

Smooth Transform

① let P, Q be two smooth density on \mathbb{R}

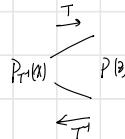
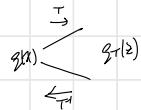
and $T = T(x)$ be an one-to-one transform on \mathbb{R} differentiable w.r.t x

g_T be the density of $Z = T(X) \sim Q$

$$\nabla_Q \text{KL}(Q||P) = -E_{x \sim P} \left[\langle D\ln g_T^{-1}, \nabla_Q \ln g_T \rangle + \langle \nabla_Q \log p(z)|_{z=T(x)}, \nabla_Q \ln g_T \rangle \right]$$

$$\begin{aligned} \text{KL}(Q||P) &= \int_{\mathbb{R}} g_T(z) \log \frac{g_T(z)}{p(z)} dz \\ &= \int_{\mathbb{R}} g_T(T(x)) \log \frac{g_T(T(x))}{p(T(x))} dx \\ &= \int_{\mathbb{R}} \frac{g_T(x)}{|DT(x)|} \log \frac{g_T(x)}{p_T(x)} |DT(x)| dx \\ &= \int_{\mathbb{R}} \frac{g_T(x)}{|DT(x)|} \log \frac{g_T(x)}{p_T(x)} |DT(x)| dx \\ &= \text{KL}(Q||P_T) \end{aligned}$$

let $z = T(x) \quad x = T^{-1}(z)$



$$\begin{aligned} \nabla_Q \text{KL}(Q||P) &= \nabla_Q \int_{\mathbb{R}} g_T(x) \log \frac{g_T(x)}{p_T(x)} dx \\ &= -E_{x \sim P} \left[\nabla_Q \log p_T(x) \right] \\ &= -E_{x \sim P} \left[\nabla_Q \log (p(T(x)) \cdot |DT(x)|) \right] \\ &= -E_{x \sim P} \left[\nabla_Q \log \det(DT(x)) + \nabla_Q \log p(T(x)) \right] \\ &= -E_{x \sim P} \left[\langle D\ln g_T^{-1}, \nabla_Q \ln g_T \rangle + \langle \nabla_Q \log p(z)|_{z=T(x)}, \nabla_Q \ln g_T \rangle \right] \end{aligned}$$

② consider the mapping $T(x) = x + \phi(x)$ where ϕ is a smooth function

③ When $|\phi|$ is sufficiently small, T is a one-to-one mapping

$$\text{④ } \nabla_Q \text{KL}(Q||P)|_{\phi=0} = -E_{x \sim P} \left[\underbrace{\text{trace}(\nabla_Q \log p(z) \phi(x) + D\phi(x)^T)}_{\text{Stein operator } (\nabla_Q \phi)(x)} \right]$$

⑤ When $|\phi|$ is sufficiently small

$D\phi(x) = I + \phi'(x)$ is full rank

By inverse function theorem, a map is locally invertible if its linearization is invertible

$$\text{⑥ } D\phi(x)|_{\phi=0} = I + \phi'(0) = I$$

$$\nabla_Q D\phi(x)|_{\phi=0} = D\phi'(x)$$

$$\nabla_Q D\phi(x)|_{\phi=0} = \phi'(x)$$

$$\nabla_Q \text{KL}(Q||P) = -E_{x \sim P} \left[\langle D\ln g_T^{-1}, \nabla_Q \ln g_T \rangle + \langle \nabla_Q \log p(z)|_{z=T(x)}, \nabla_Q \ln g_T \rangle \right]$$

$$= -E_{x \sim P} \left[\text{trace}(D\phi(x)) + \langle \nabla_Q \log p(z)|_{z=T(x)}, \phi(x) \rangle \right]$$

Let $\beta(x) = E_{x \sim p}[(A_p k(x))]$,

then $\|\beta\|_{2^2} = \sqrt{\beta(\beta)}$.

$$\langle \phi, \beta \rangle_{2^2} = E_{x \sim p}[\text{trace}(A_p \phi(x))] = -\nabla_{\beta} \text{KL}(g_p \| p)|_{\beta=0}$$

Consider all perturbations $\|\phi\|_{2^2} \leq \sqrt{S_k(p)} = \|\beta\|_{2^2}$

to make a "steepest descent".

$$\max_{\phi} \langle \phi, \beta \rangle_{2^2}$$

$$\text{So } \|\phi\|_{2^2} \leq \|\beta\|_{2^2}$$

$$\text{Then } \phi^* = \beta, \nabla_{\beta} \text{KL}(g_p \| p)|_{\beta=0} = -\langle \phi^*, \beta \rangle_{2^2} = -\|\beta\|_{2^2}^2 = -S_k(p)$$

↓

Suggests an iterative update

$$x := x + \epsilon \cdot \phi^*(x)$$

$$= x + \epsilon f(x)$$

$$= x + \epsilon E_{x \sim p}[(A_p k(x))(x)]$$

$$= x + \epsilon E_{x \sim p}[\nabla_{\beta} \text{KL}(g_p \| p)(k(x)) + \nabla_{\beta} \text{KL}(g_p \| p)]$$

next section shows that
this is also a functional gradient

Stein Variational Gradient Descent

Consider the transformation $T(x) = x + f(x)$

at each iteration, apply $x := x + f(x)$ to minimize $\text{KL}(g \| p)$

1' let $T(x) = x + f(x)$, where $f \in \mathcal{F}^2$ and g be the density of $z = T(x)$ wrt

$$\text{then } D_f \text{KL}(g \| p)|_{f=0} = -E_{x \sim p}[\nabla_{\beta} \log p(x) k(x) + \nabla_{\beta} k(x)]$$

