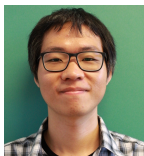


Relative Goodness-of-Fit Tests for Models with Latent Variables

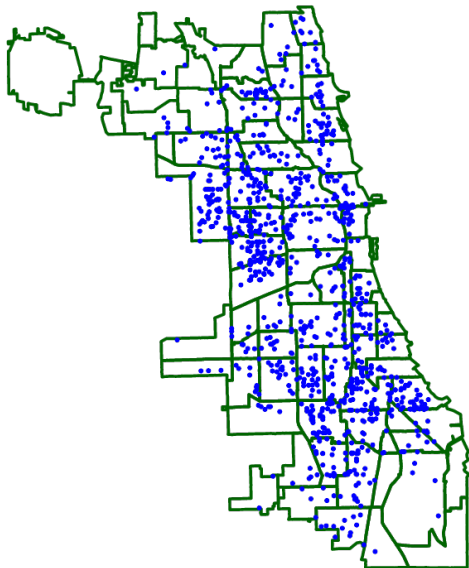
Arthur Gretton



Gatsby Computational Neuroscience Unit,
University College London

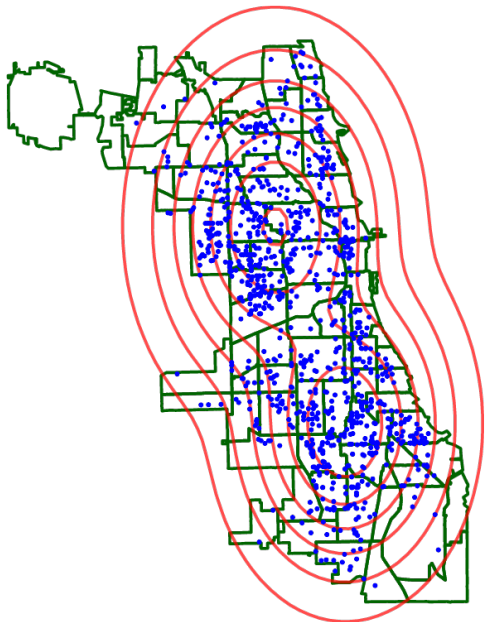
CIRM 2022

Model Criticism



Data = robbery events in
Chicago in 2016.

Model Criticism



Is this a good **model**?

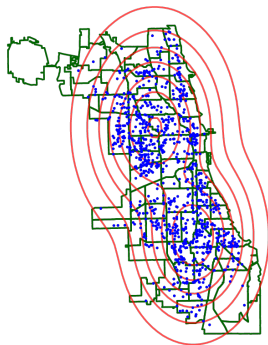
Model Criticism

"All models are wrong."

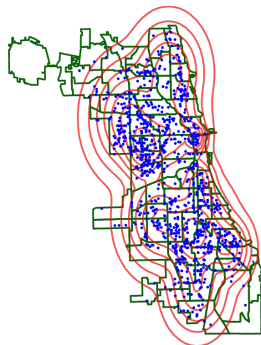
G. Box (1976)

Model comparison

- Have: two candidate models P and Q , and samples $\{x_i\}_{i=1}^n$ from reference distribution R
- Goal: which of P and Q is better?



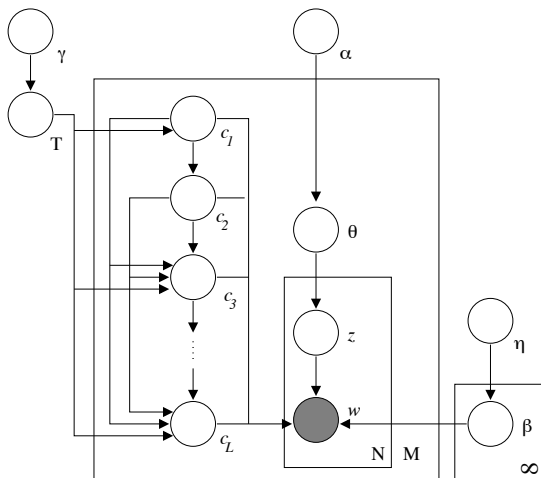
P : two components



Q : ten components

Most interesting models have latent structure

Graphical model representation of hierarchical LDA with a nested CRP prior, Blei et al. (2003)



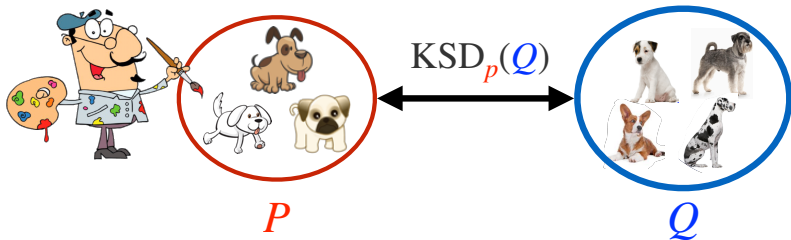
Outline

Relative goodness-of-fit tests for Models with Latent Variables

- The kernel Stein discrepancy
 - Comparing two models via samples: MMD and the witness function.
 - Comparing a sample and a model: Stein modification of the witness class
- Constructing a relative hypothesis test using the KSD
- Relative hypothesis tests with latent variables

Kernel Stein Discrepancy

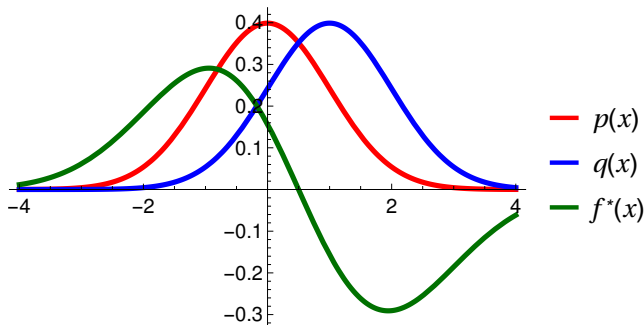
- Model P , data $\{x_i\}_{i=1}^n \sim Q$.
- “All models are wrong” ($P \neq Q$).



MMD: an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$\text{MMD}(P, Q; \mathcal{F}) := \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$



All of kernel methods

Kernels: dot products of features

Feature map $\varphi(x) \in \mathcal{F}$,

$$\varphi(x) = [\dots \varphi_\ell(x) \dots] \in \ell_2$$

For positive definite k ,

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$$

Infinitely many features $\varphi(x)$, dot product in closed form!

Features are solutions to kernel eigenvalue equation

$$\begin{aligned}\lambda_\ell e_\ell(x) &= \int k(x, x') e_\ell(x') d\mathbf{p}(x') dx' \\ \varphi_\ell(x) &= \sqrt{\lambda_\ell} e_\ell(x)\end{aligned}$$

where $\mathbf{p}(x)$ finite Borel measure satisfying Mercer (e.g. supported on \mathcal{X} where \mathcal{X} compact).

All of kernel methods

Kernels: dot products of features

Feature map $\varphi(x) \in \mathcal{F}$,

$$\varphi(x) = [\dots \varphi_\ell(x) \dots] \in \ell_2$$

For positive definite k ,

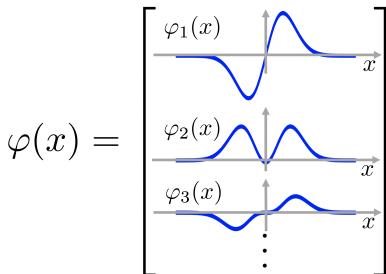
$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$$

Infinitely many features
 $\varphi(x)$, dot product in
closed form!

Exponentiated quadratic kernel

$$k(x, x') = \exp \left(-\gamma \|x - x'\|^2 \right)$$

$$p(x) = \mathcal{N}(0, 1)$$



Features: Gaussian Processes for Machine learning, Rasmussen and Williams, Ch. 4.

