Kernelized Stelh Discrepancy. Stehn Variational Gradient Descent. (Oct 27th) For E(p) = Jpln(f)dx. $\frac{dE}{dt}(p) = \int V \cdot (\nabla p + p \nabla U) dx$ $= - \| v \|_{\mathcal{H}}^{2}$ (*) $= - D(\rho \| \sigma)^{2}$ Here V(x) = [VInv(y)k(x,y) + Vyk(x,y)dp(y) $D(p||\sigma) = \max_{\substack{\varphi \in \mathcal{H}, \|\varphi\|_{\mathcal{H}} \leq I}} E_{p}(\mathcal{A}_{\sigma}(\varphi))$

 $\mathcal{A}_{\varphi}(\phi) := \nabla \ln \sigma \cdot \phi + \nabla \cdot \phi$ "K" Steh Discrepancy LAKSD for Goodness-of-Att Tests (Liu, Lee, Jordan, ICML16') O Coodness- of- At Test" Measuring how well do the observed data correspond to the fitted model. $H_0: \mu = \gamma$ $\nu \rightarrow \nu \eta$ $H_1: M \neq V$ v known. Xi ... Xn ~ M

("Two-sample test": [Xi - Xn - M [Yi - Yn - M] Hi: M≠: $\Rightarrow H_{1}: M = \mathcal{V}$ $H_{1}: M \neq \mathcal{V} .$ - [[kelihood-based approaches (Need to compute libelihoods & CDF) hidden variable models & generative models X L MCMC; variational methods error large & hard to estimate (Curse of dhuendbacky) — Traditional methods (X², -Smithor etc) Today: Evaluating KSD. A "likelohood-free" approach

with guaranteed statistick significance. (i). Preliminarles, defention of KSD (ii) Propertiles of KSD (jii). Orighel SVGD approach. (iv). Open problems. (i) » . Stell's method. (1970s) for distributional approximation. "Probability methods", "MMD," $d_{\mathcal{H}}(\mu, \nu) = \sup_{h \in \mathcal{H}} |\int h d\mu - \int h d\nu|$ Test function family H: rich enough?

1_{8. ≤ x3} / Lip / 1_{8. ∈ Borel A3}... dK dW dTV Q. Given [Xi]_{i=1} independent variables, How to bound du (1 IXi, N(0,1))? (CLT: împortant cug rate.) Mah idea: Replace the charateristic <u>function</u> typically used to show cug in distribution with a "characteristy operator". Dethe A: C'(R) ~ CCR) by $A f = f' - \pi f$ Let I be the cdf of N(0,1), then $\exists ! f_x \text{ solves } f' - x f = \mathbf{1}_{\{\cdot \leq x\}} - \overline{\Phi}(x)$

for V fixed XGR. (ODE result)
⇒[Steh lemma]
$A r.v. W ~ N(0, 1) < \Rightarrow$
$\forall a.c. f with E[f(w)] < \infty, EAf(w) = 0$
⇒ (Cor1) For Vr.v. W,
$ P(w \le x) - \overline{E}(x) = E(f'_{x}(w) - wf_{x}(w)) .$
(ok) <u>Ceneral Setup</u> :
dze(u.v) = sup fhdu - fhdu hede
For ghen LEH, let fr solves
$f_{\mu}(w) - w f_{\mu}(w) = h(w) - \overline{\Phi}(h)$, where
$\overline{\Phi}(h) = E(h(z))$ for $Z \cap N(o, 1)$.
Then dae(ma) = sup E ficus - wficus)

for Wndu. E Stein's Normal Approximation] Once ne solve and restrict our discuss on $F = \{f \mid \|_{\infty}, \|f\|_{\infty} \leq 2, \|f'\|_{\infty} \leq \overline{E}\},\$ he have $d_{W}(W,Z) \leq \sup_{f \in F} |Ef(w) - wf(w)|$ (Cor 2) For O-mean independent {Xi}i=1, W = IXit, ne have $d_{W}(W, z) \leq \frac{1}{N^{3/2}} \Sigma E[X_{i}]^{3} + \int_{\pi n^{2}} \Sigma E[X_{i}]^{4}$ Ceneral Sotup. (N. Ross) What it we are approximatily a Probability distribution with smooth density 2? Replace -x by 2/2.

