *ICDM*, 2008.
- [391] Charalampos E. Tsourakakis. Towards quantifying vertex similarity in networks. *Internet Mathematics (accepted)*, 2014.
- [392] Charalampos E Tsourakakis. Modeling intratumor gene copy number heterogeneity using fluorescence in situ hybridization data. 2013. submitted.
- [393] Charalampos E Tsourakakis, Petros Drineas, Eirinaios Michelakis, Ioannis Koutis, and Christos Faloutsos. Spectral counting of triangles in power-law networks via element-wise sparsification. In *Social Network Analysis and Mining, 2009. ASONAM'09. International Conference on Advances in*, pages 66–71. IEEE, 2009.
- [394] Charalampos E Tsourakakis, U Kang, Gary L Miller, and Christos Faloutsos. Doulion: counting triangles in massive graphs with a coin. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 837–846. ACM, 2009.
- [395] Charalampos E Tsourakakis, Petros Drineas, Eirinaios Michelakis, Ioannis Koutis, and Christos Faloutsos. Spectral counting of triangles via element-wise sparsification and triangle-based link recommendation. *Social Network Analysis and Mining*, 1(2):75–81, 2011.

- [396] Charalampos E. Tsourakakis, Mihail N. Kolountzakis, and Gary L. Miller. Triangle sparsifiers. *JGAA*, 15(6), 2011.
- [397] Charalampos E. Tsourakakis, Richard Peng, Maria A. Tsiarli, Gary L. Miller, and Russell Schwartz. Approximation algorithms for speeding up dynamic programming and denoising acgh data. *ACM Journal of Experimental Algorithmics*, 16, 2011.
- [398] Johan Ugander, Lars Backstrom, and Jon Kleinberg. Subgraph frequencies: Mapping the empirical and extremal geography of large graph collections. *WWW*, 2013.
- [399] Takeaki Uno. An efficient algorithm for solving pseudo clique enumeration problem. *Algorithmica*, 56(1), 2010.
- [400] Leslie G Valiant. A bridging model for parallel computation. *Communications of the ACM*, 33(8):103–111, 1990.
- [401] Leslie G Valiant. Memorization and association on a realistic neural model. *Neural computation*, 17(3):527–555, 2005.
- [402] Marc J. van de Vijver and et al. A gene-expression signature as a predictor of survival in breast cancer. *The New England journal of medicine*, 347(25), 2002.
- [403] Remco Van Der Hofstad. Random graphs and complex networks. Available online at <http://www.win.tue.nl/~rhofstad/>, 2009.
- [404] F. Viti, E. Mosca, I. Merelli, A. Calabria, R. Alfieri, and L. Milanese. Ontological enrichment of the Genes-to-Systems Breast Cancer Database. *Communications in Computer and Information Science*, 46:178–182, 2009.
- [405] Bert Vogelstein, Eric R Fearon, Stanley R Hamilton, Scott E Kern, Ann C Preisinger, Mark Leppert, Y Nakamura, R White, AM Smits, JL Bos, et al. Genetic alterations during colorectal-tumor development. *The New England journal of medicine*, 319(9):525, 1988.
- [406] Anja Von Heydebreck, Bastian Gunawan, and László Füzesi. Maximum likelihood estimation of oncogenetic tree models. *Biostatistics*, 5(4):545–556, 2004.
- [407] Van H Vu. On the concentration of multivariate polynomials with small expectation. *Random Structures and Algorithms*, 16(4):344–363, 2000.
- [408] P. Wang, Y. Kim, J. Pollack, B Narasimhan, and R. Tibshirani. A method for calling gains and losses in array cgh data. *Biostatistics*, 6(1):45–58, January 2005. ISSN 1465-4644. URL <http://view.ncbi.nlm.nih.gov/pubmed/15618527>.

