

The contents comes from the paper

"A kernelized Stein Discrepancy for Goodness-of-Fit Tests"

Some notations are different (eg. p and q are switched when they co-appear)

Stein's class

1° let $\{x \in \mathbb{R}^d\}$ be a continuously differentiable distribution supported on \mathcal{X}

a function $f: \mathcal{X} \rightarrow \mathbb{R}$ is in the Stein's class of p if

and $\int_{\mathcal{X}} \nabla_x (p(x) f(x)) dx = 0$

2° $f: \mathcal{X} \rightarrow \mathbb{R}$ is in the Stein's class of p if f is ∇ -integrable

3° f can be shown to be in Stein's class of p if

or $\left\{ \begin{array}{l} 1. \lim_{|x| \rightarrow \infty} p(x) f(x) = 0 \text{ if } p(x) \text{ is bounded and } \lim_{|x| \rightarrow \infty} p(x) = 0 \\ 2. \mathcal{X} \text{ is compact with piecewise continuous boundary} \\ \quad p(x) f(x) = 0 \text{ for } x \in \partial \mathcal{X} \\ \vdots \end{array} \right.$

Stein's Identity

1° let p be a smooth distribution supported on $\mathcal{X} \subseteq \mathbb{R}^d$, $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$

the Stein's operator of p is a linear operator

$$(\mathcal{A}_p f)(x) = \nabla(p(x) f(x)) + p(x) \nabla f(x)$$

$$\mathcal{A}_p f = \begin{bmatrix} \nabla f_1 \\ \vdots \\ \nabla f_m \end{bmatrix}$$

2° Assume p is a smooth density supported on \mathcal{X} , f is in Stein's class of p
then $\mathbb{E}_p[\mathcal{A}_p f(x)] = \mathbb{E}_p[\nabla(p(x) f(x)) + p(x) \nabla f(x)] = 0$

$$\begin{aligned} & \int_{\mathcal{X}} p(x) \nabla(p(x) f(x)) + p(x) \nabla f(x) dx \\ &= \int_{\mathcal{X}} \nabla(p(x) f(x)) + p(x) \nabla f(x) dx \\ &= \int_{\mathcal{X}} \nabla(p(x) f(x)) dx \\ &= 0 \end{aligned}$$

3° Suppose p and q are distributions on \mathcal{X}
 f is in Stein's class of q

$$\text{then } \mathbb{E}_{p,q}[(\mathbb{A}f)(x)] = \mathbb{E}_{q,q}[(\nabla \log p(x) - \nabla \log q(x)) f(x)^T]$$

$$\begin{aligned} \mathbb{E}_{p,q}[(\mathbb{A}f)(x)] &= \mathbb{E}_{q,q}[(\mathbb{A}f)(x) - (\mathbb{A}f)(x)] \\ &= \mathbb{E}_{q,q}[\nabla \log p(x) f(x)^T - Df(x)^T - \nabla \log q(x) f(x)^T + Df(x)^T] \\ &= \mathbb{E}_{q,q}[(\nabla \log p(x) - \nabla \log q(x)) f(x)^T] \end{aligned}$$

$$\text{Also } \mathbb{E}_{p,q}[\text{trace}((\mathbb{A}f)(x))] = \mathbb{E}_{q,q}[\underbrace{(\nabla \log p(x) - \nabla \log q(x))^T f(x)}_{\text{Expected difference in score function, measured in } f \text{ direction}}] \quad (\text{when } f: \mathbb{R}^d \mapsto \mathbb{R}^d)$$

$$4^\circ \quad p \neq q \Rightarrow \exists f \text{ s.t. } \mathbb{E}_{p,q}[(\mathbb{A}f)(x)] \neq 0$$

