**Stein Operator** $(\mathcal{J}_p f)(x) := \langle \nabla_x \log p(x), f(x) \rangle + \nabla_x \cdot f(x)$

we know $\mathbb{E}_p[\mathcal{J}_p f(x)] = \int \langle \nabla_x p, f \rangle + p \cdot \nabla_x^T f = \int \nabla(pf) = 0$, $\forall f$

$\min_{\{x_i\}} KL[\{x_i\} \| p]$, given $p$ $\qquad q(x) = \frac{1}{N}\sum_{i=1}^{N} \mathbb{1}\{x-x_i\}$    empirical dist. of $x'$

$x_i' \longleftarrow x_i + \varepsilon \phi(x_i)$, $\quad \phi = \arg\max_{\phi \in \mathcal{F}} \{ KL[q\|p] - KL[q_{[\varepsilon\phi]}\|p] \}$

$\qquad\qquad\qquad\qquad = \arg\max_{\phi \in \mathcal{F}} \{ -\frac{d}{d\varepsilon} KL[q_{[\varepsilon\phi]}\|p]|_{\varepsilon=0} \}$

$\qquad\qquad\qquad\qquad\qquad \overset{(*)}{=} \mathbb{E}_{x \sim q}[\mathcal{J}_p \phi(x)]$

**(*) Proof:** $\quad T(x) \overset{\Delta}{=} x + \varepsilon \cdot \phi(x)$

$KL[q_{[T]}\|p] = KL[q\|p_{[T^{-1}]}] = \int q \cdot (\log q - \log p_{[T^{-1}]}(x)) dx$

$\log p_{[T^{-1}]}(x) = \log \det(I + \varepsilon \nabla \phi(x)) + \log p(x + \varepsilon \phi(x))$

$\qquad\qquad = \varepsilon \cdot \text{Tr}(\nabla\phi) + \log p(x) + \varepsilon \cdot \nabla \log p(x)^T \phi(x)$

$\Rightarrow \frac{d}{d\varepsilon} KL[q_{[T]}\|p] = \frac{d}{d\varepsilon}[-\int q(x) \log p_{[T^{-1}]}(x) dx]$

$\qquad\qquad\qquad\qquad = -\int q(x) \cdot [\text{Tr}(\nabla\phi) + \nabla \log p(x)^T \phi(x)] dx$

**Functional** $T(x) := x + f(x)$ $\qquad\qquad\qquad \langle \mathcal{J}_p k(x,\cdot), f(\cdot) \rangle$

**derivative** $KL[q\|p_{[T^{-1}]}] = \int q(x) \cdot [\log q(x) - \log p(x) - \underset{\|}{(\mathcal{J}_p f)(x)}] dx$ (when $f \approx 0$)

**in RKHS** $\Rightarrow \frac{\delta}{\delta f} KL[q_{[T]}\|p]|_{f=0} = -\int q(x) \mathcal{J}_p k(x,\cdot) dx = \phi_{q,p}^*(\cdot)$

**(kernelized)** $\mathbb{D}_{\mathcal{F}}(p\|q) := \max_{\phi \in \mathcal{F}} \mathbb{E}_{x \sim q}[\mathcal{J}_p \phi(x)] \overset{\mathcal{F}=\mathcal{H}}{=\!=} \max_{\phi \in \mathcal{H}} \{ \mathbb{E}_{x \sim q} \mathcal{J}_p \phi(x) \mid \|\phi\|_{\mathcal{H}} \leq 1 \}$

**Stein Discrepancy**
$\Rightarrow \phi^*(\cdot) \propto \mathbb{E}_{x \sim q}[\mathcal{J}_p k(x,\cdot)]$ $\qquad\qquad \underset{\|}{\mathbb{E}_{x \sim q} \langle k(x,\cdot), \mathcal{J}_p \phi(\cdot) \rangle_{\mathcal{H}}}$

$\Rightarrow \mathbb{D}^2(p\|q) = \mathbb{E}_{x,x' \sim q}[\mathcal{J}_p^x \mathcal{J}_p^{x'} k(x,x')]$ $\qquad \underset{\|}{\langle \mathbb{E}_{x \sim q} \mathcal{J}_p k(x,\cdot), \phi(\cdot) \rangle_{\mathcal{H}}}$

**Goodness** $\Rightarrow \mathbb{D}(\{x_i\}\|p) = \frac{1}{n(n-1)} \sum_{i \neq j} \mathcal{J}_p^x \mathcal{J}_p^{x'} k(x_i, x_j)$ $\quad (\langle k(x,\cdot), \phi(\cdot) \rangle = \phi(x))$

**of fit test:** whether $\mathbb{D}(\{x_i\}\|p) \geq \cdots$

**SVGD** $\quad x_j \longleftarrow x_j + \varepsilon \cdot \hat{\mathbb{E}}_{x \sim \{x_i\}_{i=1}^{n}}[\mathcal{J}_p k(x, x_j)]$ $\qquad$ repulsive force