All of kernel methods

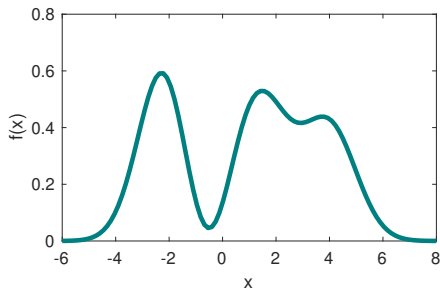
Functions are linear combinations of features:

$$f(x) = \langle f, \varphi(x) \rangle_{\mathcal{F}} = \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \end{bmatrix}^{\top} \begin{bmatrix} \varphi_1(x) \\ \varphi_2(x) \\ \varphi_3(x) \\ \vdots \end{bmatrix}$$
$$\|f\|_{\mathcal{F}}^2 := \sum_{i=1}^{\infty} f_i^2$$

All of kernel methods

“The kernel trick”

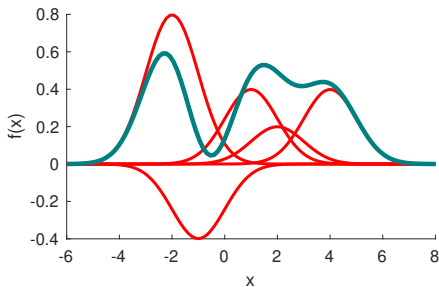
$$\begin{aligned} f(x) &= \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) \\ &= \sum_{i=1}^m \alpha_i \underbrace{k(x_i, x)}_{\langle \varphi(x_i), \varphi(x) \rangle_{\mathcal{F}}} \end{aligned}$$



All of kernel methods

“The kernel trick”

$$\begin{aligned} f(x) &= \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) \\ &= \sum_{i=1}^m \alpha_i \underbrace{k(x_i, x)}_{\langle \varphi(x_i), \varphi(x) \rangle_{\mathcal{F}}} \end{aligned}$$



$$f_l := \sum_{i=1}^m \alpha_i \varphi_l(x_i)$$

Function of infinitely many features expressed using m coefficients.

MMD: an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$\text{MMD}(P, Q; \mathcal{F}) := \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$

For characteristic RKHS \mathcal{F} , $\text{MMD}(P, Q; \mathcal{F}) = 0$ iff $P = Q$

Other choices for witness function class:

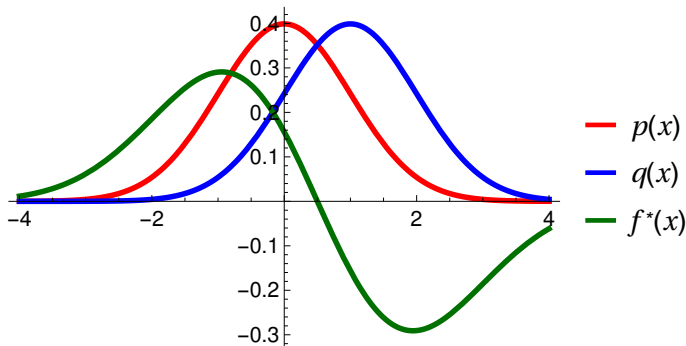
- Bounded continuous [Dudley, 2002]
- Bounded variation 1 (Kolmogorov metric) [Müller, 1997]
- 1-Lipschitz (Wasserstein distances) [Dudley, 2002]

Statistical model criticism: toy example

Can we compute MMD with samples from Q and a model P ?

Problem: usually can't compute $\mathbb{E}_{p f}$ in closed form.

$$\text{MMD}(P, Q) = \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbb{E}_q f - \mathbb{E}_p f]$$



Stein idea

To get rid of $\mathbb{E}_p f$ in

$$\sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbb{E}_q f - \mathbb{E}_p f]$$

we use the (1-D) **Langevin Stein operator**

$$\begin{aligned} [\mathcal{A}_p f](x) &= \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x)) \\ &= f(x) \frac{d}{dx} \log p(x) + \frac{d}{dx} f(x) \end{aligned}$$

Then

$$\mathbb{E}_p \mathcal{A}_p f = 0$$

subject to appropriate boundary conditions.

$$\mathbb{E}_p [\mathcal{A}_p f] = \int \left[\frac{1}{\cancel{p(x)}} \frac{d}{dx} (f(x)p(x)) \right] \cancel{p(x)} dx = [f(x)p(x)]_{-\infty}^{\infty}$$

Kernel Stein Discrepancy

Stein operator

$$\mathcal{A}_p f = f(x) \frac{d}{dx} \log p(x) + \frac{d}{dx} f(x)$$

Kernel Stein Discrepancy (KSD)

$$\text{KSD}_p(Q) = \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_q \mathcal{A}_p g - \mathbb{E}_p \mathcal{A}_p g$$

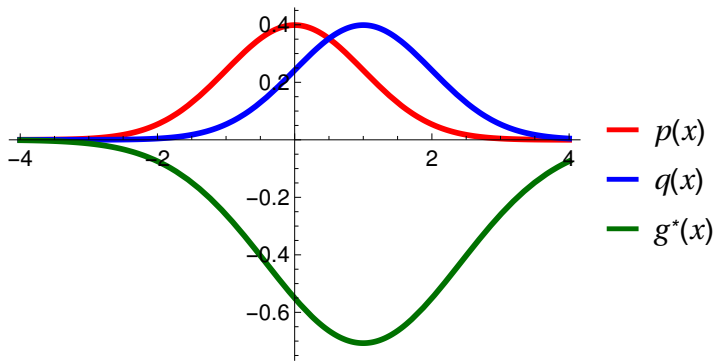
Kernel Stein Discrepancy

Stein operator

$$\mathcal{A}_p f = f(x) \frac{d}{dx} \log p(x) + \frac{d}{dx} f(x)$$

Kernel Stein Discrepancy (KSD)

$$\text{KSD}_p(Q) = \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_q \mathcal{A}_p g - \cancel{\mathbb{E}_p \mathcal{A}_p g} = \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_q \mathcal{A}_p g$$



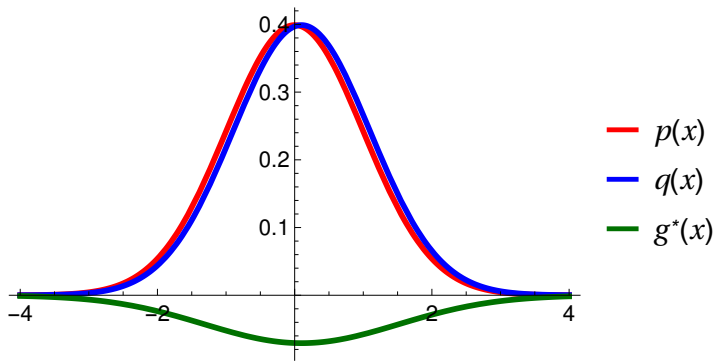
Kernel Stein Discrepancy

Stein operator

$$\mathcal{A}_p f = f(x) \frac{d}{dx} \log p(x) + \frac{d}{dx} f(x)$$

Kernel Stein Discrepancy (KSD)

$$\text{KSD}_p(Q) = \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_q \mathcal{A}_p g - \cancel{\mathbb{E}_p \mathcal{A}_p g} = \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_q \mathcal{A}_p g$$



Computing the kernel Stein discrepancy

How do we get the KSD in closed form (with kernels)?

Can we define “Stein features”?

$$\begin{aligned} [\mathcal{A}_p f](x) &= f(x) \frac{d}{dx} \log p(x) + \frac{d}{dx} f(x) \\ &\stackrel{?}{=} \langle f, \underbrace{\xi(x)}_{\text{stein features}} \rangle_{\mathcal{F}} \end{aligned}$$

where $\mathbb{E}_{x \sim p} \xi(x) = 0$.