• Stein Discrepancy

1° define the Stein Discrepancy between p and q

$$JS(p,q) = \max_{f \in \mathcal{F}} \mathbb{E}_{p,q}[\text{trace}((\mathbb{A}f)(x))]$$

$$\text{eg. } f(x) = \sum_i w_i f_i(x)$$

f is a linear combination of simple f_i 's

$$\begin{aligned} \mathbb{E}_{p,q}[\text{trace}((\mathbb{A}f)(x))] &= \mathbb{E}_{p,q}[(\nabla \log p(x) - \nabla \log q(x))^T f(x)] \\ &= \sum_i w_i \underbrace{\mathbb{E}_{p,q}[(\nabla \log p(x) - \nabla \log q(x))^T f_i(x)]}_{\beta_i} \end{aligned}$$

$$\mathcal{F} = \left\{ \sum_i w_i f_i(x) : \|w\| \leq 1 \right\}$$

$$\text{then } JS(p,q) = \max_{w \text{ s.t. } \|w\| \leq 1} \sum_i w_i \beta_i$$

Kernelized Stein Discrepancy

1° A kernel $k(x, x')$ is integrally strictly positive definite if

$$\forall g \text{ s.t. } 0 < \|g\|_{L^2}^2 < \infty$$

$$\text{we have } \int \int g(x) k(x, x') g(x') dx dx' > 0$$

$$\text{think this as a continuous version of } g^T K g = \sum_i \sum_j g_i k_{ij} g_j > 0$$

2° the kernelized Stein discrepancy is defined as

$$S_k(q, p) = E_{x, x' \sim q} [(\nabla g(p, x) - \nabla g(q, x))^T k(x, x') (\nabla g(p, x) - \nabla g(q, x))]$$

3° let $g_{q,p}(x) = g(x) (\nabla g(p, x) - \nabla g(q, x))$

assume $\left\{ \begin{array}{l} k(x, x') \text{ is integrally strictly positive definite} \\ p, q \text{ are continuous density s.t. } \|g_{q,p}\|_{L^2}^2 < \infty \end{array} \right.$

then $\left\{ \begin{array}{l} S_k(q, p) > 0 \\ S_k(q, p) = 0 \Leftrightarrow p = q \end{array} \right.$ eg. satisfied when tail of q decays exponentially
eg. not satisfied when q has a heavy tail
(when q is a Cauchy distribution, p is a Gaussian)

if $p = q$, then $\nabla g(p, x) - \nabla g(q, x) = 0$, $g_{q,p}(x) = 0$, then $S_k(q, p) = 0$ obviously

if $S_k(q, p) = 0$

$$\begin{aligned} S_k(q, p) &= E_{x, x' \sim q} [(\nabla g(p, x) - \nabla g(q, x))^T k(x, x') (\nabla g(p, x) - \nabla g(q, x))] \\ &= \iint g(x) (\nabla g(p, x) - \nabla g(q, x))^T k(x, x') (\nabla g(p, x) - \nabla g(q, x)) g(x') dx dx' \\ &= \iint g_{q,p}(x)^T k(x, x') g_{q,p}(x') dx dx' \\ &= 0 \end{aligned}$$

then $\|g_{q,p}\|_{L^2}^2 = 0$, then $p = q$

4° A kernel $k(x, x')$ is in the Stein class of p if

and $\left\{ \begin{array}{l} k \text{ has continuous 2nd partial derivatives} \\ \text{both } k(x, \cdot) \text{ and } k(\cdot, x) \text{ are in the Stein class of } p \text{ s.t. } x \end{array} \right.$

eg. RBF kernel is in the Stein class for smooth densities supported on \mathbb{R}^n

$$\forall u, \lim_{n \rightarrow \infty} k(x, u) = 0 \quad \text{since } \left\{ \begin{array}{l} k(x, u) \text{ bounded s.t.} \\ \lim_{n \rightarrow \infty} p(x) = 0 \end{array} \right.$$

5° If $k(x): \mathbb{R}^d \rightarrow \mathbb{R}$ is in the Stein class of \mathbb{P} & x'

