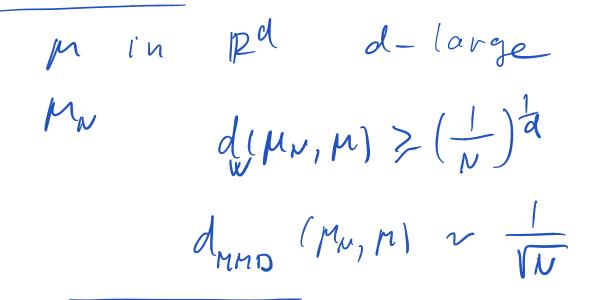
Stein Variational descent

Two objectives

- o deterministic, particle-based approximation for Fokker-Planck and related eq.
- methods that are robust wrt dimension.

(connection to MMD - max. mean discreparcy)



F-P

 $\partial_t g = 3g - dN(gPU)$

Sf dp-d) $d(\mu, \nu) = svp$ f & Lip(1)

sup FERN Sf dp-dy

Reproducing Kernel Hilbert Spaces (RKHS)
Let H be a Hilbert space of functions from

$$S \subseteq \mathbb{R}^{d}$$
 to \mathbb{R} .
Def H is an RKHS if fore SL the
evaluation at x is a continuous function
on H, that is $\sigma_{x} : H \rightarrow \mathbb{R}$ defined by
 $f \downarrow \xrightarrow{\sigma_{x}} f(x)$
is a continuous functional.
Consequence: By Piesz Representation Theorem
there exists $f_{x} \in H$ such that
 $\forall f \in H$ $f(x) = \delta_{x}(f) = \langle f, \phi_{x} \rangle_{H}$
The mepping from $StoH$, $X \rightarrow \phi_{x}$ is called
the feature map.
Let $K(x,y) = \langle \phi_{x}, \phi_{y} \rangle = \phi_{x}(r) = \phi_{y}(x)$
 K is called the Kernel. Inicon
 K is positive definite. For presigned measure
 $\int K(x,y) d\mu(x) d\mu(r) = \int \langle \phi_{x}, \phi_{y} \rangle = 0$
 $= \int \langle f_{x} d\mu(x), \phi_{y} \rangle d\mu(r) \rangle \geq 0$

A bit more precisely. For 6 probability measure

$$(m = a6, -b6_2)$$
 consider
 $f \mapsto S f(x) df(x)$
we note $|Sf(x) df(x)| = |S < \phi_x, f > df(x)|$
 $\leq S ||\phi_x||_H ||f||_H df(x)$
 $\leq S ||\phi_x||_H df(x) \cdot ||f||_H$
 $\leq V S ||f|_H df(x) \cdot ||f||_H$

From $(\not A)$ follows that if $g(x) = \int k(x, y) dy(y) = E_{pl}(d_{x})$ then $\|g\|_{H}^{2} = \int K(x, y) dy(x) dy(y)$

Theorem (Moore - Aronszajn) If $K : \mathcal{D} \times \mathcal{S} \longrightarrow \mathcal{R}$ is symmetric, positive definite then $\mathcal{I}_{o} = span \{ K(\cdot, y) : y \in \mathcal{S} \}$ is an inner product space with

 $< f_i g \rangle_{\mathcal{X}_0} = \overline{\mathcal{Z}} \overline{\mathcal{Z}} d_i (\beta_i k(x_i, y_i))$ where $f = \overline{z} d_i k(x_i, \cdot)$, $g = \overline{z} \beta_i k(\cdot, y_i)$. let Il be a completion of Ho. They His an PKHS. H is the set of pointwise limits of fu in Ro. [Notes by Sejdinović and Gretton]. [Bob mentioned : Paulsen & Raghupethi "An Intro to the Theory of RKHS" available online via Chu Library.] Examples of RKHS () K(x, 7) = y(x-7) y - Gaussian Il = { 2 * f : f - measure, 1 f 1 < 00 } < n*f, n*g> = SS n(x-7) fixigiyi dxdy ② M - compact manifold of dimension d. H^s - fractional Sobolev space. If s> d H³ C (M) and H' is PKHS