Stein RKHS features

Reproducing property for the derivative: for differentiable $k(x, x')$,

$$\frac{d}{dx}f(x) = \left\langle f, \frac{d}{dx}\varphi(x) \right\rangle_{\mathcal{F}} \quad \left\langle \frac{d}{dx}\varphi(x), \varphi(x') \right\rangle_{\mathcal{F}} = \frac{d}{dx}k(x, x')$$

Stein RKHS features

Reproducing property for the derivative: for differentiable $k(x, x')$,

$$\frac{d}{dx}f(x) = \left\langle f, \frac{d}{dx}\varphi(x) \right\rangle_{\mathcal{F}} \quad \left\langle \frac{d}{dx}\varphi(x), \varphi(x') \right\rangle_{\mathcal{F}} = \frac{d}{dx}k(x, x')$$

Using kernel derivative trick in (a),

$$\begin{aligned} [\mathcal{A}_p f](x) &= \left(\frac{d}{dx} \log p(x) \right) f(x) + \frac{d}{dx} f(x) \\ &= \left\langle f, \left(\frac{d}{dx} \log p(x) \right) \varphi(x) + \underbrace{\frac{d}{dx} \varphi(x)}_{(a)} \right\rangle_{\mathcal{F}} \\ &=: \langle f, \xi(x) \rangle_{\mathcal{F}}. \end{aligned}$$

Kernel Stein discrepancy: derivation

Closed-form expression for KSD: given independent $x, x' \sim Q$, then

$$\begin{aligned}\text{KSD}_p(Q) &= \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_{x \sim q}([\mathcal{A}_p g](x)) \\ &= \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_{x \sim q} \langle g, \xi_x \rangle_{\mathcal{F}} \\ &\stackrel{(a)}{=} \sup_{\|g\|_{\mathcal{F}} \leq 1} \langle g, \mathbb{E}_{x \sim q} \xi_x \rangle_{\mathcal{F}} = \|\mathbb{E}_{x \sim q} \xi_x\|_{\mathcal{F}}\end{aligned}$$

Kernel Stein discrepancy: derivation

Closed-form expression for KSD: given independent $x, x' \sim Q$, then

$$\begin{aligned}\text{KSD}_p(Q) &= \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_{x \sim q}([\mathcal{A}_p g](x)) \\ &= \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_{x \sim q} \langle g, \xi_x \rangle_{\mathcal{F}} \\ &\stackrel{(a)}{=} \sup_{\|g\|_{\mathcal{F}} \leq 1} \langle g, \mathbb{E}_{x \sim q} \xi_x \rangle_{\mathcal{F}} = \|\mathbb{E}_{x \sim q} \xi_x\|_{\mathcal{F}}\end{aligned}$$

Kernel Stein discrepancy: derivation

Closed-form expression for KSD: given independent $x, x' \sim Q$, then

$$\begin{aligned}\text{KSD}_p(Q) &= \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_{x \sim q}([\mathcal{A}_p g](x)) \\ &= \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_{x \sim q} \langle g, \xi_x \rangle_{\mathcal{F}} \\ &\stackrel{(a)}{=} \sup_{\|g\|_{\mathcal{F}} \leq 1} \langle g, \mathbb{E}_{x \sim q} \xi_x \rangle_{\mathcal{F}} = \|\mathbb{E}_{x \sim q} \xi_x\|_{\mathcal{F}}\end{aligned}$$

Caution: (a) requires a condition for Bochner integrability,

$$\mathbb{E}_{x \sim q} \left(\frac{d}{dx} \log p(x) \right)^2 < \infty.$$

Kernel Stein discrepancy: derivation

Closed-form expression for KSD: given independent $x, x' \sim Q$, then

$$\begin{aligned}\text{KSD}_p(Q) &= \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_{x \sim q}([\mathcal{A}_p g](x)) \\ &= \sup_{\|g\|_{\mathcal{F}} \leq 1} \mathbb{E}_{x \sim q} \langle g, \xi_x \rangle_{\mathcal{F}} \\ &\stackrel{(a)}{=} \sup_{\|g\|_{\mathcal{F}} \leq 1} \langle g, \mathbb{E}_{x \sim q} \xi_x \rangle_{\mathcal{F}} = \|\mathbb{E}_{x \sim q} \xi_x\|_{\mathcal{F}}\end{aligned}$$

Kernel expression:

$$\begin{aligned}&\|\mathbb{E}_{x \sim q} \xi_x\|_{\mathcal{F}}^2 \\ &= \left\| \mathbb{E}_{x \sim q} \left(\varphi(x) \frac{d}{dx} \log p(x) + \frac{d}{dx} \varphi(x) \right) \right\|_{\mathcal{F}}^2 \\ &= \mathbb{E}_{x, x' \sim Q} \left(k(x, x') \frac{\partial p(x)}{p(x)} \frac{\partial p(x')}{p(x')} + \partial_1 k(x, x') \frac{\partial p(x')}{p(x')} \right. \\ &\quad \left. + \partial_2 k(x, x') \frac{\partial p(x)}{p(x)} + \partial_{12} k(x, x') \right)\end{aligned}$$

Does the Bochner condition matter?

Consider the **standard normal**,

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp(-x^2/2).$$

Then

$$\frac{d}{dx} \log p(x) = -x.$$

If q is a **Cauchy distribution**, then the integral

$$\mathbb{E}_{x \sim q} \left(\frac{d}{dx} \log p(x) \right)^2 = \int_{-\infty}^{\infty} x^2 q(x) dx$$

is undefined.

Does the Bochner condition matter?

Consider the **standard normal**,

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-x^2/2\right).$$

Then

$$\frac{d}{dx} \log p(x) = -x.$$

If q is a **Cauchy distribution**, then the integral

$$\mathbb{E}_{x \sim q} \left(\frac{d}{dx} \log p(x) \right)^2 = \int_{-\infty}^{\infty} x^2 q(x) dx$$

is undefined.

Kernel Stein discrepancy: population expression

Population kernel Stein discrepancy (in \mathbb{R}^D):

$$\text{KSD}_p^2(\mathcal{Q}) = \mathbb{E}_{x, x' \sim \mathcal{Q}} h_p(x, x')$$

where

$$\begin{aligned} h_p(x, x') &= \mathbf{s}_p(x)^\top \mathbf{s}_p(x') k(x, x') + \mathbf{s}_p(x)^\top k_2(x, x') \\ &\quad + \mathbf{s}_p(x')^\top k_1(x, x') + \text{tr}[k_{12}(x, x')] \end{aligned}$$

- $\mathbf{s}_p(x) \in \mathbb{R}^D = \frac{\nabla p(x)}{p(x)}$
- $k_1(a, b) := \nabla_x k(x, x')|_{x=a, x'=b} \in \mathbb{R}^D,$
 $k_2(a, b) := \nabla_{x'} k(x, x')|_{x=a, x'=b} \in \mathbb{R}^D,$
- $k_{12}(a, b) := \nabla_x \nabla_{x'} k(x, x')|_{x=a, x'=b} \in \mathbb{R}^{D \times D}$

Kernel Stein discrepancy: population expression

Population kernel Stein discrepancy (in \mathbb{R}^D):

$$\text{KSD}_{\mathbf{p}}^2(\mathcal{Q}) = \mathbb{E}_{x, x' \sim \mathcal{Q}} h_{\mathbf{p}}(x, x')$$

where

$$\begin{aligned} h_{\mathbf{p}}(x, x') &= \mathbf{s}_{\mathbf{p}}(x)^\top \mathbf{s}_{\mathbf{p}}(x') k(x, x') + \mathbf{s}_{\mathbf{p}}(x)^\top k_2(x, x') \\ &\quad + \mathbf{s}_{\mathbf{p}}(x')^\top k_1(x, x') + \text{tr} [k_{12}(x, x')] \end{aligned}$$

- $\mathbf{s}_{\mathbf{p}}(x) \in \mathbb{R}^D = \frac{\nabla \mathbf{p}(x)}{\mathbf{p}(x)}$
- $k_1(a, b) := \nabla_x k(x, x')|_{x=a, x'=b} \in \mathbb{R}^D$,
 $k_2(a, b) := \nabla_{x'} k(x, x')|_{x=a, x'=b} \in \mathbb{R}^D$,
- $k_{12}(a, b) := \nabla_x \nabla_{x'} k(x, x')|_{x=a, x'=b} \in \mathbb{R}^{D \times D}$

