

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

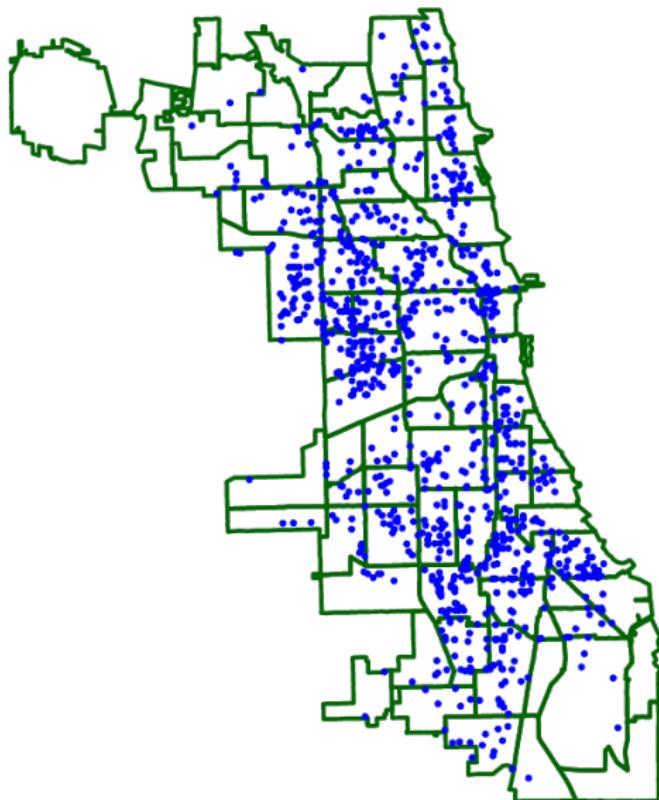
Arthur Gretton



Gatsby Computational Neuroscience Unit,  
University College London

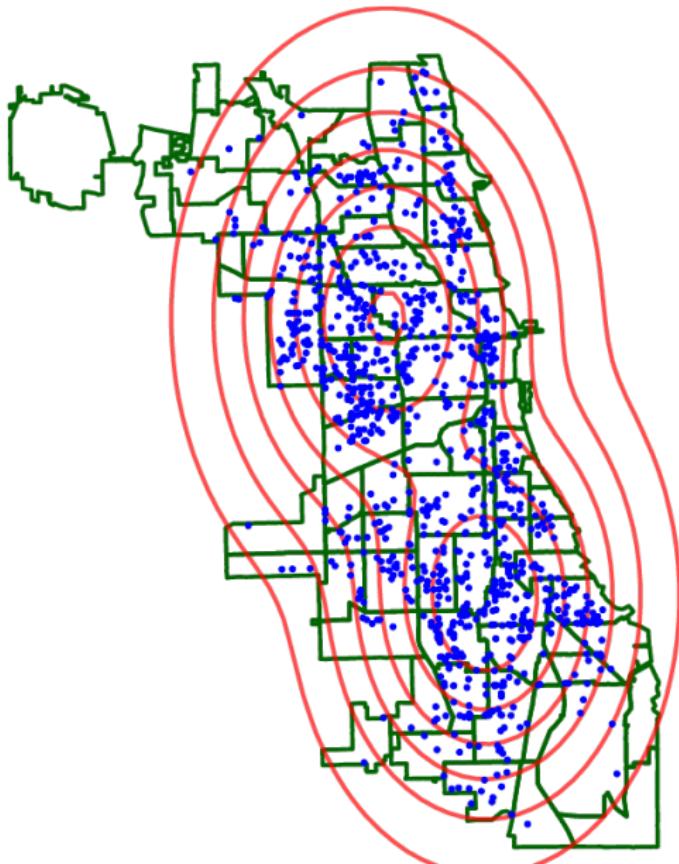
Department of Statistics, Harvard, 2023

## Model Criticism



Data = robbery events in Chicago in 2016.

## Model Criticism



Is this a good **model**?

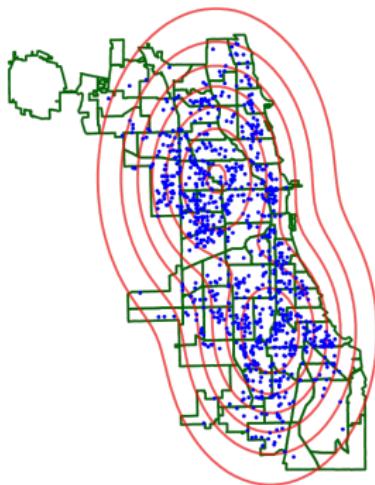
## Model Criticism

"All models are wrong."

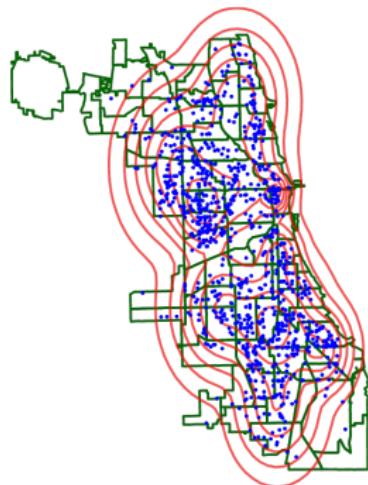
G. Box (1976)

## 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?



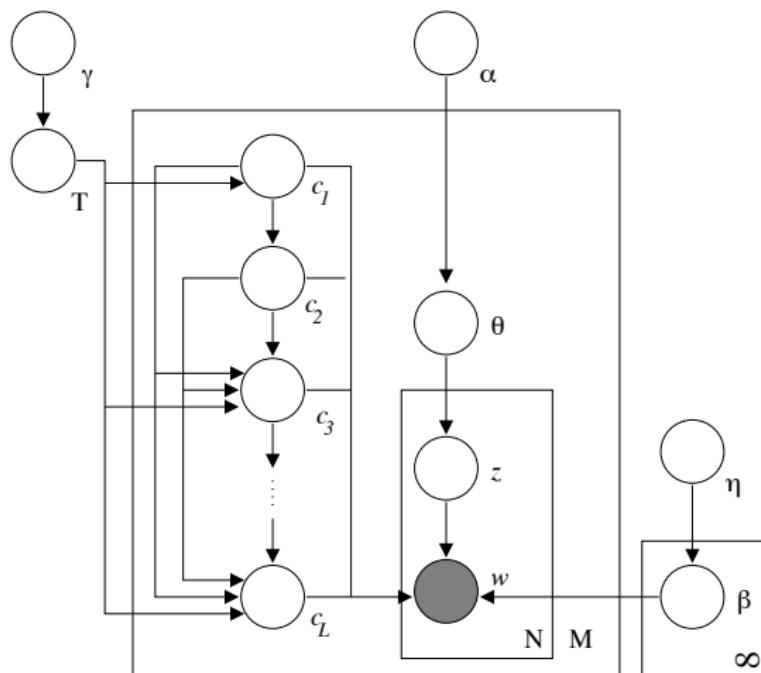
$P$  : two components



$Q$  : ten components

## Most interesting models have latent structure

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



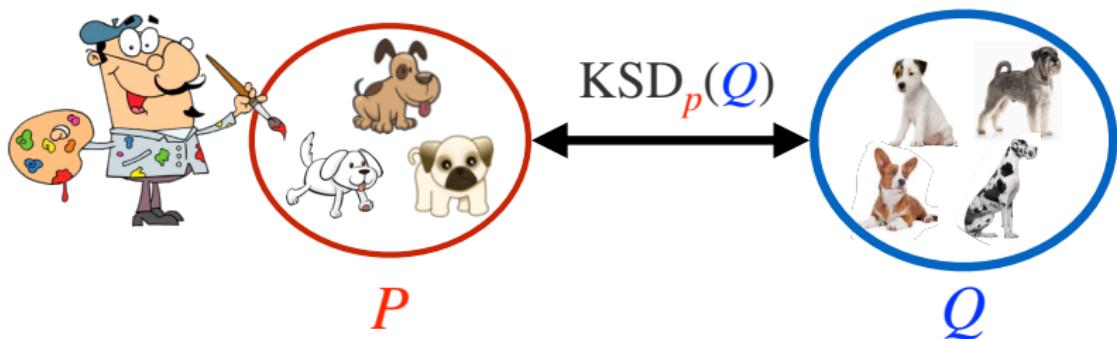
# Outline

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

- The Maximum Mean Discrepancy: an integral probability metric
  - maximize difference in expectations using an RKHS witness class
- The kernel Stein discrepancy
  - 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