Can show $\nabla k(x): \mathbb{R}^d \rightarrow \mathbb{R}^d$ is also in the Stein class

$$\begin{aligned} & \int_{\mathbb{R}^d} \nabla_x \cdot (\mathbb{P}(x) \nabla \nabla k(x)) dx \\ &= \int_{\mathbb{R}^d} \nabla \mathbb{P}(x) \nabla k(x)' + \mathbb{P}(x) \nabla \nabla k(x) dx \\ &= \int_{\mathbb{R}^d} \nabla_x \cdot [\nabla \mathbb{P}(x) k(x) + \mathbb{P}(x) \nabla k(x)] dx \\ &= \nabla_x \cdot \int_{\mathbb{R}^d} \nabla_x [\mathbb{P}(x) k(x)] dx \\ &= \nabla_x \cdot 0 = 0 \end{aligned}$$

6° Assume $\begin{cases} \mathbb{P}, \mathbb{Q} \text{ are smooth densities} \\ k(x) \text{ is in the Stein's class of } \mathbb{Q} \end{cases}$

then $S_k(\mathbb{Q}, \mathbb{P}) = \mathbb{E}_{x \sim \mathbb{Q}} [\nabla \mathbb{P}(x)' k(x) \nabla \mathbb{P}(x) + \nabla \mathbb{P}(x)' \nabla k(x) + \nabla k(x)' \mathbb{P}(x) + \text{trace}(\nabla \nabla k(x))]$

evaluate $S_k(\mathbb{Q}, \mathbb{P})$ in a way so only $\nabla \mathbb{P}(x)$ is used, no $\nabla \mathbb{Q}(x)$

$$\begin{aligned} S_k(\mathbb{Q}, \mathbb{P}) &= \mathbb{E}_{x \sim \mathbb{Q}} [(\nabla \mathbb{P}(x) - \nabla \mathbb{Q}(x))' k(x) (\nabla \mathbb{P}(x) - \nabla \mathbb{Q}(x))] \\ &= \mathbb{E}_{x \sim \mathbb{Q}} [(\nabla \mathbb{P}(x) - \nabla \mathbb{Q}(x))' (\nabla \mathbb{P}(x) k(x) + \nabla k(x) - \nabla \mathbb{Q}(x) k(x) - \nabla k(x))] \quad \text{use the notation } k(x) = k(x') \\ &= \mathbb{E}_{x \sim \mathbb{Q}} [(\nabla \mathbb{P}(x) - \nabla \mathbb{Q}(x))' ((A \nabla k)(x) - (A \nabla k)(x'))] \\ &= \mathbb{E}_{x \sim \mathbb{Q}} [(\nabla \mathbb{P}(x) - \nabla \mathbb{Q}(x))' (A \nabla k)(x)] \quad \mathbb{E}_{x \sim \mathbb{Q}} [(A \nabla k)(x)] = 0 \\ &= \mathbb{E}_{x \sim \mathbb{Q}} [\nabla \mathbb{P}(x)' (\nabla \mathbb{Q}(x) k(x) + \nabla k(x))] - \mathbb{E}_{x \sim \mathbb{Q}} [\nabla \mathbb{Q}(x)' (\nabla \mathbb{P}(x) k(x) + \nabla k(x))] \\ &= \underbrace{\mathbb{E}_{x \sim \mathbb{Q}} [\nabla \mathbb{P}(x)' k(x) \nabla \mathbb{P}(x)]}_{\Phi} + \mathbb{E}_{x \sim \mathbb{Q}} [\nabla \mathbb{P}(x)' \nabla k(x)] - \mathbb{E}_{x \sim \mathbb{Q}} [\nabla \mathbb{Q}(x)' k(x) \nabla \mathbb{P}(x)] - \mathbb{E}_{x \sim \mathbb{Q}} [\nabla \mathbb{Q}(x)' \nabla k(x)] \\ &= \Phi - \mathbb{E}_{x \sim \mathbb{Q}} [(\nabla \mathbb{Q}(x)' k(x) + \nabla k(x)' - \nabla k(x)') \nabla \mathbb{P}(x)] + \underbrace{\mathbb{E}_{x \sim \mathbb{Q}} [\text{trace}(\nabla \nabla k(x))]}_{\nabla k(x) \text{ in Stein class of } \mathbb{Q}} \\ &= \Phi + \mathbb{E}_{x \sim \mathbb{Q}} [\nabla k(x)' \nabla \mathbb{P}(x)] + \mathbb{E}_{x \sim \mathbb{Q}} [\text{trace}(\nabla \nabla k(x))] \\ &= \mathbb{E}_{x \sim \mathbb{Q}} \left[\begin{aligned} & \nabla \mathbb{P}(x)' k(x) \nabla \mathbb{P}(x) \\ & + \nabla \mathbb{P}(x)' \nabla k(x) + \nabla k(x)' \nabla \mathbb{P}(x) \\ & + \text{trace}(\nabla \nabla k(x)) \end{aligned} \right] \end{aligned}$$