Spectral representation: (λ_i, Ψ_i) Laplacian eigenval. and eigenfunctions, orthonormal in $L^2(M)$. $f \in H^S$ $f = \sum_{i=1}^{\infty} a_i \Psi_i$

$$\begin{split} \|f\|_{H^{s}}^{z} &= \sum_{i=1}^{s} (1+\lambda_{i})^{s} a_{i}^{2} \qquad (1-\Delta)^{s} \\ \text{let us compute the Kernel} \\ (\varphi_{\overline{x}}, f) &= f(x) \\ \varphi_{\overline{x}} &= \overline{z} b_{i} \ \Psi_{i} \qquad f = \overline{z} a_{j} \ \Psi_{i} \\ <\varphi_{\overline{x}}, f) &= \sum_{i} (1+\lambda_{i})^{s} a_{i} b_{i} = \sum_{i} a_{i} \ \Psi_{i}(\overline{x}) \\ \text{thus } b_{i} &= \frac{1}{(1+\lambda_{i})^{s}} \ \Psi_{i}(\overline{x}) \\ So \qquad \overline{\Phi_{\overline{x}}}(x) &= \sum_{i} \frac{1}{(1+\lambda_{i})^{s}} \ \Psi_{i}(\overline{x}) \ \Psi_{i}(x) \\ [\sim Mercer theorem] \end{split}$$

Stein Variational Descent for Relative Entropy (Liv, Wang Neur IPS '16, Liv Neur IPS '17, Lu, Lu, Noley 18.] $E(g) = \int_{\mathbb{P}^d} g \ln\left(\frac{g}{F}\right) dx = \int g \ln g + g U dx$ $F F \sim e^{-U}$ F is a.c. wrt lobesgue measure.• Wasserstein gradient flow dE(r) g'10)=V - On a manifold - grad E is a minimizer of $R(v) = \frac{1}{2}g(y,v) + diff E[v] = \frac{1}{2}g(y,v) + diff E[v]$. For Wass: Tangent vector 2yg = -div(gv) $\mathcal{R}(v) = \frac{1}{2} \int |v|^2 g \, dx + \int V \cdot (\nabla g - g \nabla U) \, dx$ $\frac{dE}{dt} = \int S_t \ln g + S_t + S_t \ln F dt$ $= \int \oint V \frac{Pg}{P} + g V \frac{T}{P} dx = \int V (\nabla g - g PU) dr$ Minimizing R(V) gives $gV = -\nabla g + g \nabla U$ so $\partial_+g = -dN(gV) = \Delta g - dN(gVV)$ (FP)

· Stern Variational descent Consider a different metric for velocity $P(v) = \frac{1}{2} \|v\|_{H}^{2} + \int v \cdot (\nabla g - g R V) dx$ minimizing gives the $< \vee, w >_{H} = - \int (\nabla g - g \tau U) - w dt$ By above intro we have that V is the mean embedding for measure $(g \nabla U - \nabla g).$ So V(x) = S K(x, y) (S PU-RS) dy (f K(x,y) = K(x-y) then $V_s = K * (g P U - \nabla g) = K * (g P U) - D(K * g)$ Gradient flow $(SVD) \quad \exists_{tg} = dN(K * (\nabla p - g\nabla V))$ So by penalizing the velocity in higher norm we get a smoother velocity. Particle methods Consider $S_N = \frac{1}{N} \sum_{r=1}^{N} \delta_{x_i}$ Vwass, and relative entropy do not make sense. But V_{stein} does $SVD_{N} \begin{cases} V_{\tilde{e}} = \frac{1}{N} \sum_{\tilde{s}} (k(X_{\tilde{e}} - X_{\tilde{s}}) \nabla U(X_{\tilde{s}}) - \nabla k(X_{\tilde{e}} - X_{\tilde{s}})) \\ \frac{dX_{\tilde{e}}}{dt} = V_{\tilde{e}} \end{cases}$ Lu, Lu, Nolen study well posedness and asymptotics for SVD. (if K = J * J)

Solutions of SVDN are weak sol of SVD

In T2.7 Hey show stability for SVD using a coupling technique. This implies discrete to continuum convergence. Alternative functional (carrillo, Craig, Patacchini) $E_{k}(g) = \int \ln\left(\frac{k * g}{g}\right) dg = \int \ln(k * g) dg + \int U dg$ $\partial_{t}g = -div(gv)$ $diff E_{k}(v) = \frac{dE}{dt}$ $P(v) = \frac{1}{2} \int |v|^{2} dg + diff E_{k}(v)$ minimizer $V = -\nabla U - \nabla K * \left(\frac{g}{k * g}\right) - \frac{1}{k * g} \nabla K * g$ [CCP] show that as $N \to \infty$ particle

approximations converge, using Sandrer-Serfaty framework.