On $X \in \mathbb{R}^d$, 2 smooth densities p, 2 are "identical" $E_p(A_q(f)) = E_p(S_{q(x)} + \nabla_x f(x))$ Here Sq = Valagia is the "Stelh score function". linear operator Aq is called the "Stell's operator", acts on the "Stell class of q", i.e. $f \in C'(\mathcal{X})$ $\int_{\mathcal{X}}^{\prime} \nabla_{x} (f(x) p(x)) dx = 0$ For f= [fi, ... fd], Apf ERdid'.

(Corham& Mackey 15') $S(p,q) = sup(E_p(A_qf))^{2}$ $f \in F(E_p(A_qf))^{2}$ - regulies a difficult variational optimization So ... [RKHS Settings] Hd = Hx ... xH, Hmk. RKHS for vector-valued functions k is strictly p.d. Then in (Liu, Lee, Jordan), (Defluttion) For p. 2 EP(Rd). KSQ(p,q) := S(p,q)

 $= \mathbb{E}_{x,x'} \left[(S_2 - S_p)(x) k(x,x') (S_2 - S_p)(x) \right]$ "Score différence" pd. $p(S_2 - S_p) \in L^2(X)$, KSD(p, 2) = 0(=) J=q. (ii). Characterszatton & proporties. [Thm 3.6]. $S(p,q) = E_{x,x'-p}(u_{q}(x,x'))$ with kernel $u(x,x') = S_{e}(x)k(x,x')S_{e}(x') +$ $S_q(x) \nabla_{x'} k(x_x') + \nabla_{x} k(x_x') S_q(x') + tr(\nabla_{x,x} k(x,x))$ Not symmetric w.r.t. (p, ?)! A "special" kernelded MMD. (later). [Proof]. Notice $E_p(A_2f) = E_p((S_2-S_p)f')$. i.e. $E_p(tr(A_{2}f)) = E_p(S_{2}-S_p)'f)$

Apply it on k(x. .) (Ared x), we have $S(p,q) = E_{x,xinp}[(S_{2}-S_{p})^{T}A_{2}k_{x}(x')]$ Apply it again for fixed x'.

[Thm 3.8] For B(x') = Ep(Agkx'(x)) $S(p,q) = \|\beta\|_{\mathcal{H}^d}^2 = \max \left(\operatorname{E_per}(A_q f) \right)^2$ $f \in \mathcal{H}^d$ || f ||_{2€} ≤ | [Proof] $S(p,q) = E_{x,x'} \left[S_{q(x)} - S_{p(x)} \right]^{T} k(x,x') (S_{q(x)} - S_{q(x)})$ $= E_{x,x'} [(S_{2}(x) - S_{p}(x))' < k(x, \cdot), k(\cdot, x') > (\cdots)]$ $= \sum_{l=1}^{6} \langle E_{x} [S_{2}^{l} - S_{p}^{l}] k(x, \cdot)], E_{x'} [k(\cdot, x')(S_{2}^{l} - S_{p}^{l})] \rangle$

= nBlbed, same trick as above

While $\langle f, \beta \rangle_{\mathcal{H}^d} = \sum_{l=1}^d \langle f_l, E_{x_l}, A_{\ell}^l k(x, \cdot) \rangle_{\mathcal{H}^d}$ $= \sum_{l=1}^{d} \langle f_{l}, E_{mp} [S_{q}^{l}(x) k(x, \cdot) + \nabla_{x_{l}} k(x, \cdot)] \rangle_{\mathcal{H}}$ $= \sum_{l=1}^{d} E_{mp} \left[S_{q}^{l}(x) < f_{l}, k(x, \cdot) \right]_{\ell} + < f_{l}, \nabla_{q_{l}} k(x, \cdot) \right]_{\ell}$ $(\nabla_{x}f(x) = \langle f(\cdot), \nabla_{x}k(x, \cdot) \rangle_{\mathcal{H}})$ $= \sum_{l=1}^{d} E_{x-p} \left[S_{2}^{l}(x) f_{l}(x) + \nabla_{x_{l}} f_{l}(x) \right]$ = $E_p(tr(Aqf))$. which Allishes the proof.