$$\text{let } T(x) = x + f(x), \quad \hat{T}(x) = x + f(x) + \epsilon g(x)$$

$$\text{let } T(f) = \text{KL}(g \| p) = \text{KL}(g \| p_{\hat{T}}) \quad F(f+eg) = F(f) + \epsilon \langle Df(f), g \rangle_{2^2} + O(\epsilon^2)$$

$$F(f+eg) = \text{KL}(g \| p_{\hat{T}})$$

$$= \int_x g(x) \log \frac{g(x)}{p_{\hat{T}}(x)} dx$$

$$= E_{x \sim p}[\log g(x) - \log p(\hat{T}(x)) - \log \det(D\hat{T}(x))]$$

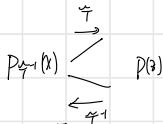
$$= E_{x \sim p}[\log g(x) - \log p(T(x)) - \log \det(DT(x))]$$

$$+ \log p(T(x)) + \log \det(DT(x))$$

$$- \log p(T(x)) - \log \det(DT(x))]$$

$$= \text{KL}(g \| p_{\hat{T}}) + E_{x \sim p}[\log p(T(x)) - \log p(\hat{T}(x))] \quad \textcircled{1}$$

$$+ E_{x \sim p}[\log \det(DT(x)) - \log \det(D\hat{T}(x))] \quad \textcircled{2}$$



$$0 = \mathbb{E}_{x,y} [\log p(x, f(x)) - \log p(x, f(x) + \epsilon g(x))]$$

$$= -\epsilon \mathbb{E}_{x,y} [\nabla_{\epsilon} \log p(x)|_{\epsilon=x, f(x)}^T g(x)] + o(\epsilon)$$

$$= -\epsilon \mathbb{E}_{x,y} [\langle \nabla_{\epsilon} \log p(x)|_{\epsilon=x, f(x)}, \begin{bmatrix} \langle g, k(x_i) \rangle \\ \vdots \end{bmatrix} \rangle] + o(\epsilon)$$

$$= -\epsilon \cdot \left\langle \mathbb{E}_{x,y} [\nabla_{\epsilon} \log p(x)|_{\epsilon=x, f(x)}], g \right\rangle_{\mathcal{H}} + o(\epsilon)$$

$$0 = \mathbb{E}_{x,y} [\log \det(I + Df(x)) - \log \det(I + Df(x) + \epsilon Dg(x))] + o(\epsilon)$$

$$= -\epsilon \mathbb{E}_{x,y} [\langle (I + Df(x))^T, Dg(x) \rangle] + o(\epsilon)$$

$$= -\epsilon \mathbb{E}_{x,y} \left[\sum_j (I + Df(x))_j^T \frac{\partial}{\partial x_j} g(x) \right] + o(\epsilon)$$

$$= -\epsilon \mathbb{E}_{x,y} \left[\sum_j (I + Df(x))_j^T \langle \frac{\partial}{\partial x_j} k(x_i), g \rangle_{\mathcal{H}} \right] + o(\epsilon)$$

$$= -\epsilon \mathbb{E}_{x,y} \left[\sum_j \left\langle \sum_i (I + Df(x))_j^T \frac{\partial}{\partial x_j} k(x_i), g \right\rangle_{\mathcal{H}} \right] + o(\epsilon)$$

$$= -\epsilon \mathbb{E}_{x,y} \left[\sum_i \langle (I + Df(x))_i^T \nabla_x k(x_i), g \rangle_{\mathcal{H}} \right] + o(\epsilon)$$

$$= -\epsilon \langle \mathbb{E}_{x,y} [(I + Df(x))^T \nabla_x k(x_i)], g \rangle_{\mathcal{H}} + o(\epsilon)$$

$$F(f+eg) = F(f) - \epsilon \langle \mathbb{E}_{x,y} [\nabla_{\epsilon} \log p(x)|_{\epsilon=x, f(x)} k(x_i)], g \rangle$$

$$-\epsilon \langle \mathbb{E}_{x,y} [(I + Df(x))^T \nabla_x k(x_i)], g \rangle + o(\epsilon)$$

$$DF(f) = -\mathbb{E}_{x,y} [\nabla_{\epsilon} \log p(x)|_{\epsilon=x, f(x)} k(x_i) + (I + Df(x))^T \nabla_x k(x_i)]$$

$$\text{at } f=0, \quad DF(f)|_{f=0} = -\mathbb{E}_{x,y} [\nabla_{\epsilon} \log p(x)|_{\epsilon=x, f(x)} + \nabla_x k(x_i)]$$

Algorithm:

$$Df_kl(g||p)|_{f=0} = -\mathbb{E}_{x,y} [\nabla_x \log p(x) k(x_i) + \nabla_x k(x_i)]$$

$$x^{(t)} := x^{(t)} - t \cdot Df_kl(g||p)|_{f=0}$$

$$= x^{(t)} + t \mathbb{E}_{x,y} [\nabla_x \log p(x) k(x, x^{(t)}) + \nabla_x k(x, x^{(t)})]$$

$$\approx x^{(t)} + t \frac{1}{m} \sum_{j=1}^m [\nabla_x \log p(x^{(t)}) k(x^{(t)}, x^{(t)}) + \nabla_x k(x^{(t)}, x^{(t)})]$$

$$\text{Eq. } N(x; u, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-u)^T \Sigma^{-1} (x-u)\right)$$

$$\nabla_x N(x; u, \Sigma) = -\frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-u)^T \Sigma^{-1} (x-u)\right) \Sigma^{-1} (x-u)$$

$$k(x_i, x_i) = \exp(-b \|x_i - x_i\|^2)$$

$$\nabla_x k(x_i, x_i) = \exp(-b \|x_i - x_i\|^2) \cdot -2b(x_i - x_i)$$