$\qquad\qquad\qquad\qquad = [\nabla \log p(x) \cdot k(x, x_j) + \nabla_x k(x, x_j)]$

$\quad$ (MAP: $x_j \leftarrow x_j + \varepsilon \nabla \log p(x_j)$) $\qquad$ move to high $p(x)$, general case of MAP

**de-Bruijin's Identity**

If $\phi_{q,p}(x) := \nabla_x \log \frac{p(x)}{q(x)}$, $T(x) = x + \varepsilon \phi_{q,p}(x)$

$\frac{d}{d\varepsilon} KL[q_{[T]} \| p]|_{\varepsilon=0} = -\mathbb{E}_{x \sim q}[\mathcal{J}_p \phi(x)] = -\int \frac{q}{p} \nabla(p\phi) = \int \nabla(\frac{q}{p}) \cdot p\phi$

$= \int \frac{\nabla q \cdot p - q \nabla p}{p} \phi = \int q(\nabla \log q - \nabla \log p) \phi = -\mathbb{E}_{x \sim q}[\| \nabla \log \frac{p(x)}{q(x)} \|_2^2]$

$\underset{\text{Fisher divergence}}{\underbrace{}}$

---

**Fokker-Planck Eq. derivation**

**(Continuity equation)**

$dx/dt = \phi(x) \implies \frac{dp(x)}{dt} = -\nabla \cdot (p(x) \phi(x))$

$\qquad\qquad \hookrightarrow \neq dp_t(x_t)/dt$

$T(x) := x + \varepsilon \phi(x) \implies T^{-1}(x) \simeq x - \varepsilon \phi(x) + O(\varepsilon)$

$\log \tilde{p}(x) = \log p(T^{-1}(x)) + \log \det(\nabla_x T^{-1}(x))$

$\qquad = \log p(x - \varepsilon \phi(x)) + \log \det(I - \varepsilon \nabla \phi(x)) + O(\varepsilon)$

$\qquad = \log p(x) - \varepsilon \nabla \log p^T \phi(x) - \varepsilon \operatorname{Tr}(\nabla \phi) + O(\varepsilon)$

$\qquad\qquad\qquad\qquad \hookrightarrow -(\mathcal{J}_p \phi)(x)$

$\implies \frac{d}{dt} \log p_t(x) = -(\mathcal{J}_p \phi)(x)$

$\implies \frac{d}{dt} p_t(x) = -p(x) \cdot \mathcal{J}_p \phi(x) = -\nabla \cdot (p(x) \phi(x))$

---

**Regularized Stein Discrepancy**

$\text{RSD}(p \| q; f) = \mathbb{D}(p \| q; f) - \frac{1}{2} \|f\|^2_{L^2(q)}$

$\qquad = \mathbb{E}_q[\nabla \log p^T f + \operatorname{div}(f)] - \frac{1}{2} \mathbb{E}_{x \sim q}[\| f(x) \|^2]$

$\qquad = \mathbb{E}_q[(\nabla \log \frac{p}{q})^T f] - \frac{1}{2} \mathbb{E}_q[\| f(x) \|^2]$

$\qquad \implies f^* = \nabla \log \frac{p}{q}$

# Gradient Flows

$$P_2(\mathcal{M}) = \{\rho : \mathcal{M} \to \mathbb{R}_{\geq 0} \mid \int_{\mathcal{M}} d\rho = 1, \quad \int_{\mathcal{M}} |x|^2 \rho(x) dx < +\infty\}$$

$$D_{KL}(\rho, \pi) = \int_{\mathcal{M}} (\log \rho(x) - \log \pi(x)) \cdot \rho(x) dx$$

$$W_2^2(\mu, \nu) = \inf_{p \in \Pi(\mu, \nu)} \int |x - y|^2 \, dp(x, y)$$

**SVGD**
**Gradient Flows**
**mean-field limits**

$$\frac{dx}{dt} = \int [k(x', x) \nabla_{x'} \log \pi(x') + \nabla_{x'} k(x', x)] \rho(x') dx'$$

$$= \int k(x', x) \nabla_{x'} \log \pi(x') \rho(x') dx' + \int \nabla_{x'} k(x', x) \rho(x') dx'$$

$$\hookrightarrow -\int k(x', x) \nabla_{x'} \rho(x') dx'$$

$$= \int k(x', x) \nabla_{x'} [\log \pi(x') - \log \rho(x')] \rho(x') dx'$$

$$\frac{dx}{dt} = \mathbb{E}_{x' \sim \rho}[k(x', x) \nabla_{x'}(\log \pi(x') - \log \rho(x'))] = \mathcal{K}_\rho \nabla_x (\log \pi(x) - \log \rho(x))$$

$$(\mathcal{K}_\rho \phi)(x) := \mathbb{E}_{x' \sim \rho}[k(x', x) \phi(x')]$$

**Liouville Eq.**

$$\frac{\partial \rho(x)}{\partial t} = -\nabla \cdot [\quad \downarrow \quad]$$

$$= \nabla \cdot [\int k(x', x) \nabla_{x'} \frac{\delta}{\delta \rho} D_{KL}(\rho, \pi) \rho(x') dx'] = \nabla \cdot (\rho(x) \mathcal{K}_\rho \nabla_{x'} \frac{\delta}{\delta \rho} D_{KL}(\rho, \pi))$$