Kernel Stein discrepancy: population expression

Population kernel Stein discrepancy (in \mathbb{R}^D):

$$\text{KSD}_{\mathbf{p}}^2(\mathcal{Q}) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim \mathcal{Q}} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}')$$

where

$$\begin{aligned} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}') &= \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top \mathbf{s}_{\mathbf{p}}(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') + \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top k_2(\mathbf{x}, \mathbf{x}') \\ &\quad + \mathbf{s}_{\mathbf{p}}(\mathbf{x}')^\top k_1(\mathbf{x}, \mathbf{x}') + \text{tr}[k_{12}(\mathbf{x}, \mathbf{x}')] \end{aligned}$$

- $\mathbf{s}_{\mathbf{p}}(\mathbf{x}) \in \mathbb{R}^D = \frac{\nabla \mathbf{p}(\mathbf{x})}{\mathbf{p}(\mathbf{x})}$
- $k_1(a, b) := \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{x}')|_{\mathbf{x}=a, \mathbf{x}'=b} \in \mathbb{R}^D$,
 $k_2(a, b) := \nabla_{\mathbf{x}'} k(\mathbf{x}, \mathbf{x}')|_{\mathbf{x}=a, \mathbf{x}'=b} \in \mathbb{R}^D$,
- $k_{12}(a, b) := \nabla_{\mathbf{x}} \nabla_{\mathbf{x}'} k(\mathbf{x}, \mathbf{x}')|_{\mathbf{x}=a, \mathbf{x}'=b} \in \mathbb{R}^{D \times D}$

Do not need to normalize \mathbf{p} , or sample from it.

Kernel Stein discrepancy: population expression

Population kernel Stein discrepancy (in \mathbb{R}^D):

$$\text{KSD}_{\mathbf{p}}^2(\mathcal{Q}) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim \mathcal{Q}} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}')$$

where

$$\begin{aligned} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}') &= \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top \mathbf{s}_{\mathbf{p}}(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') + \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top k_2(\mathbf{x}, \mathbf{x}') \\ &\quad + \mathbf{s}_{\mathbf{p}}(\mathbf{x}')^\top k_1(\mathbf{x}, \mathbf{x}') + \text{tr}[k_{12}(\mathbf{x}, \mathbf{x}')] \end{aligned}$$

- $\mathbf{s}_{\mathbf{p}}(\mathbf{x}) \in \mathbb{R}^D = \frac{\nabla \mathbf{p}(\mathbf{x})}{\mathbf{p}(\mathbf{x})}$
- $k_1(a, b) := \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{x}')|_{\mathbf{x}=a, \mathbf{x}'=b} \in \mathbb{R}^D$,
 $k_2(a, b) := \nabla_{\mathbf{x}'} k(\mathbf{x}, \mathbf{x}')|_{\mathbf{x}=a, \mathbf{x}'=b} \in \mathbb{R}^D$,
- $k_{12}(a, b) := \nabla_{\mathbf{x}} \nabla_{\mathbf{x}'} k(\mathbf{x}, \mathbf{x}')|_{\mathbf{x}=a, \mathbf{x}'=b} \in \mathbb{R}^{D \times D}$

If kernel is C_0 -universal and \mathcal{Q} satisfies $\mathbb{E}_{\mathbf{x} \sim \mathcal{Q}} \left\| \nabla \left(\log \frac{\mathbf{p}(\mathbf{x})}{\mathbf{q}(\mathbf{x})} \right) \right\|^2 < \infty$,
then $\text{KSD}_{\mathbf{p}}^2(\mathcal{Q}) = 0$ iff $\mathbf{P} = \mathcal{Q}$.

KSD for discrete-valued variables

Discrete domains: $\mathcal{X} = \{1, \dots, L\}^D$ with $L \in \mathbb{N}$.

The population KSD (discrete):

$$\text{KSD}_{\mathbf{p}}^2(\mathcal{Q}) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim \mathcal{Q}} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}')$$

where

$$h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}') = \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top \mathbf{s}_{\mathbf{p}}(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') - \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top k_2(\mathbf{x}, \mathbf{x}') \\ - \mathbf{s}_{\mathbf{p}}(\mathbf{x}')^\top k_1(\mathbf{x}, \mathbf{x}') + \text{tr}[k_{12}(\mathbf{x}, \mathbf{x}')]$$

$$k_1(\mathbf{x}, \mathbf{x}') = \Delta_{\mathbf{x}}^{-1} k(\mathbf{x}, \mathbf{x}'), \Delta_{\mathbf{x}}^{-1} \text{ is difference on } \mathbf{x}, \mathbf{s}_{\mathbf{p}}(\mathbf{x}) = \frac{\Delta \mathbf{p}(\mathbf{x})}{\mathbf{p}(\mathbf{x})}$$

KSD for discrete-valued variables

Discrete domains: $\mathcal{X} = \{1, \dots, L\}^D$ with $L \in \mathbb{N}$.

The population KSD (discrete):

$$\text{KSD}_{\mathbf{p}}^2(\mathcal{Q}) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim \mathcal{Q}} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}')$$

where

$$h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}') = \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top \mathbf{s}_{\mathbf{p}}(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') - \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top k_2(\mathbf{x}, \mathbf{x}') \\ - \mathbf{s}_{\mathbf{p}}(\mathbf{x}')^\top k_1(\mathbf{x}, \mathbf{x}') + \text{tr}[k_{12}(\mathbf{x}, \mathbf{x}')]$$

$$k_1(\mathbf{x}, \mathbf{x}') = \Delta_{\mathbf{x}}^{-1} k(\mathbf{x}, \mathbf{x}'), \Delta_{\mathbf{x}}^{-1} \text{ is difference on } \mathbf{x}, \mathbf{s}_{\mathbf{p}}(\mathbf{x}) = \frac{\Delta_{\mathbf{p}}(\mathbf{x})}{\mathbf{p}(\mathbf{x})}$$

A discrete kernel: $k(\mathbf{x}, \mathbf{x}') = \exp(-d_H(\mathbf{x}, \mathbf{x}'))$, where $d_H(\mathbf{x}, \mathbf{x}') = D^{-1} \sum_{d=1}^D \mathbb{I}(x_d \neq x'_d)$.

KSD for discrete-valued variables

Discrete domains: $\mathcal{X} = \{1, \dots, L\}^D$ with $L \in \mathbb{N}$.

The population KSD (discrete):

$$\text{KSD}_{\mathbf{p}}^2(\mathcal{Q}) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim \mathcal{Q}} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}')$$

where

$$h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}') = \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top \mathbf{s}_{\mathbf{p}}(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') - \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top k_2(\mathbf{x}, \mathbf{x}') \\ - \mathbf{s}_{\mathbf{p}}(\mathbf{x}')^\top k_1(\mathbf{x}, \mathbf{x}') + \text{tr}[k_{12}(\mathbf{x}, \mathbf{x}')]$$

$$k_1(\mathbf{x}, \mathbf{x}') = \Delta_{\mathbf{x}}^{-1} k(\mathbf{x}, \mathbf{x}'), \Delta_{\mathbf{x}}^{-1} \text{ is difference on } \mathbf{x}, \mathbf{s}_{\mathbf{p}}(\mathbf{x}) = \frac{\Delta \mathbf{p}(\mathbf{x})}{\mathbf{p}(\mathbf{x})}$$

A discrete kernel: $k(\mathbf{x}, \mathbf{x}') = \exp(-d_H(\mathbf{x}, \mathbf{x}'))$, where $d_H(\mathbf{x}, \mathbf{x}') = D^{-1} \sum_{d=1}^D \mathbb{I}(x_d \neq x'_d)$.

$\text{KSD}_{\mathbf{p}}^2(\mathcal{Q}) = 0$ iff $\mathbf{P} = \mathcal{Q}$ if

- Gram matrix over all the configurations in \mathcal{X} is strictly positive definite,
- $\mathbf{P} > 0$ and $\mathcal{Q} > 0$.