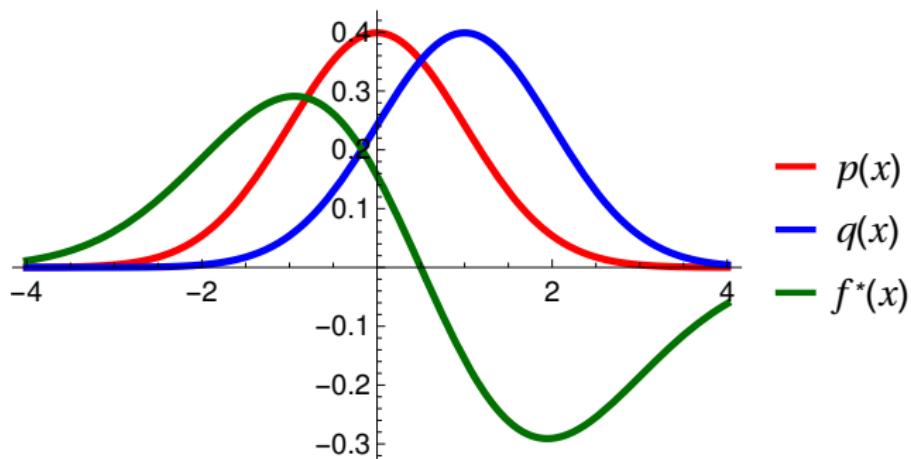
- Model  $P$ , data  $\{x_i\}_{i=1}^n \sim Q$ .
- “All models are wrong” ( $P \neq Q$ ).



## MMD: an integral probability metric

Maximum mean discrepancy: for  $P$  vs  $Q$  and "smooth" (Reproducing Kernel Hilbert Space) witness class  $\mathcal{F}$

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



## All of kernel methods

Kernels: dot products of features

Features are solutions to kernel eigenvalue equation

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

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

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

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!

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

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

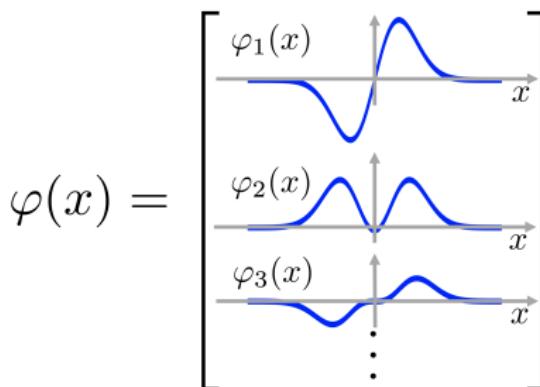
For positive definite  $k$ ,

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

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

Exponentiated quadratic kernel

$$k(x, x') = \exp(-\gamma \|x - x'\|^2)$$
$$p(x) = \mathcal{N}(0, 1)$$



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

## All of kernel methods

Functions are linear combinations of features:

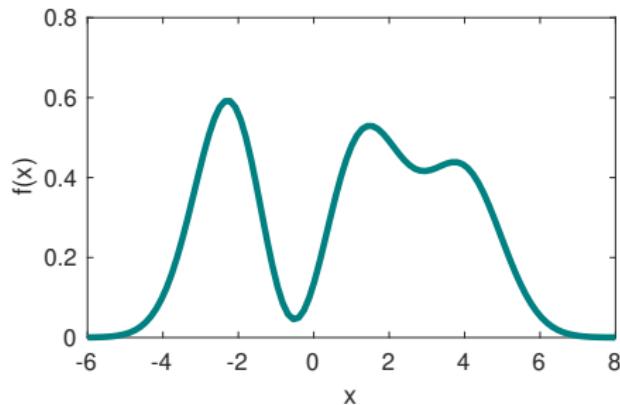
$$f(x) = \langle \mathbf{f}, \varphi(x) \rangle_{\mathcal{F}} = \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \end{bmatrix}^T \begin{bmatrix} \varphi_1(x) \\ \varphi_2(x) \\ \varphi_3(x) \\ \vdots \end{bmatrix}$$

$\|\mathbf{f}\|_{\mathcal{F}}^2 := \sum_{i=1}^{\infty} f_i^2$

# All of kernel methods

“The kernel trick”

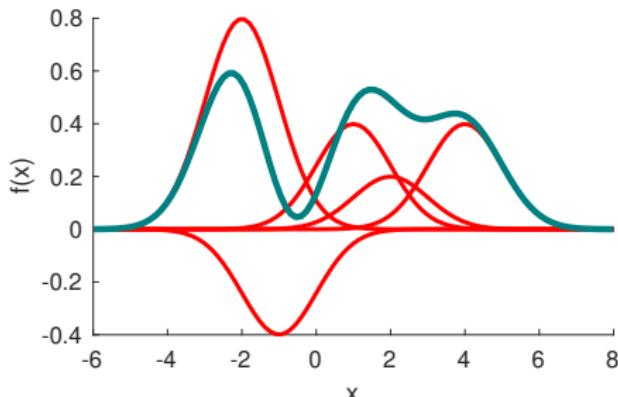
$$\begin{aligned} f(x) &= \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) \\ &= \sum_{i=1}^m \alpha_i \underbrace{k(x_i, x)}_{\langle \varphi(x_i), \varphi(x) \rangle_{\mathcal{F}}} \end{aligned}$$



# All of kernel methods

“The kernel trick”

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



$$f_{\ell} := \sum_{i=1}^m \alpha_i \varphi_{\ell}(x_i)$$

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  $P = Q$

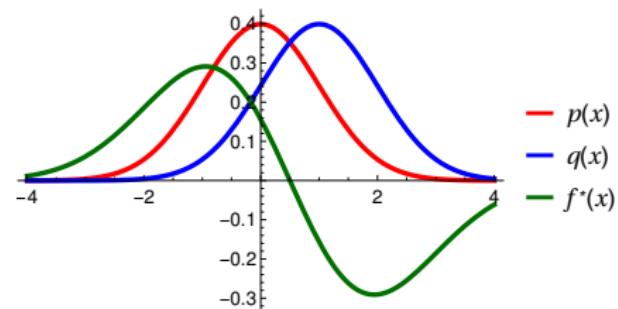
Other choices for witness function class:

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

## The MMD in closed form

The MMD:

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



## The MMD in closed form

The MMD:

$$MMD(\mathcal{P}, \mathcal{Q}; \mathcal{F})$$

$$= \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} [\mathbb{E}_{\mathcal{P}} \mathbf{f}(\mathcal{X}) - \mathbb{E}_{\mathcal{Q}} \mathbf{f}(\mathcal{Y})]$$

$$= \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} \langle \mathbf{f}, \boldsymbol{\mu}_{\mathcal{P}} - \boldsymbol{\mu}_{\mathcal{Q}} \rangle_{\mathcal{F}}$$

use

$$\begin{aligned} \mathbb{E}_{\mathcal{P}} \mathbf{f}(\mathcal{X}) &= \mathbb{E}_{\mathcal{P}} \langle \varphi(\mathcal{X}), \mathbf{f} \rangle_{\mathcal{F}} \\ &= \langle \mathbb{E}_{\mathcal{P}} [\varphi(\mathcal{X})], \mathbf{f} \rangle_{\mathcal{F}} \\ &= \langle \boldsymbol{\mu}_{\mathcal{P}}, \mathbf{f} \rangle_{\mathcal{F}} \end{aligned}$$

## The MMD in closed form

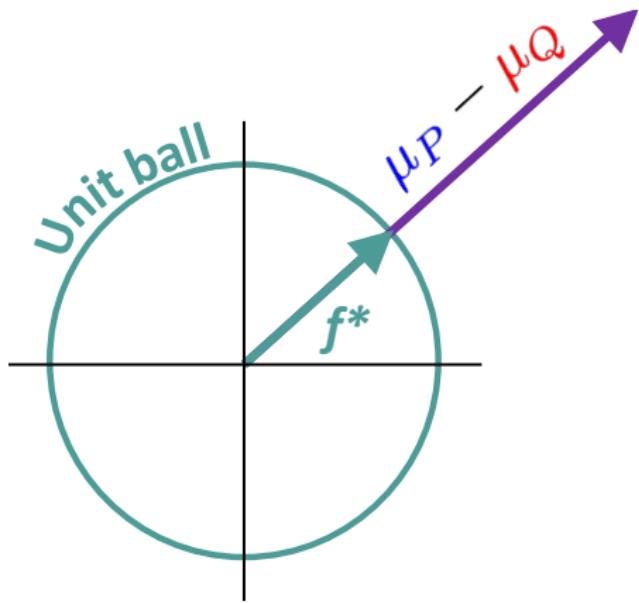
The MMD:

$$MMD(P, Q; \mathcal{F})$$

$$= \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$

$$= \sup_{\|f\|_{\mathcal{F}} \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$

$$= \|\mu_P - \mu_Q\|$$



$$f^* = \frac{\mu_P - \mu_Q}{\|\mu_P - \mu_Q\|}$$

## The MMD in closed form

The MMD:

$$\begin{aligned}MMD(\mathcal{P}, \mathcal{Q}; \mathcal{F}) &= \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} [\mathbb{E}_{\mathcal{P}} f(\mathbf{X}) - \mathbb{E}_{\mathcal{Q}} f(\mathbf{Y})] \\&= \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} \langle \mathbf{f}, \boldsymbol{\mu}_{\mathcal{P}} - \boldsymbol{\mu}_{\mathcal{Q}} \rangle_{\mathcal{F}} \\&= \|\boldsymbol{\mu}_{\mathcal{P}} - \boldsymbol{\mu}_{\mathcal{Q}}\|\end{aligned}$$

In terms of kernels:

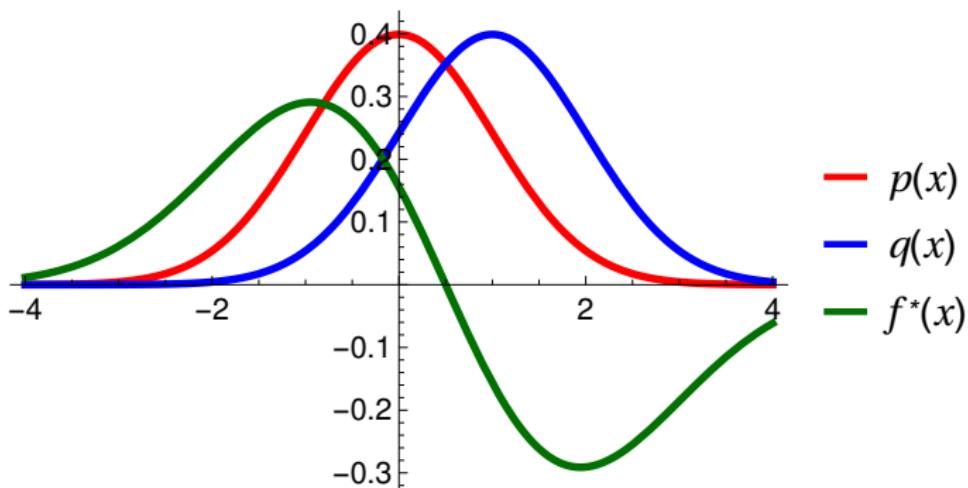
$$\begin{aligned}MMD^2(\mathcal{P}, \mathcal{Q}) &= \|\boldsymbol{\mu}_{\mathcal{P}} - \boldsymbol{\mu}_{\mathcal{Q}}\|_{\mathcal{F}}^2 \\&= \underbrace{\mathbb{E}_{\mathcal{P}} k(\mathbf{x}, \mathbf{x}')}_{(a)} + \underbrace{\mathbb{E}_{\mathcal{Q}} k(\mathbf{y}, \mathbf{y}')}_{(a)} - 2 \underbrace{\mathbb{E}_{\mathcal{P}, \mathcal{Q}} k(\mathbf{x}, \mathbf{y})}_{(b)}\end{aligned}$$

## Comparing a sample and model

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

Problem: usually can't compute  $E_{p,f}$  in closed form.

$$\text{MMD}(P, Q) = \sup_{\|f\|_{\mathcal{F}} \leq 1} [E_q f - E_p f]$$



## Stein idea

To get rid of  $\mathbb{E}_{\textcolor{red}{p}} f$  in

$$\sup_{\|\textcolor{teal}{f}\|_{\mathcal{F}} \leq 1} [\mathbb{E}_{\textcolor{blue}{q}} \textcolor{teal}{f} - \mathbb{E}_{\textcolor{red}{p}} f]$$

we use the (1-D) Langevin Stein operator

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

Then

$$\mathbb{E}_{\textcolor{red}{p}} \mathcal{A}_{\textcolor{red}{p}} f = 0$$

subject to appropriate boundary conditions.

Gorham and Mackey (NeurIPS 15), Oates, Girolami, Chopin (JRSS B 2016)

## Stein idea

To get rid of  $E_{\textcolor{red}{p}} \textcolor{teal}{f}$  in

$$\sup_{\|\textcolor{teal}{f}\|_{\mathcal{F}} \leq 1} [E_q f - E_{\textcolor{red}{p}} \textcolor{teal}{f}]$$

we use the (1-D) Langevin Stein operator

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

Then

$$E_{\textcolor{red}{p}} \mathcal{A}_{\textcolor{red}{p}} \textcolor{teal}{f} = 0$$

subject to appropriate boundary conditions.

$$E_{\textcolor{red}{p}} [\mathcal{A}_{\textcolor{red}{p}} \textcolor{teal}{f}] = \int \left[ \frac{1}{\textcolor{red}{p}(x)} \frac{d}{dx} (\textcolor{teal}{f}(x) \textcolor{red}{p}(x)) \right] \textcolor{red}{p}(x) dx = [\textcolor{teal}{f}(x) \textcolor{red}{p}(x)]_{-\infty}^{\infty}$$

## Stein idea

To get rid of  $E_{\textcolor{red}{p}} f$  in

$$\sup_{\|f\|_{\mathcal{F}} \leq 1} [E_q f - E_{\textcolor{red}{p}} f]$$

we use the (1-D) Langevin Stein operator

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

Then

$$E_{\textcolor{red}{p}} \mathcal{A}_{\textcolor{red}{p}} f = 0$$

subject to appropriate boundary conditions.