We are in a situation to give an estimation of S(p, g): From [Thm 36], "U-statistics" $\widehat{S}_{n}(p,q) = \frac{1}{n(n-i)} \sum_{i \neq j} U_{q}(x_{i},x_{j})$ can be used in application (Boststrapple

The connection with . Fisher duergence & MMD. 1. Fisher FCP, q) = Ep | Sp - Sql² KSD is a "bernetized" Fisher duergence. $\frac{1}{\|S_{2}-S_{1}\|_{2}^{2}}F(P, 2) \leq S(P, 2) \leq E_{X, x', p}k(x, x') F(P, 2)$

11-2-1.0e $E_{p}[tr(A_{q}f)] = E_{p}[(S_{q}-S_{p})^{T}f] \quad \text{for}$ $\|f\|_{\mathcal{H}^{d}} \leq 1.$ 2. MMD kernelided vorsion: $\sup_{f \in \mathcal{H}} \{ E_{p}f - E_{q}f, \|f\|_{\mathcal{H}} \leq l \}$ $= E_{x,x',p} \left[k(x,x') + k(y,y') - 2k(x,y) \right].$ y. y'n q (Cretton /o') Notice that Vr. $E_{y,y'-2} u_{2}(y,y') - 2u_{2}(x,y') = 0$. from Stehn identity. Thus KSD is a MMD with asymmetric kernel Uq. MMD: 2 sample fest KSD: Goodness-of-Att test.

Tor
$$p \neq q$$
, "asymptotically normal"
 $Jn(\hat{S}_{u}(p,q) - Scp,q)) \stackrel{d}{\rightarrow} N(0, \sigma_{u}^{2})$
where $\sigma_{u}^{2} = Var_{x-p} E_{x,p} U_{q}(x, x') > 0$
 $For p = q$, $\Rightarrow \sigma_{u}^{2} = 0$
 $n \hat{S}_{u}(p,q) \stackrel{d}{\rightarrow} \sum_{j=1}^{\infty} q_{j}(z_{j}^{2}-1)$
 $n \hat{S}_{u}(p,q) \stackrel{d}{\rightarrow} \sum_{j=1}^{\infty} q_{j}(z_{j}^{2}-1)$
 $eigen$.
At standard results of U-stat.



(Lu&Lu: A Universal Approximation Thm (2020) of DNN for expressing distributions $MMD(p,\pi) = \sup_{\|f\|_{\infty}} |E_pf - E_{\pi}f|$ $KSD(p,\pi) := \sup_{\substack{\|f\| \leq 1 \\ y \in \mathcal{I}}} E_p(\nabla \ln \pi \cdot f + \nabla \cdot f)$ [Three 4.1] $\{X_i\}_{i=1}^n \sim \pi, P_n = \frac{i}{n} \Sigma S_{X_i}$. Then there exists realizations of Pn s.t. followly magnalithes hold.

(1) for $\pi \in P_3(\mathbb{R}^d)$. $W_1(P_{n},\pi) \leq \begin{cases} \frac{C}{\sqrt{n}} & d=1\\ \frac{C\log n}{\sqrt{n}} & d=2\\ \frac{\sqrt{n}}{\sqrt{n}} & d=2 \end{cases}$

J(2). for pd, bounded k, $MMD(P_n, \pi) \leq \frac{C}{\sqrt{n}}$

(3). for k has bold derNatilles X sub-Gaussian smooth T with Vhr L-Lip, $KSD(P_n, \pi) \in C_n^{d}$ (Or you can understand above results as: "with prob. at least 1-S, (\leq) holds for C = C(S...)independent of n) [Troop]. (1). (2). Well-known Sketch of (2): Here, (K) MMD $(P_n, \pi) = \|\int_{K}^{\infty} k(\pi, \cdot) d(P_n - \pi)\|_{\mathcal{H}}$ $=: \varphi(\chi_1, \dots, \chi_n). \qquad (\frac{\text{Sriperumbulur 16}}{\text{or survey G}})$ y satisfies (g(x1,...xn) - g(x1,...xn)) < 2 The sup k(xix). Then from $\frac{McDiarmid's Ineq}{gloefflipg-type} , with prob. 1-e^{-z}}{gloefflipg-type} = \frac{1-e^{-z}}{1-e^{-z}}$ $\|\int_{P} e^{(x,x)} d(P_n - \overline{x})\|_{\mathcal{H}} \leq E(\dots) + \frac{1-e^{-z}}{n}$