Empirical statistic, asymptotic normality for $P \neq Q$

The empirical statistic:

$$\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{Q}) := \frac{1}{n(n-1)} \sum_{i \neq j} h_{\textcolor{red}{p}}(\textcolor{blue}{x}_i, \textcolor{blue}{x}_j).$$

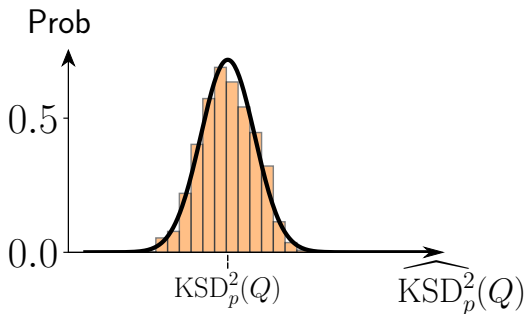
Empirical statistic, asymptotic normality for $P \neq Q$

The empirical statistic:

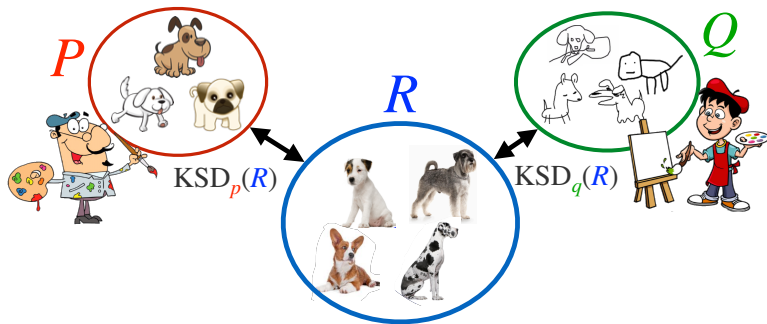
$$\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{Q}) := \frac{1}{n(n-1)} \sum_{i \neq j} h_{\textcolor{red}{p}}(\textcolor{blue}{x}_i, \textcolor{blue}{x}_j).$$

Asymptotic distribution when $\textcolor{red}{P} \neq \textcolor{blue}{Q}$:

$$\sqrt{n} \left(\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{Q}) - \text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{Q}) \right) \xrightarrow{d} \mathcal{N}(0, \sigma_{h_{\textcolor{red}{p}}}^2) \quad \sigma_{h_{\textcolor{red}{p}}}^2 = 4 \text{Var}[\mathbb{E}_{\textcolor{blue}{x}'}[h_{\textcolor{red}{p}}(\textcolor{blue}{x}, \textcolor{blue}{x}')]].$$



Relative goodness-of-fit testing



- Two latent variable models P and Q , data $\{x_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} R$.
- Distinct models $p \neq q$

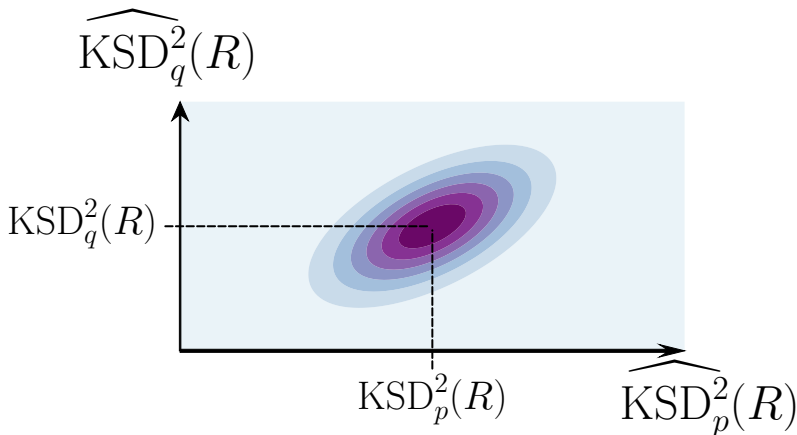
Hypotheses:

$$H_0 : KSD_p(R) \leq KSD_q(R) \text{ vs. } H_1 : KSD_p(R) > KSD_q(R) \\ (H_0 : 'P \text{ is as good as } Q, \text{ or better}' \text{ vs. } H_1 : 'Q \text{ is better}'))$$

Relative GOF testing: joint asymptotic normality

Joint asymptotic normality when $P \neq R$ and $Q \neq R$

$$\sqrt{n} \begin{bmatrix} \widehat{\text{KSD}}_P^2(R) - \text{KSD}_P(R) \\ \widehat{\text{KSD}}_Q^2(R) - \text{KSD}_Q(R) \end{bmatrix} \xrightarrow{d} \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{h_P}^2 & \sigma_{h_P h_Q} \\ \sigma_{h_P h_Q} & \sigma_{h_Q}^2 \end{bmatrix} \right)$$



Relative GOF testing: joint asymptotic normality

Joint asymptotic normality when $P \neq R$ and $Q \neq R$

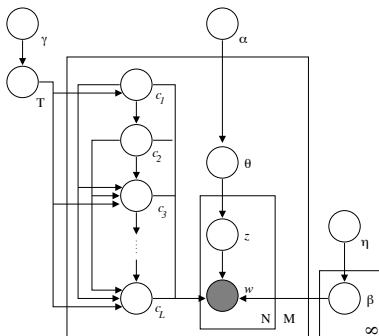
$$\sqrt{n} \begin{bmatrix} \widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{red}{p}}(\textcolor{blue}{R}) \\ \widehat{\text{KSD}}_{\textcolor{teal}{q}}^2(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{teal}{q}}(\textcolor{blue}{R}) \end{bmatrix} \xrightarrow{d} \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{h_{\textcolor{red}{p}}}^2 & \sigma_{h_{\textcolor{red}{p}}h_{\textcolor{teal}{q}}} \\ \sigma_{h_{\textcolor{red}{p}}h_{\textcolor{teal}{q}}} & \sigma_{h_{\textcolor{teal}{q}}}^2 \end{bmatrix} \right)$$

Difference in statistics is asymptotically normal:

$$\begin{aligned} \sqrt{n} \left[\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) - \widehat{\text{KSD}}_{\textcolor{teal}{q}}^2(\textcolor{blue}{R}) - (\text{KSD}_{\textcolor{red}{p}}(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{teal}{q}}(\textcolor{blue}{R})) \right] \\ \xrightarrow{d} \mathcal{N} \left(0, \sigma_{h_{\textcolor{red}{p}}}^2 + \sigma_{h_{\textcolor{teal}{q}}}^2 - 2\sigma_{h_{\textcolor{red}{p}}h_{\textcolor{teal}{q}}} \right) \end{aligned}$$

\implies a statistical test with **null hypothesis** $\text{KSD}_{\textcolor{red}{p}}(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{teal}{q}}(\textcolor{blue}{R}) \leq 0$ is straightforward.

Latent variable models

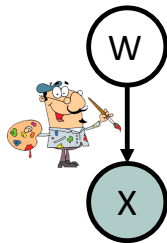
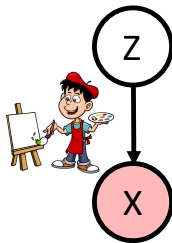


Latent variable models

Can we compare latent variable models with KSD?

$$p(x) = \int p(x|z)p(z)dz$$

$$q(x) = \int q(x|w)p(w)dw$$



Multi-dimensional Stein operator:

$$[T_{\textcolor{red}{p}}f](x) = \left\langle f(x), \underbrace{\frac{\nabla \textcolor{red}{p}(x)}{\textcolor{red}{p}(x)}}_{(\textcolor{brown}{a})} \right\rangle + \langle \nabla, f(x) \rangle.$$

Expression $(\textcolor{brown}{a})$ requires **marginal $p(x)$** , **often intractable**...