Do not need to normalize  $p$ , or sample from it.

Gorham and Mackey (NeurIPS 15), Oates, Girolami, Chopin (JRSS B 2016)

## Kernel Stein Discrepancy

Stein operator

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

Kernel Stein Discrepancy (KSD)

$$\text{KSD}_{\textcolor{red}{p}}(\textcolor{blue}{Q}) = \sup_{\|\textcolor{teal}{g}\|_{\mathcal{F}} \leq 1} \mathbb{E}_{\textcolor{blue}{q}} \mathcal{A}_{\textcolor{red}{p}} \textcolor{teal}{g} - \mathbb{E}_{\textcolor{red}{p}} \mathcal{A}_{\textcolor{red}{p}} \textcolor{teal}{g}$$

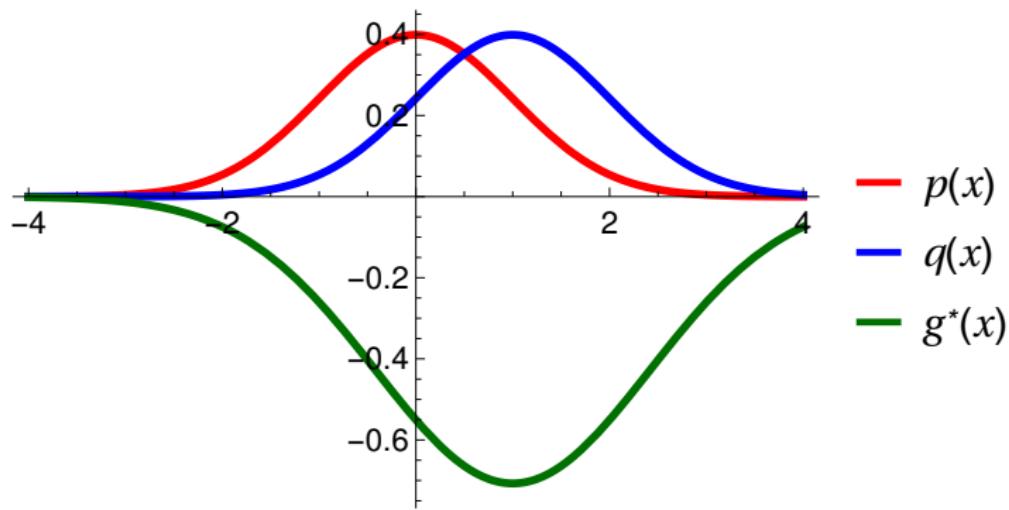
# Kernel Stein Discrepancy

Stein operator

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

Kernel Stein Discrepancy (KSD)

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



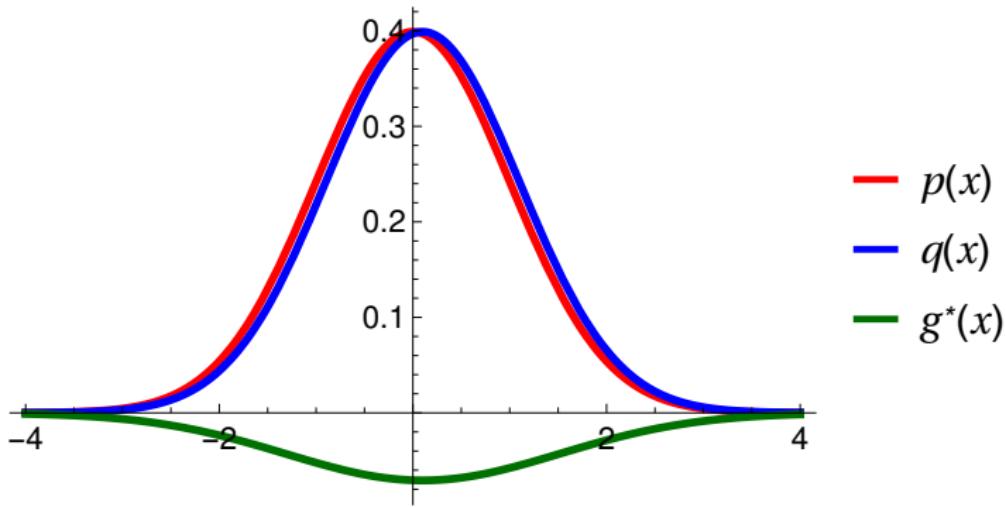
# Kernel Stein Discrepancy

Stein operator

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

Kernel Stein Discrepancy (KSD)

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



## Computing the kernel Stein discrepancy

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

Can we define “Stein features”?

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

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

## 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  $(a)$ ,

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

## Kernel Stein discrepancy: derivation

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

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

## Kernel Stein discrepancy: derivation

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

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

## Kernel Stein discrepancy: derivation

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

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

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

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

## Kernel Stein discrepancy: derivation

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

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

Kernel expression:

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

## Does the Bochner condition matter?

Consider the standard normal,

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

Then

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

If  $\textcolor{blue}{q}$  is a Cauchy distribution, then the integral

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

is undefined.

## Does the Bochner condition matter?

Consider the standard normal,

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

Then

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

If  $\textcolor{blue}{q}$  is a Cauchy distribution, then the integral

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

is undefined.

## Kernel Stein discrepancy: population expression

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

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

where

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

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

## Kernel Stein discrepancy: population expression

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

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

where

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

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

## Kernel Stein discrepancy: population expression

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

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

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

## 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

$$\begin{aligned} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}') &= \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top \mathbf{s}_{\mathbf{p}}(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') - \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top k_2(\mathbf{x}, \mathbf{x}') \\ &\quad - \mathbf{s}_{\mathbf{p}}(\mathbf{x}')^\top \mathbf{k}_1(\mathbf{x}, \mathbf{x}') + \text{tr}[k_{12}(\mathbf{x}, \mathbf{x}')] \end{aligned}$$

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

## 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

$$\begin{aligned} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}') &= \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top \mathbf{s}_{\mathbf{p}}(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') - \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top k_2(\mathbf{x}, \mathbf{x}') \\ &\quad - \mathbf{s}_{\mathbf{p}}(\mathbf{x}')^\top \mathbf{k}_1(\mathbf{x}, \mathbf{x}') + \text{tr}[k_{12}(\mathbf{x}, \mathbf{x}')] \end{aligned}$$

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

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):

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

where

$$\begin{aligned} h_{\mathbf{p}}(\mathbf{x}, \mathbf{x}') &= \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top \mathbf{s}_{\mathbf{p}}(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') - \mathbf{s}_{\mathbf{p}}(\mathbf{x})^\top k_2(\mathbf{x}, \mathbf{x}') \\ &\quad - \mathbf{s}_{\mathbf{p}}(\mathbf{x}')^\top \mathbf{k}_1(\mathbf{x}, \mathbf{x}') + \text{tr}[k_{12}(\mathbf{x}, \mathbf{x}')] \end{aligned}$$