- [409] S. Wasserman and K. Faust. *Social Network Analysis: Methods and Applications*. Cambridge University Press, 1994.
- [410] Stanley Wasserman and Philippa Pattison. Logit models and logistic regressions for social networks: I. an introduction to markov graphs andp. *Psychometrika*, 61(3):401–425, September 1996. doi: 10.1007/BF02294547. URL <http://dx.doi.org/10.1007/BF02294547>.
- [411] M. S. Waterman and T.F. Smith. Rapid dynamic programming algorithms for rna secondary structure. *Adv. Appl. Math.*, 7:455–464, December 1986. ISSN 0196-8858. doi: 10.1016/0196-8858(86)90025-4. URL <http://portal.acm.org/citation.cfm?id=9300.9303>.
- [412] Duncan J. Watts and Steven H. Strogatz. Collective dynamics of “small-world” networks. *Nature*, 393(6684):440–442, 1998. URL <http://dx.doi.org/10.1038/30918>.
- [413] Robert A. Weinberg. *The Biology of Cancer HB*. Garland Science, 2006.
- [414] Howard T. Welser, Eric Gleave, Danyel Fisher, and Marc Smith. Visualizing the Signatures of Social Roles in Online Discussion Groups. *The Journal of Social Structure*, 8(2), 2007. URL <http://www.cmu.edu/joss/content/articles/volume8/Welser/>.
- [415] Hanni Willenbrock and Jane Fridlyand. A comparison study: applying segmentation to array cgh data for downstream analyses. *Bioinformatics*, 21(22):4084–4091, November 2005. ISSN 1367-4803. doi: 10.1093/bioinformatics/bti677. URL <http://dx.doi.org/10.1093/bioinformatics/bti677>.
- [416] V Vassilevska Williams. Breaking the coppersmith-winograd barrier. *Unpublished manuscript*, Nov, 2011.
- [417] Andreas Wimmer and Kevin Lewis. Beyond and Below Racial Homophily: ERG Models of a Friendship Network Documented on Facebook. *American Journal of Sociology*, 116(2):583–642, September 2010. doi: 10.1086/653658. URL <http://dx.doi.org/10.1086/653658>.
- [418] David Woodruff. Optimal space lower bounds for all frequency moments. In *Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms, SODA '04*, pages 167–175, Philadelphia, PA, USA, 2004. Society for Industrial and Applied Mathematics. ISBN 0-89871-558-X. URL <http://dl.acm.org/citation.cfm?id=982792.982817>.

- [419] Nicholas C Wormald. Models of random regular graphs. *London Mathematical Society Lecture Note Series*, pages 239–298, 1999.
- [420] Zhi-Xi Wu, Xin-Jian Xu, and Ying-Hai Wang. Comment on “maximal planar networks with large clustering coefficient and power-law degree distribution”. *Physical Review E*, 73(5):58101, 2006.
- [421] B. Xing, C. Greenwood, and S. Bull. A hierarchical clustering method for estimating copy number variation. *Biostatistics*, 8:632–653, 2007.
- [422] F. F. Yao. Efficient dynamic programming using quadrangle inequalities. In *Proceedings of the twelfth annual ACM symposium on Theory of computing*, STOC ’80, pages 429–435, New York, NY, USA, 1980. ACM. ISBN 0-89791-017-6. doi: <http://doi.acm.org/10.1145/800141.804691>. URL <http://doi.acm.org/10.1145/800141.804691>.
- [423] F. F. Yao. Speed-up in dynamic programming. *SIAM Journal on Algebraic and Discrete Methods*, 3(4):532–540, 1982. doi: 10.1137/0603055. URL <http://link.aip.org/link/?SML/3/532/1>.
- [424] F. Zhang, W. Gu, M. E. Hurles, and J. R. Lupski. Copy number variation in human health, disease, and evolution. *Annual Review of Genomics and Human Genetics*, 10:451–481, 2009.
- [425] Zhongzhi Zhang, Francesc Comellas, Guillaume Fertin, and Lili Rong. High-dimensional apollonian networks. *Journal of Physics A: Mathematical and General*, 39(8):1811, 2006.
- [426] Tao Zhou, Gang Yan, and Bing-Hong Wang. Maximal planar networks with large clustering coefficient and power-law degree distribution. *Physical Review E*, 71(4):046141, 2005.
- [427] Y. Zhou and S. Suri. Algorithms for minimum volume enclosing simplex in R^3 . *Proceedings of the Eleventh Annual ACM/SIAM Symposium on Discrete Algorithms*, pages 500–509, 2000.
- [428] Zeyuan Allen Zhu, Silvio Lattanzi, and Vahab Mirrokni. A local algorithm for finding well-connected clusters. *arXiv preprint arXiv:1304.8132*, 2013.