What not to do

Approximate the integral using $\{z_j\}_{j=1}^m \sim p(z)$:

$$\begin{aligned} p(x) &= \int p(x|z)p(z)dz \\ &\approx p_m(x) = \frac{1}{m} \sum_{j=1}^m p(x|z_j) \end{aligned}$$

Estimate KSD with approximate density:

$$\widehat{\text{KSD}}_p^2(R) \approx \widehat{\text{KSD}}_{p_m}^2(R)$$

What not to do

Approximate the integral using $\{z_j\}_{j=1}^m \sim p(z)$:

$$\begin{aligned} p(x) &= \int p(x|z)p(z)dz \\ &\approx p_m(x) = \frac{1}{m} \sum_{j=1}^m p(x|z_j) \end{aligned}$$

Estimate KSD with approximate density:

$$\widehat{\text{KSD}}_p^2(R) \approx \widehat{\text{KSD}}_{p_m}^2(R)$$

Problem: $\widehat{\text{KSD}}_{p_m}^2(R)$ asymptotically normal but slow bias decay.

MCMC approximation of score function

Result we use:

$$\mathbf{s}_{\mathbf{p}}(\mathbf{x}) = \mathbb{E}_{z|x}[\mathbf{s}_{\mathbf{p}}(\mathbf{x}|z)]$$

Proof:

$$\begin{aligned}\mathbf{s}_{\mathbf{p}}(\mathbf{x}) &= \frac{\nabla \mathbf{p}(\mathbf{x})}{\mathbf{p}(\mathbf{x})} = \frac{1}{\mathbf{p}(\mathbf{x})} \int \nabla \mathbf{p}(\mathbf{x}|z) d\mathbf{p}(z) \\ &= \int \frac{\nabla \mathbf{p}(\mathbf{x}|z)}{\mathbf{p}(\mathbf{x}|z)} \cdot \frac{\mathbf{p}(\mathbf{x}|z) d\mathbf{p}(z)}{\mathbf{p}(\mathbf{x})} = \mathbb{E}_{z|x}[\mathbf{s}_{\mathbf{p}}(\mathbf{x}|z)],\end{aligned}$$

Friel, N., Mira, A. and Oates, C. J. (2016) Exploiting multi-core architectures for reduced-variance estimation with intractable likelihoods. *Bayesian Analysis*, 11, 215–245.

MCMC approximation of score function

Result we use:

$$\mathbf{s}_{\textcolor{red}{p}}(x) = \mathbb{E}_{z|x}[\mathbf{s}_{\textcolor{red}{p}}(x|z)]$$

Proof:

$$\begin{aligned}\mathbf{s}_{\textcolor{red}{p}}(x) &= \frac{\nabla \textcolor{red}{p}(x)}{\textcolor{red}{p}(x)} = \frac{1}{\textcolor{red}{p}(x)} \int \nabla \textcolor{red}{p}(x|z) dp(z) \\ &= \int \frac{\nabla \textcolor{red}{p}(x|z)}{\textcolor{red}{p}(x|z)} \cdot \frac{\textcolor{red}{p}(x|z) dp(z)}{\textcolor{red}{p}(x)} = \mathbb{E}_{z|x}[\mathbf{s}_{\textcolor{red}{p}}(x|z)],\end{aligned}$$

Approximate intractable posterior $\mathbb{E}_{z|x_i}[\mathbf{s}_{\textcolor{red}{p}}(x_i|z)]$

$$\bar{\mathbf{s}}_{\textcolor{red}{p}}(x_i; z_i^{(t)}) := \frac{1}{m} \sum_{j=1}^m \mathbf{s}_{\textcolor{red}{p}}(x_i|z_{i,j}^{(t)}) \approx \mathbf{s}_{\textcolor{red}{p}}(x_i)$$

with $z_i^{(t)} = (z_{i,1}^{(t)}, \dots, z_{i,m}^{(t)})$ via **MCMC** (after t burn-in steps)

Friel, N., Mira, A. and Oates, C. J. (2016) Exploiting multi-core architectures for reduced-variance estimation with intractable likelihoods. Bayesian Analysis, 11, 215–245.

KSD for latent variable models

Recall earlier KSD estimate:

$$U_n(\textcolor{red}{P}) = \frac{1}{n(n-1)} \sum_{i \neq j} h_{\textcolor{red}{p}}(x_i, x_j) \ (\approx \text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}))$$

KSD for latent variable models

Recall earlier KSD estimate:

$$U_n(\textcolor{red}{P}) = \frac{1}{n(n-1)} \sum_{i \neq j} h_{\textcolor{red}{p}}(x_i, x_j) \ (\approx \text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}))$$

KSD estimate for latent variable models:

$$U_n^{(t)}(\textcolor{red}{P}) := \frac{1}{n(n-1)} \sum_{i \neq j} \bar{H}_{\textcolor{red}{p}}[(x_i, z_i^{(t)}), (x_j, z_j^{(t)})] \ (\approx \text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}))$$

where $\bar{H}_{\textcolor{red}{p}}$ is the Stein kernel $h_{\textcolor{red}{p}}$ with $s_{\textcolor{red}{p}}(x_i)$ replaced with $\bar{s}_{\textcolor{red}{p}}(x_i; z_i^{(t)})$.

Return to relative GOF test, latent variable models

Hypotheses:

$$H_0 : \text{KSD}_{\textcolor{red}{p}}(\textcolor{blue}{R}) \leq \text{KSD}_{\textcolor{teal}{q}}(\textcolor{blue}{R}) \text{ vs. } H_1 : \text{KSD}_{\textcolor{red}{p}}(\textcolor{blue}{R}) > \text{KSD}_{\textcolor{teal}{q}}(\textcolor{blue}{R})$$

(H_0 : ' $\textcolor{red}{P}$ is as good as $\textcolor{teal}{Q}$, or better' vs. H_1 : ' $\textcolor{teal}{Q}$ is better')

Return to relative GOF test, latent variable models

Hypotheses:

$$H_0 : \text{KSD}_{\textcolor{red}{p}}(\textcolor{blue}{R}) \leq \text{KSD}_{\textcolor{teal}{q}}(\textcolor{blue}{R}) \text{ vs. } H_1 : \text{KSD}_{\textcolor{red}{p}}(\textcolor{blue}{R}) > \text{KSD}_{\textcolor{teal}{q}}(\textcolor{blue}{R})$$

(H_0 : ' $\textcolor{red}{P}$ is as good as $\textcolor{teal}{Q}$, or better' vs. H_1 : ' $\textcolor{teal}{Q}$ is better')

Strategy:

- Estimate the difference $\text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{teal}{q}}^2(\textcolor{blue}{R})$ by

$$D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q}) = U_n^{(t)}(\textcolor{red}{P}) - U_n^{(t)}(\textcolor{teal}{Q}).$$

- If $D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q})$ is sufficiently large, reject H_0 .
 - “Sufficient”: control type-I error (falsely rejecting H_0)
 - Requires the (asymptotic) behaviour of $D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q})$

Asymptotic distribution for relative KSD test

Asymptotic distribution of approximate KSD estimate $n, t \rightarrow \infty$:

$$\sqrt{n} \left[D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q}) - \mu_{\textcolor{red}{P}\textcolor{teal}{Q}} \right] \xrightarrow{d} \mathcal{N}(0, \sigma_{\textcolor{red}{P}\textcolor{teal}{Q}}^2)$$

where

$$\begin{aligned} \mu_{\textcolor{red}{P}\textcolor{teal}{Q}} &= \text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{teal}{q}}^2(\textcolor{blue}{R}), \\ \sigma_{\textcolor{red}{P}\textcolor{teal}{Q}}^2 &= \lim_{n, t \rightarrow \infty} n \cdot \text{Var} \left[D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q}) \right]. \end{aligned}$$

Fine print:

- The double limit requires fast bias decay

$$\sqrt{n} [\mathbb{E}\{D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q})\} - \mu_{\textcolor{red}{P}\textcolor{teal}{Q}}] \rightarrow 0$$

- The fourth moment of $\bar{H}_{\textcolor{red}{p}}^{(t)} - \bar{H}_{\textcolor{teal}{q}}^{(t)}$ has finite limit sup. ($t \rightarrow \infty$).