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

A discrete kernel:  $k(\mathbf{x}, \mathbf{x}') = \exp(-d_H(\mathbf{x}, \mathbf{x}'))$ , where

$$d_H(\mathbf{x}, \mathbf{x}') = D^{-1} \sum_{d=1}^D \mathbb{I}(x_d \neq x'_d).$$

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

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

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

The empirical statistic:

$$\widehat{\text{KSD}}_p^2(Q) := \frac{1}{n(n-1)} \sum_{i \neq j} h_p(x_i, x_j).$$

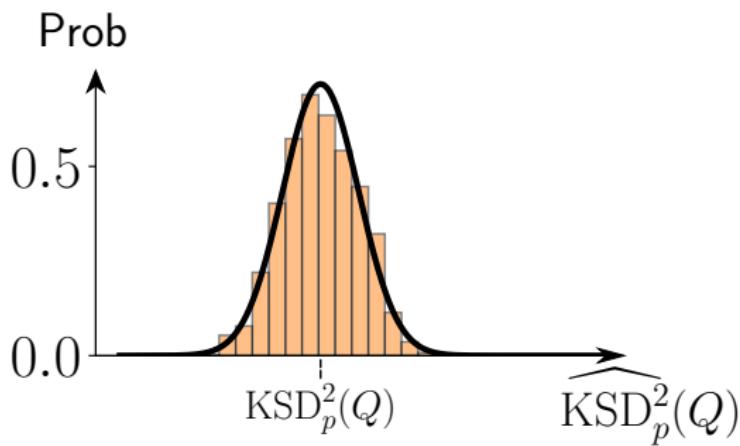
## Empirical statistic, asymptotic normality for $P \neq Q$

The empirical statistic:

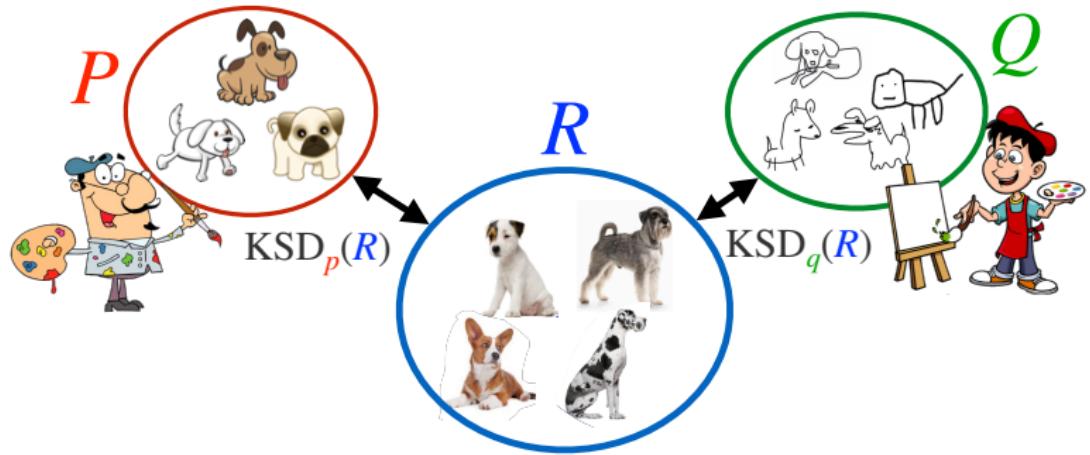
$$\widehat{\text{KSD}}_p^2(Q) := \frac{1}{n(n-1)} \sum_{i \neq j} h_p(x_i, x_j).$$

Asymptotic distribution when  $P \neq Q$ :

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



## Relative goodness-of-fit testing



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

Hypotheses:

$$H_0 : \text{KSD}_p(R) \leq \text{KSD}_q(R) \text{ vs. } H_1 : \text{KSD}_p(R) > \text{KSD}_q(R)$$

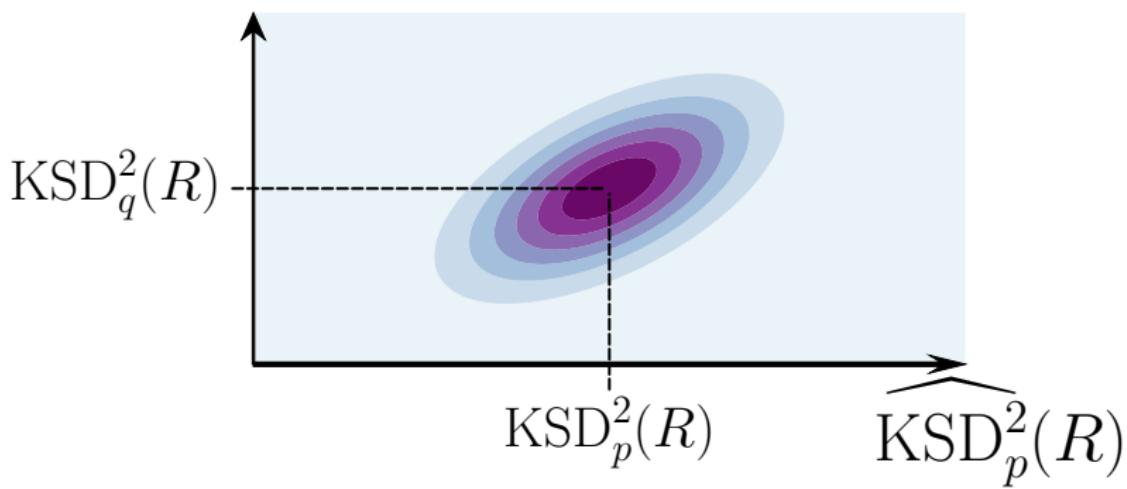
( $H_0$  : ' $P$  is as good as  $Q$ , or better' vs.  $H_1$  : ' $Q$  is better' )

## Relative GOF testing: joint asymptotic normality

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

$$\sqrt{n} \begin{bmatrix} \widehat{\text{KSD}}_p^2(R) - \text{KSD}_p(R) \\ \widehat{\text{KSD}}_q^2(R) - \text{KSD}_q(R) \end{bmatrix} \xrightarrow{d} \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{h_p}^2 & \sigma_{h_p h_q} \\ \sigma_{h_p h_q} & \sigma_{h_q}^2 \end{bmatrix} \right)$$

$$\widehat{\text{KSD}}_q^2(R)$$



## Relative GOF testing: joint asymptotic normality

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

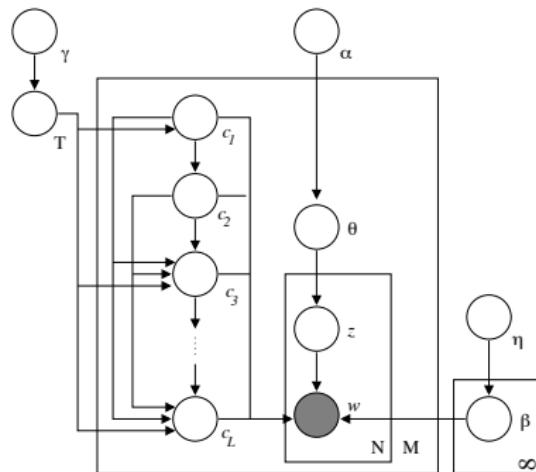
$$\sqrt{n} \begin{bmatrix} \widehat{\text{KSD}}_p^2(R) - \text{KSD}_p(R) \\ \widehat{\text{KSD}}_q^2(R) - \text{KSD}_q(R) \end{bmatrix} \xrightarrow{d} \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{h_p}^2 & \sigma_{h_p h_q} \\ \sigma_{h_p h_q} & \sigma_{h_q}^2 \end{bmatrix} \right)$$

**Difference** in statistics is asymptotically normal:

$$\begin{aligned} & \sqrt{n} \left[ \widehat{\text{KSD}}_p^2(R) - \widehat{\text{KSD}}_q^2(R) - (\text{KSD}_p(R) - \text{KSD}_q(R)) \right] \\ & \xrightarrow{d} \mathcal{N} \left( 0, \sigma_{h_p}^2 + \sigma_{h_q}^2 - 2\sigma_{h_p h_q} \right) \end{aligned}$$

$\implies$  a statistical test with **null hypothesis**  $\text{KSD}_p(R) - \text{KSD}_q(R) \leq 0$  is straightforward.