By standard symmetrization argument. Ell $\int k(\cdot, X) d(\widehat{R}_{1} - \overline{X}) \|_{\mathcal{H}} \leq 2 E E \| \frac{1}{n} \sum_{i=1}^{n} E_{i} k(\cdot, X_{i}) \|_{\mathcal{H}_{k}}$ "Rademacher averages". i.i.d. $E_{i} - \pm 1$ Use McDiarmid agah. RHS $\leq E_{2} \| \frac{1}{n} \sum_{i=1}^{n} E_{i} k(\cdot, X_{i}) \|_{\mathcal{H}_{i}} + \frac{2 |\widehat{R}_{2} \overline{Z}_{i}}{n}$ $\leq (E_{2} \| \frac{1}{n} \sum_{i=1}^{n} k(\cdot, X_{i}) \|_{\mathcal{H}_{i}})^{2} + \sqrt{2} \leq \frac{C(K_{0}, C)}{\sqrt{2}}$

(3) Sketch : KSD(Pn, Tc) = 12 Un(Xi, Xj) Use Bernstelh-type Meg. for V-statistics. (Symmetric bernel)

Check this bernel us satisfies the conditions of (Thm C.I): degenerate; $|u_{\mathcal{R}}(x,y)| \leq g(x)g(y)$; $E[g(x)^{k}] \leq g^{2}J^{k}E!/2;$ $\Rightarrow P(|KgD^{2}| \geq t) \leq Cexp(-\frac{C_{x}nt}{g^{2}+JFE}).$

(iii) Orighel SVGD (Lu. Wang 16 NourIPS) · Varlational Inference Problem. Find 2^{*} = arandh EKL(911P)} GeQ, Qasimpler distribution set. 9[T](Z) be the density of Z=T(X), $Q_{TT}(Z) = q(T(Z)) \cdot |det(\nabla_Z T(Z))|.$ Previous methods: consider T with certaln parametric form, then optimize parameters. In [Thm 3.1], they noticed that if "Directly applying gradient descent to solve above problem." For $x \sim q \cdot q_{TT}$ dethed above, $T(x) = x + 2\phi x$ $\Rightarrow d \in [KL(P_{TT}]|P) = -E_{q}(+(A_{p}\phi))$ (K) $\left(\xrightarrow{dE} (2) = \int V(\nabla 2 + 2\nabla \ln p) dx \right)$

 $\partial_{t} \mathcal{Q} = -\mathcal{U}(\mathcal{Q}\mathcal{V})$ Agath, they were able to minimize (*) in {II \$ 11 \$ 12 \$ < Sq. p) }: the "direction of steepest descent" is $V(x) = E_{y \sim q}(k(x,y)\nabla_y \ln p(y) + \nabla_y k(x,y))$ i.e. for this v, $\frac{\partial}{\partial t} q_t = -grad_{\mathcal{H}} K L(q_t \| p)$ The metric We will be discussed later. Advantages, "Stem" Geometric structure --- (17-20?) Duncan et.al. 19]

D. Under which condition we need to consider this gradient flow (or use

"Stell distance between measures)? What can ue benefit from moving particles smoothly? (2). Can ve find some alternative "positive" flows to improve SVGD (w.r.t. Cvg rate; "outliers")? 3. Understand the bias and variance of SVGD particles, or combine it with traditional MC. [from LN17'].

(Nathan Ross) - Fundamentals of

Stehs method.

(I Bortson) - Approximation of distri. of V-stat. with multidhuensional kornels.

Thank you o