似乎是KL为 objective, 在 Wasserstein metric 下 ρ 的 梯度流          **?** $\nabla_{W_2} D_{KL}(\rho, \pi)$

**Liouville Eq.**

$$\frac{\partial \rho(x)}{\partial t} = \nabla \cdot [\rho(x) \nabla \frac{\delta}{\delta \rho} D_{KL}(\rho, \pi)] = \nabla \cdot [\rho(x) \nabla (\log \rho(x) - \log \pi(x))]$$

$$= -\nabla \cdot [\rho(x) \nabla \log \pi(x)] + \nabla^2 \rho(x)$$

**mean-field**
**Wasserstein**
**dynamics**

$$\frac{dx}{dt} = -\nabla \cdot [\log \rho(x) - \log \pi(x)]$$

$$\tilde{x}_i \leftarrow x_i - \varepsilon \cdot [\nabla \log \rho(x_i) - \nabla \log \pi(x_i)]$$

$$\rho(x) \approx \frac{1}{n} \sum_{i=1}^{n} k(x, x_i) \Rightarrow \nabla \log \rho(x) = \nabla \rho(x)/\rho(x) \approx \frac{\sum_i \nabla_x k(x, x_i)}{\sum_i k(x, x_i)}$$

$\text{div}(\cdot)$

**$W_2$ grad flow** $\quad \partial_t \rho_t = \nabla(\rho_t \nabla_{W_2} F[\rho_t]) \qquad \nabla_{W_2} F[\rho] := \nabla \frac{\delta}{\delta \rho} F[\rho] \quad ?$

**continuity Eq.** $\quad \partial_t \rho_t = \text{div}(\rho_t \frac{dx}{dt}) \implies \frac{dx}{dt} = -\nabla_{W_2} F[\rho_t]$

not sure

$\partial_t F[\rho_t] = \partial_t \rho_t \cdot \delta F[\rho_t] = \nabla(\rho_t \frac{dx}{dt}) \cdot \delta F[\rho_t] = \mathbb{E}_{\rho_t}[\langle \nabla_{W_2} F[\rho_t], \frac{dx}{dt} \rangle]$ ?  divergence积分符号对不上

$\qquad\qquad\qquad\qquad\qquad\qquad = -\mathbb{E}_{\rho_t}[\|\nabla_{W_2} F[\rho_t]\|^2]$

**$f$-divergence** $\quad D_f[\rho \| \pi] = \mathbb{E}_\pi[f(\rho/\pi)]$

$D_f[\rho + \varepsilon v \| \pi] - D_f[\rho \| \pi] = \int \pi \cdot [f(\frac{\rho}{\pi} + \varepsilon \frac{v}{\pi}) - f(\frac{\rho}{\pi})] \approx \int \pi \cdot f'(\frac{\rho}{\pi}) \cdot \frac{v}{\pi} \cdot \varepsilon$

$\implies \frac{\delta}{\delta \rho} D_f[\rho \| \pi] = f'(\rho/\pi)$

$\implies \nabla_{W_2} D_f[\rho \| \pi] = \nabla \frac{\delta}{\delta \rho} D_f[\rho \| \pi] = \nabla f'(\rho/\pi)$

$KL: f(t) = t \log t, \quad f'(t) = \log t + 1 \qquad \chi^2: f(t) = t^2 - 1, \quad f'(t) = 2t,$

**结论:** $\quad \nabla_{W_2} KL[\rho \| \pi] = \nabla_x \log \frac{\rho}{\pi} \qquad\qquad \nabla_{W_2} \chi^2[\rho \| \pi] = 2 \nabla(\frac{\rho}{\pi})$

We know SVGD is kernelized KL grad flow

**kernelized** $-\frac{dx}{dt} = \mathcal{K}_\pi \nabla_{W_2} \chi^2[\rho \| \pi] = \mathbb{E}_{x' \sim \pi}[k(x, x') \, 2 \nabla(\frac{\rho(x')}{\pi(x')})]$
**$\chi^2$ grad flow**
$\qquad\qquad\qquad\qquad = 2 \int k(x, x') [\nabla \rho(x') - \nabla \log \pi(x') \rho(x')] dx'$

$\qquad\qquad\qquad\qquad = 2 \int k(x, x') \nabla \log \frac{\rho(x')}{\pi(x')} \rho(x') dx' = 2 \mathcal{K}_\rho \nabla_{W_2} KL[\rho \| \pi]$

本质上是 $\int \pi \nabla(\frac{\rho}{\pi}) = \int \rho \nabla(\log \frac{\rho}{\pi})$