# Latent variable models

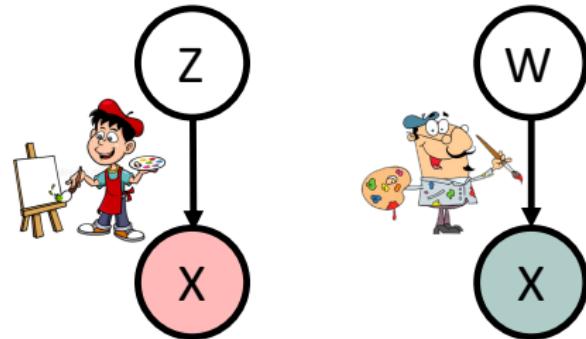


## Latent variable models

Can we compare latent variable models with KSD?

$$\textcolor{red}{p}(x) = \int \textcolor{red}{p}(x|z)p(z)dz$$

$$\textcolor{teal}{q}(x) = \int \textcolor{teal}{q}(x|w)p(w)dw$$



Multi-dimensional Stein operator:

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

Expression  $(a)$  requires marginal  $p(x)$ , often intractable...

## What not to do

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

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

Estimate KSD with approximate density:

$$\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) \approx \widehat{\text{KSD}}_{\textcolor{red}{p}_m}^2(\textcolor{blue}{R})$$

## What not to do

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

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

Estimate KSD with approximate density:

$$\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) \approx \widehat{\text{KSD}}_{\textcolor{red}{p}_m}^2(\textcolor{blue}{R})$$

Problem:  $\widehat{\text{KSD}}_{\textcolor{red}{p}_m}^2(\textcolor{blue}{R})$  asymptotically normal but slow bias decay.

## MCMC approximation of score function

Result we use:

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

Proof:

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

## MCMC approximation of score function

Result we use:

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

Proof:

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

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

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

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:

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

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

## Return to relative GOF test, latent variable models

Hypotheses:

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

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

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

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

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

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- $\alpha$  test:

Reject  $H_0$  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  $\textcolor{blue}{R}$  : Probabilistic Principal Component Analysis PPCA( $A$ ):

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

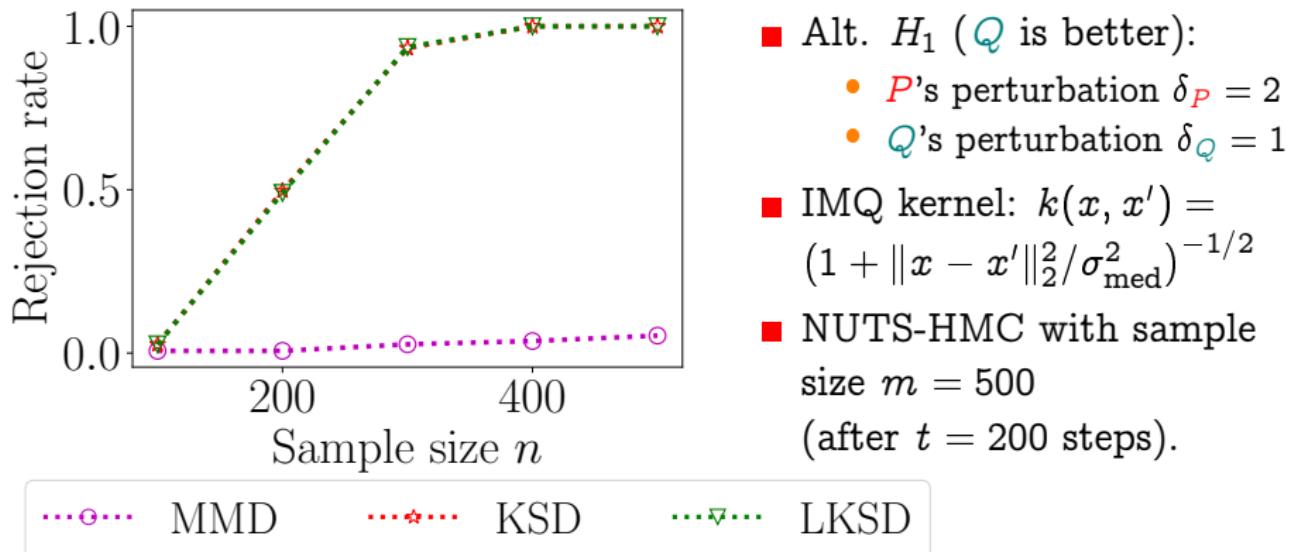
- Generate  $\textcolor{red}{P}$ ,  $\textcolor{teal}{Q}$  : perturb (1, 1)-entry :  $A_\delta = A + \delta E_{1,1}$