• RKHS Interpretation

- 1° if $k(x, x')$ is in the Stein class of p , \mathcal{H} is the RKHS induced by k
 then $\forall f \in \mathcal{H}$, f also satisfies $\mathbb{E}_p[\langle \text{Apf}, w \rangle] = 0$

if $f \in \mathcal{H}$, then $\mathbb{E}_p[f] = \langle f, \mathbb{V}k(x, \cdot) \rangle_{\mathcal{H}}$

refer to theorem 11b) in the paper

"Derive reproducing property for kernel method in machine learning"

$\forall f \in \mathcal{H}$

$$\begin{aligned}\mathbb{E}_p[\langle \text{Apf}, w \rangle] &= \mathbb{E}_p[\langle \nabla \log p(w) f(w) + \mathbb{V}f(w) \rangle] \\ &= \mathbb{E}_p[\langle \nabla \log p(w) \langle f, k(x, \cdot) \rangle_{\mathcal{H}} + \langle f, \mathbb{V}k(x, \cdot) \rangle_{\mathcal{H}} \rangle] \\ &= \langle f, \mathbb{E}_p[\nabla \log p(w) k(x, \cdot) + \mathbb{V}k(x, \cdot)] \rangle_{\mathcal{H}} \\ &= \langle f, \mathbb{E}_p[\text{Afk}] \rangle_{\mathcal{H}} \\ &= 0\end{aligned}$$

A more clear notation, let $f = \sum_i \alpha_i k(x_i, \cdot)$

$$\begin{aligned}\mathbb{E}_p[\langle \nabla \log p(w) f(w) + \mathbb{V}f(w) \rangle] &= \mathbb{E}_p[\langle \nabla \log p(w) \sum_i \alpha_i k(x_i, w) + \sum_i \alpha_i \nabla k(x_i, w) \rangle] \\ &= \sum_i \alpha_i \mathbb{E}_p[\langle \nabla \log p(w) k(x_i, w) + \mathbb{V}k(x_i, w) \rangle] \\ &= 0\end{aligned}$$

- 2° let k be a kernel in Stein's class of p , \mathcal{H} be the RKHS induced by k
 let $\beta(x) = \mathbb{E}_p[\langle \text{Afk}, w \rangle]$, then $S_k(p) = \|\beta\|_{\mathcal{H}}^2$

$$\begin{cases} \beta(x) = \mathbb{E}_p[\langle \nabla \log p(w) k(x, w) + \mathbb{V}k(x, w) \rangle] \\ \beta_1(x) = \mathbb{E}_p[\langle \frac{\partial}{\partial x_1} \log p(w) k(x, w) + \frac{\partial}{\partial x_1} k(x, w) \rangle] \end{cases}$$