Asymptotic distribution for relative KSD test

Asymptotic distribution of approximate KSD estimate $n, t \rightarrow \infty$:

$$\sqrt{n} \left[D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q}) - \mu_{\textcolor{red}{P}\textcolor{teal}{Q}} \right] \xrightarrow{d} \mathcal{N}(0, \sigma_{\textcolor{red}{P}\textcolor{teal}{Q}}^2)$$

where

$$\begin{aligned} \mu_{\textcolor{red}{P}\textcolor{teal}{Q}} &= \text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{teal}{q}}^2(\textcolor{blue}{R}), \\ \sigma_{\textcolor{red}{P}\textcolor{teal}{Q}}^2 &= \lim_{n, t \rightarrow \infty} n \cdot \text{Var} \left[D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q}) \right]. \end{aligned}$$

Fine print:

- The double limit requires fast bias decay

$$\sqrt{n} [\mathbb{E}\{D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q})\} - \mu_{\textcolor{red}{P}\textcolor{teal}{Q}}] \rightarrow 0$$

- The fourth moment of $\bar{H}_{\textcolor{red}{p}}^{(t)} - \bar{H}_{\textcolor{teal}{q}}^{(t)}$ has finite limit sup. ($t \rightarrow \infty$).

Asymptotic distribution for relative KSD test

Asymptotic distribution of approximate KSD estimate $n, t \rightarrow \infty$:

$$\sqrt{n} \left[D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q}) - \mu_{\textcolor{red}{P}\textcolor{teal}{Q}} \right] \xrightarrow{d} \mathcal{N}(0, \sigma_{\textcolor{red}{P}\textcolor{teal}{Q}}^2)$$

where

$$\mu_{\textcolor{red}{P}\textcolor{teal}{Q}} = \text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{teal}{q}}^2(\textcolor{blue}{R}),$$

$$\sigma_{\textcolor{red}{P}\textcolor{teal}{Q}}^2 = \lim_{n, t \rightarrow \infty} n \cdot \text{Var} \left[D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q}) \right].$$

Level- α test:

$$\text{Reject } H_0 \text{ if } D_n^{(t)}(\textcolor{red}{P}, \textcolor{teal}{Q}) \geq \frac{\hat{\sigma}_{\textcolor{red}{P}\textcolor{teal}{Q}}}{\sqrt{n}} c_{1-\alpha}$$

- $c_{1-\alpha}$ is $(1 - \alpha)$ -quantile of $\mathcal{N}(0, 1)$.
- $\hat{\sigma}_{\textcolor{red}{P}\textcolor{teal}{Q}}$ estimated via jackknife

Experiments

Experiment 1: sensitivity to model difference

- Data R : Probabilistic Principal Component Analysis PPCA(A):

$$x_i \in \mathbb{R}^{100} \sim \mathcal{N}(Az_i, I), \quad z_i \in \mathbb{R}^{10} \sim \mathcal{N}(0, I_z)$$

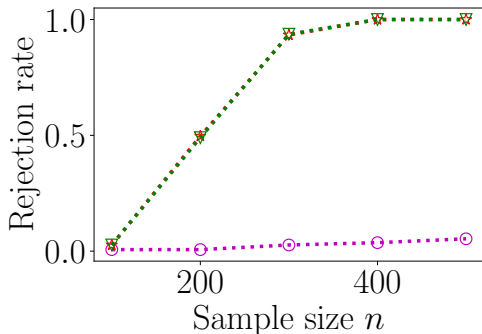
- Generate P , Q : perturb (1, 1)-entry : $A_\delta = A + \delta E_{1,1}$

Experiment 1: sensitivity to model difference

- Data R : Probabilistic Principal Component Analysis PPCA(A):

$$x_i \in \mathbb{R}^{100} \sim \mathcal{N}(Az_i, I), \quad z_i \in \mathbb{R}^{10} \sim \mathcal{N}(0, I_z)$$

- Generate P , Q : perturb (1, 1)-entry : $A_\delta = A + \delta E_{1,1}$



- Alt. H_1 (Q is better):

- P 's perturbation $\delta_P = 2$
- Q 's perturbation $\delta_Q = 1$

- IMQ kernel: $k(x, x') = (1 + \|x - x'\|_2^2 / \sigma_{\text{med}}^2)^{-1/2}$

- NUTS-HMC with sample size $m = 500$ (after $t = 200$ steps).

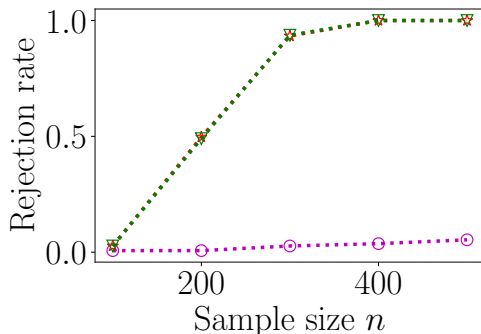
.....○..... MMD ☆..... KSD ▽..... LKSD

Experiment 1: sensitivity to model difference

- Data R : Probabilistic Principal Component Analysis PPCA(A):

$$x_i \in \mathbb{R}^{100} \sim \mathcal{N}(Az_i, I), \quad z_i \in \mathbb{R}^{10} \sim \mathcal{N}(0, I_z)$$

- Generate P , Q : perturb (1, 1)-entry : $A_\delta = A + \delta E_{1,1}$



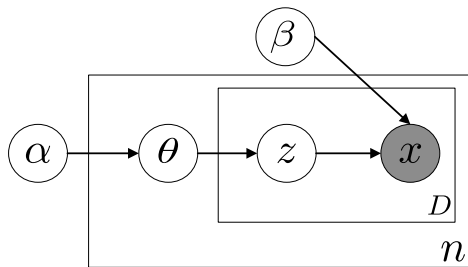
(L)KSD = higher power

- Sample-wise difference in models = subtle (MMD fails)
- Model information is helpful

.....○..... MMD ☆..... KSD ▽..... LKSD

Experiment 2: topic models for arXiv articles

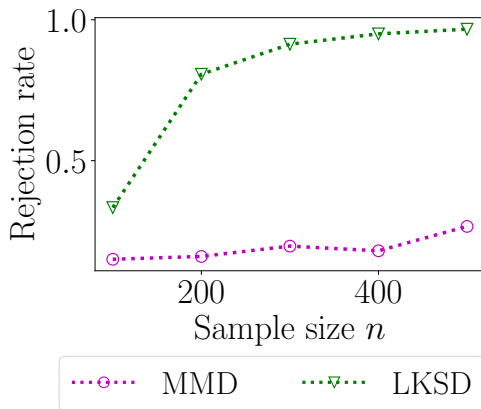
- Data R : arXiv articles from category stat.TH (stat theory) :
- Models P , Q : LDAs trained on articles from different categories
 - P : math.PR (math probability theory)
 - Q : stat.ME (stat methodology)



Graphical model of LDA

Experiment 2: topic models for arXiv articles

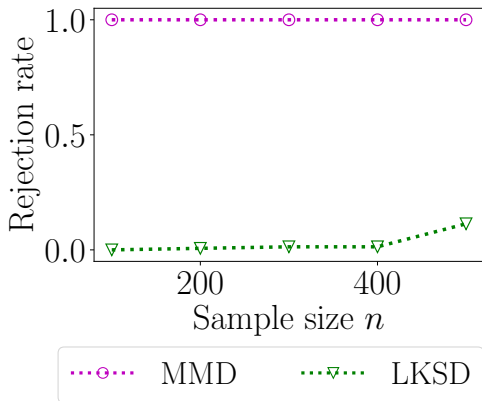
- Data R : arXiv articles from category stat.TH (stat theory):
- Models P , Q : LDAs trained on articles from different categories (100 topics)
 - P : math.PR (math probability theory)
 - Q : stat.ME (stat methodology)



- $\mathcal{X} = \{1, \dots, L\}^D$, $D = 100$, $L = 126, 190$.
- IMQ kernel in BoW rep.:
$$k(x, x') = (1 + \|B(x) - B(x')\|_2^2)^{-1/2}$$
- MCMC size $m = 5000$ (after $t = 500$ steps).