## Experiment 1: sensitivity to model difference

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

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

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

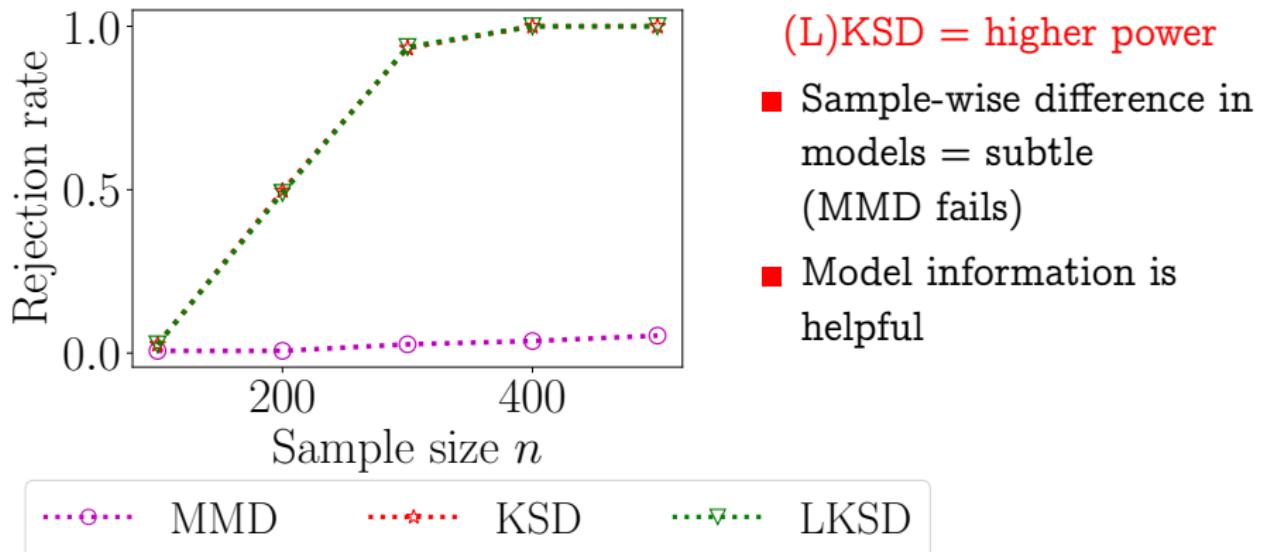


## Experiment 1: sensitivity to model difference

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

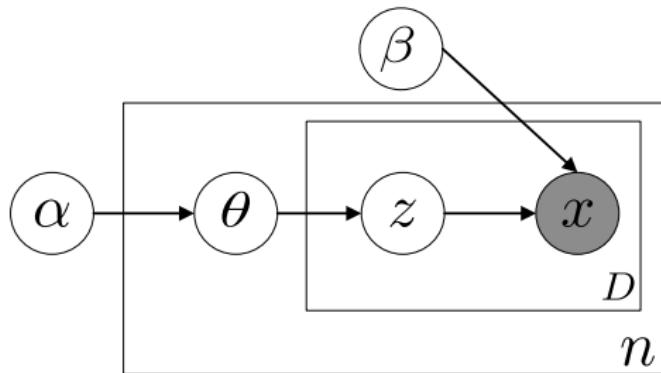
$$x_i \in \mathbb{R}^{100} \sim \mathcal{N}(Az_i, I), z_i \in \mathbb{R}^{10} \sim \mathcal{N}(0, I_z)$$

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



## Experiment 2: topic models for arXiv articles

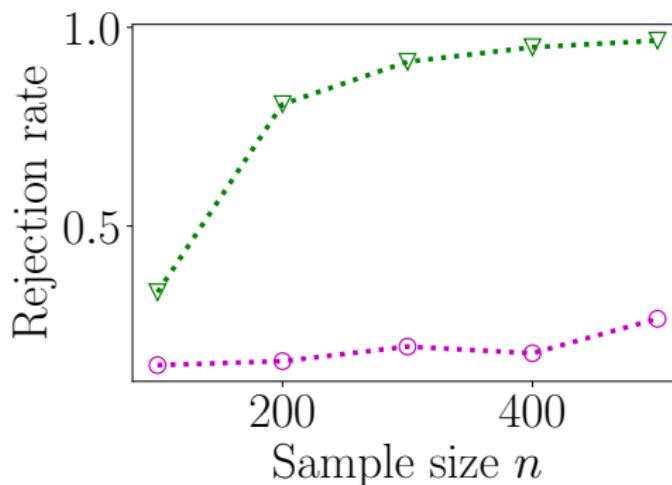
- Data  $R$  : arXiv articles from category stat.TH (stat theory) :
- Models  $P, Q$  : LDAs trained on articles from different categories
  - $P$  : math.PR (math probability theory)
  - $Q$  : stat.ME (stat methodology).  $H_1$ :  $Q$  is better



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).  $H_1$ :  $Q$  is better

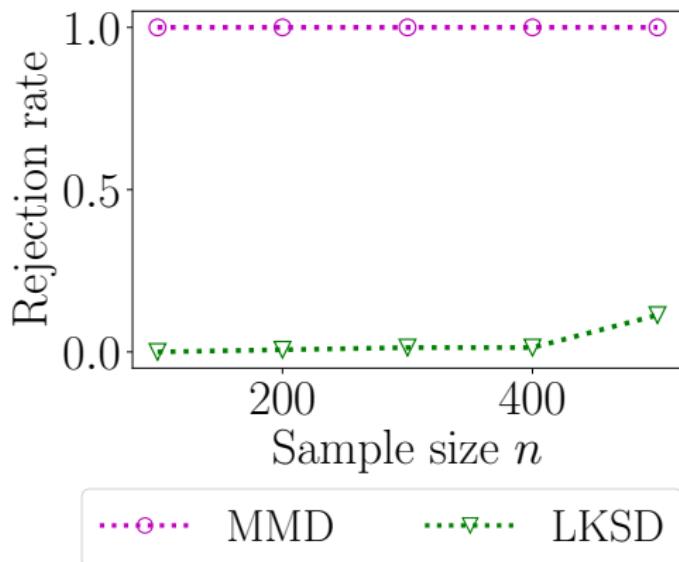


- $\mathcal{X} = \{1, \dots, L\}^D$ ,  $D = 100$ ,  
 $L = 126, 190$ .
- IMQ kernel in BoW rep.:  
 $k(x, x') = (1 + \|B(x) - B(x')\|_2^2)^{-1/2}$
- MCMC size  $m = 5000$   
(after  $t = 500$  steps).

.....○..... MMD      .....▽..... LKSD

## 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).  $H_1$ :  $Q$  is better



- $\mathcal{X} = \{1, \dots, L\}^D$ ,  $D = 100$ ,  
 $L = 208,671$ .
- 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):

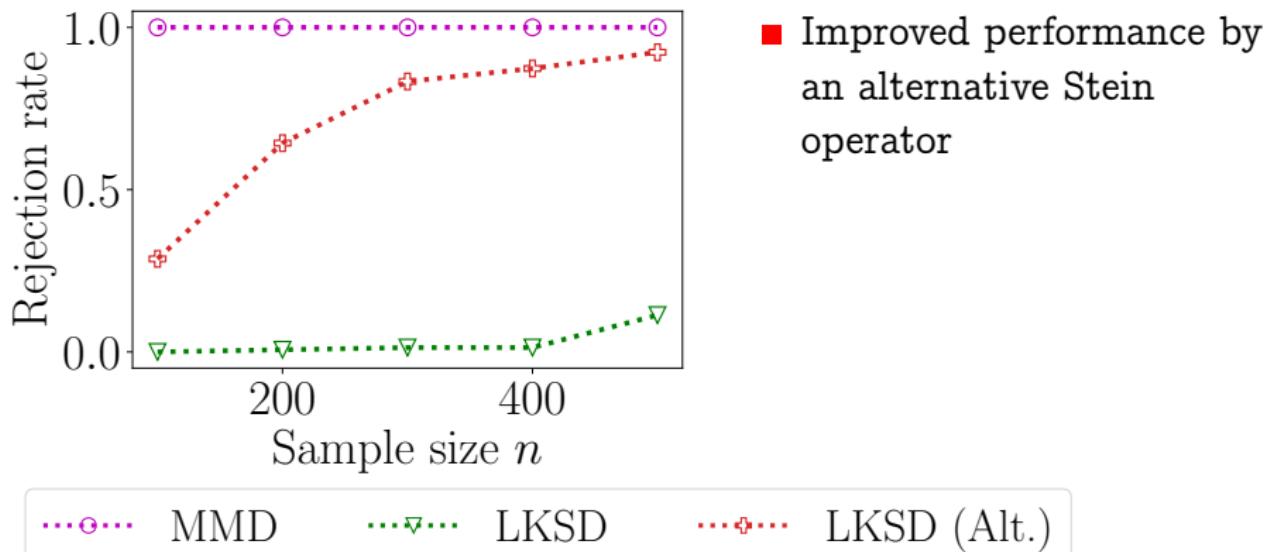
$$\mathcal{A}_p^{\text{ZB}} f(x) = \sum_{\tilde{x} \in \{x+1, x-1\}} \frac{p(\tilde{x})}{p(\tilde{x}) + p(x)} \cdot \{f(\tilde{x}) - f(x)\}$$

- 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).  $H_1$ :  $Q$  is better