$$\beta(x) = \mathbb{E}_p[\langle (\nabla \log p(w) - \nabla \log q(w)) k(x, w) \rangle]$$

$$\begin{aligned}S_k(p) &= \mathbb{E}_{x, x'}[(\nabla \log p(x) - \nabla \log q(x))^T k(x, x) (\nabla \log p(x') - \nabla \log q(x'))] \\ &= \mathbb{E}_{x, x'}[(\nabla \log p(x) - \nabla \log q(x))^T \langle k(x, \cdot), k(x', \cdot) \rangle_{\mathcal{H}} (\nabla \log p(x') - \nabla \log q(x'))] \\ &= \sum_{i=1}^d \mathbb{E}_{x, x'}[(\frac{\partial}{\partial x_i} \log p(x) - \frac{\partial}{\partial x_i} \log q(x)) \langle k(x, \cdot), k(x', \cdot) \rangle_{\mathcal{H}} (\frac{\partial}{\partial x'_i} \log p(x') - \frac{\partial}{\partial x'_i} \log q(x'))] \\ &= \sum_{i=1}^d \langle \mathbb{E}_{x, x'}[(\frac{\partial}{\partial x_i} \log p(x) - \frac{\partial}{\partial x_i} \log q(x)) k(x, \cdot)], \mathbb{E}_{x, x'}[(\frac{\partial}{\partial x'_i} \log p(x') - \frac{\partial}{\partial x'_i} \log q(x')) k(x', \cdot)] \rangle_{\mathcal{H}} \\ &= \sum_{i=1}^d \langle \beta_i, \beta_i \rangle_{\mathcal{H}} \\ &= \|\beta\|_{\mathcal{H}}^2\end{aligned}$$

$$3: \langle f, \beta \rangle_{\mathcal{H}^d} = E_{\pi_0}[\text{trace}((\mathcal{A}f)(x))] \quad \forall f \in \mathcal{H}^d$$

$$\begin{aligned} E_{\pi_0}[\text{trace}((\mathcal{A}f)(x))] &= E_{\pi_0}[\text{trace}(\sqrt{q(x)} f(x)^T + \mathcal{A}f(x)^T)] \\ &= \sum_{i=1}^d E_{\pi_0} \left[\frac{\partial}{\partial x_i} \log p(x) f(x) + \frac{\partial f(x)}{\partial x_i} \right] \\ &= \sum_{i=1}^d E_{\pi_0} \left[\frac{\partial}{\partial x_i} \log p(x) \langle f, k(x, \cdot) \rangle_{\mathcal{H}} + \langle f, \frac{\partial}{\partial x_i} k(x, \cdot) \rangle_{\mathcal{H}} \right] \\ &= \sum_{i=1}^d \left\langle f, E_{\pi_0} \left[\frac{\partial}{\partial x_i} \log p(x) k(x, \cdot) + \frac{\partial}{\partial x_i} k(x, \cdot) \right] \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^d \langle f, \beta_i \rangle_{\mathcal{H}} \\ &= \langle f, \beta \rangle_{\mathcal{H}^d} \end{aligned}$$

$$4: B(x) = E_{\pi_0}[(\mathcal{A}f)(x)]$$

$$S_k(g, p) = \|\beta\|_{\mathcal{H}^d}^2$$

$$E_{\pi_0}[(\mathcal{A}f)(x)] = \langle f, \beta \rangle_{\mathcal{H}^d} \quad \forall f \in \mathcal{H}^d$$

$$\begin{aligned} \text{then } JS(g, p) = \|\beta\|_{\mathcal{H}^d}^2 &= \max \{ \langle f, \beta \rangle_{\mathcal{H}^d} : \|f\|_{\mathcal{H}^d} \leq 1 \} \\ &= \max \{ E_{\pi_0}[(\mathcal{A}f)(x)] : \|f\|_{\mathcal{H}^d} \leq 1 \} \end{aligned}$$