A failure mode

- Data R : arXiv articles from category stat.TH (stat theory) :
- Models P , Q : LDAs trained on articles from different categories (100 topics)
 - P : cs.LG (CS machine learning)
 - Q : stat.ME (stat methodology)



- $\mathcal{X} = \{1, \dots, L\}^D$, $D = 100$, $L = 208,671$.
- IMQ kernel in BoW rep.:
$$k(x, x') = (1 + \|B(x) - B(x')\|_2^2)^{-1/2}$$
- MCMC size $m = 5000$ (after $t = 500$ steps).

What went wrong?

Recall (one-dimension, informally)

$$s_p(x) = \frac{p(x+1)}{p(x)} - 1$$

Numerical instability arises when

- Observed word x has low probability
- Word next to x in vocabulary has non-negligible probability

Zanella-Barker Stein operator

Zanella-Barker Stein operator (1-D):

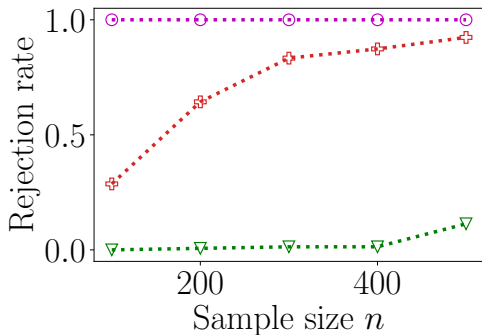
$$\mathcal{A}_p^{\text{ZB}} f(x) = \sum_{\tilde{x} \in \{x+1, x-1\}} \frac{p(\tilde{x})}{p(\tilde{x}) + p(x)} \cdot \{f(\tilde{x}) - f(x)\}$$

- More stable: the ratio $p(\tilde{x})/\{p(\tilde{x}) + p(x)\}$ is always between 0 and 1.
- Similarly applies to latent variable models.

Hodgkinson, Salomone, and Roosta (2020); Shi, Zhou, Hwang, Titsias, and Mackey. (2022)

A resolution to the failure mode

- Data R : arXiv articles from category stat.TH (stat theory) :
- Models P , Q : LDAs trained on articles from different categories (100 topics)
 - P : cs.LG (CS machine learning)
 - Q : stat.ME (stat methodology)



- Improved performance by an alternative Stein operator

.....○..... MMD ▽..... LKSD +..... LKSD (Alt.)

References

A Kernel Test of Goodness of Fit

Kacper Chwialkowski, Heiko Strathmann, Arthur Gretton

<https://arxiv.org/abs/1602.02964>

A Kernel Stein Test for Comparing Latent Variable Models

Heishiro Kanagawa, Wittawat Jitkrittum, Lester Mackey,

Kenji Fukumizu, Arthur Gretton

<https://arxiv.org/abs/1907.00586>

Questions?



Poor MCMC hurts the test

How important is the quality of $\frac{1}{m} \sum_{j=1}^m s_{\textcolor{red}{P}}(x|z_j^{(t)})$?

Experiment with PPCA:

- $\textcolor{red}{P}$: MALA with a bad step size (poor sampler)
- $\textcolor{teal}{Q}$: NUTS-HMC (good sampler)

Expectation:

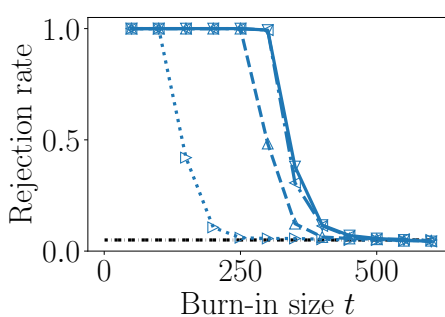
If poor, the test would reject even if $\textcolor{red}{P}$ and $\textcolor{teal}{Q}$ are equally good

Poor MCMC hurts the test

How important is the quality of $\frac{1}{m} \sum_{j=1}^m s_{\textcolor{red}{P}}(x|z_j^{(t)})$?

Experiment with PPCA:

- $\textcolor{red}{P}$: MALA with a bad step size (poor sampler)
- $\textcolor{teal}{Q}$: NUTS-HMC (good sampler)



- Null H_0 (should not reject)
- Significance level $\alpha = 0.05$
- Sample size $n = 100$

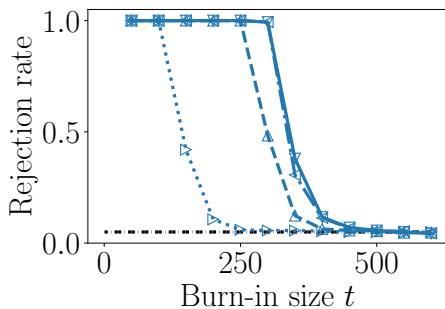
$m = 1$ $m = 10$ $m = 100$ $m = 1000$

Poor MCMC hurts the test

How important is the quality of $\frac{1}{m} \sum_{j=1}^m \mathbf{s}_{\textcolor{red}{P}}(x|z_j^{(t)})$?

Experiment with PPCA:

- $\textcolor{red}{P}$: MALA with a bad step size (poor sampler)
- $\textcolor{teal}{Q}$: NUTS-HMC (good sampler)



- Null H_0 (should not reject)
- Significance level $\alpha = 0.05$
- Sample size $n = 100$

Sufficient burn-in
→ correct type-I error

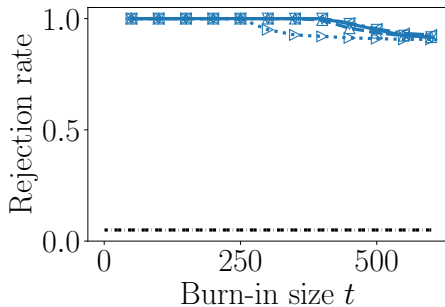
$\textcolor{teal}{\nabla}$ — $m = 1$ $\textcolor{teal}{\triangleleft}$ - - $m = 10$ - $\textcolor{teal}{\triangle}$ - - $m = 100$ $\textcolor{teal}{\triangleright}$... $m = 1000$

Poor MCMC hurts the test

How important is the quality of $\frac{1}{m} \sum_{j=1}^m \mathbf{s}_{\textcolor{red}{P}}(x|z_j^{(t)})$?

Experiment with PPCA:

- $\textcolor{red}{P}$: MALA with a bad step size (poor sampler)
- $\textcolor{teal}{Q}$: NUTS-HMC (good sampler)



- Null H_0 (should not reject)
- Significance level $\alpha = 0.05$
- Sample size $n = 300$

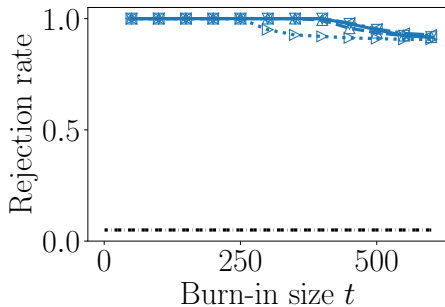
— ∇ — $m = 1$ — \triangleleft — $m = 10$ - - \triangle - - $m = 100$... \triangleright ... $m = 1000$

Poor MCMC hurts the test

How important is the quality of $\frac{1}{m} \sum_{j=1}^m s_{\textcolor{red}{P}}(x|z_j^{(t)})$?

Experiment with PPCA:

- $\textcolor{red}{P}$: MALA with a bad step size (poor sampler)
- $\textcolor{teal}{Q}$: NUTS-HMC (good sampler)



- Null H_0 (should not reject)
- Significance level $\alpha = 0.05$
- Sample size $n = 300$

Larger $n \implies$ more
sensitive to mismatch

— ∇ — $m = 1$ - \triangleleft - $m = 10$ - \triangle - $m = 100$ $\cdots\triangleright\cdots$ $m = 1000$