# Conclusion

## 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

## 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?



## Can sampler influence test power?

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

Experiment with PPCA:

- $P$  : MALA with a bad step size (poor sampler)
- $Q$  : NUTS-HMC (good sampler)

Expectation:

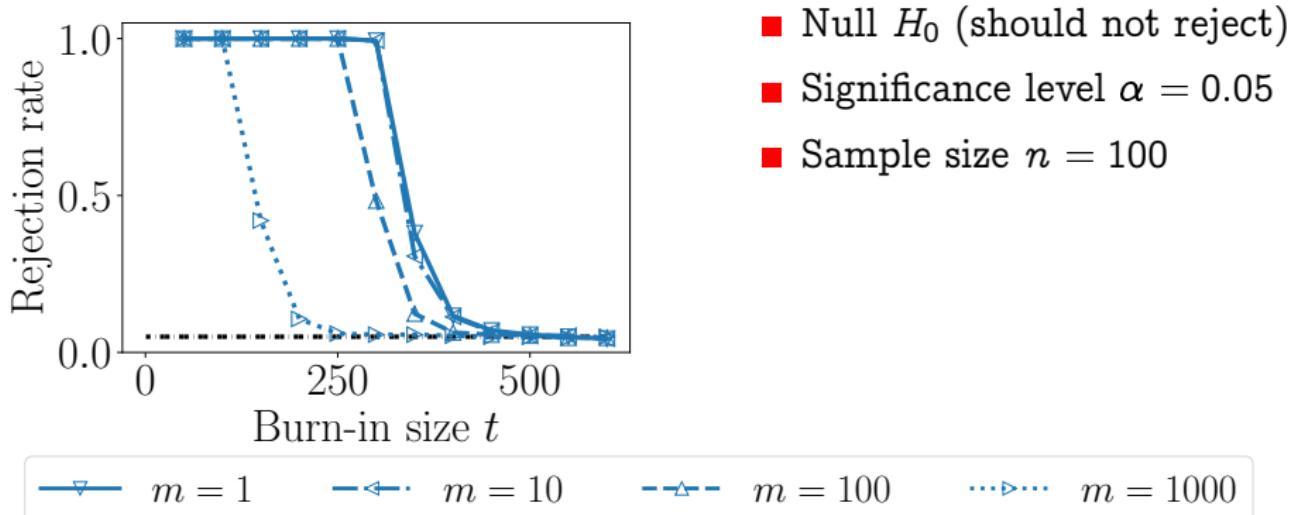
If poor, the test would reject even if  $P$  and  $Q$  are equally good

## Can sampler influence test power?

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

Experiment with PPCA:

- $P$  : MALA with a bad step size (poor sampler)
- $Q$  : NUTS-HMC (good sampler)

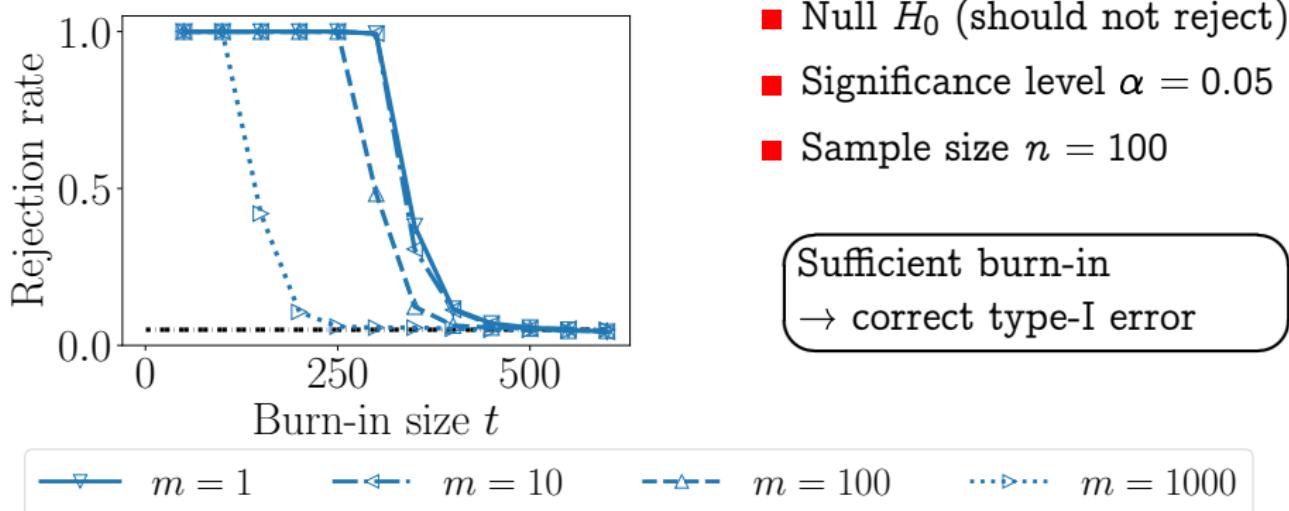


## Can sampler influence test power?

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

Experiment with PPCA:

- $P$  : MALA with a bad step size (poor sampler)
- $Q$  : NUTS-HMC (good sampler)

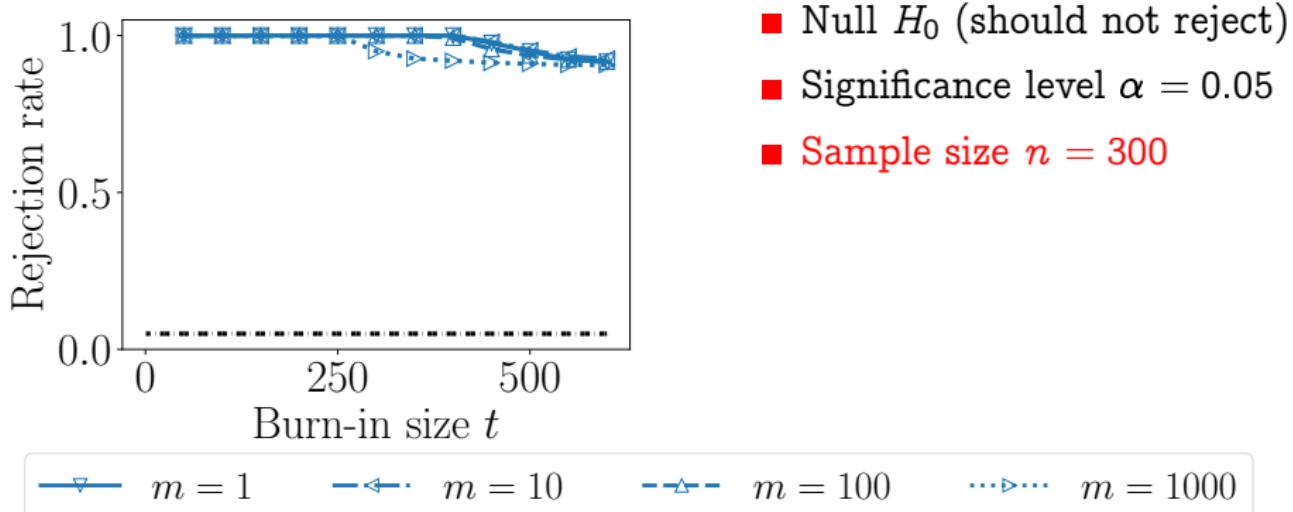


## Can sampler influence test power?

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

Experiment with PPCA:

- $P$  : MALA with a bad step size (poor sampler)
- $Q$  : NUTS-HMC (good sampler)



## Can sampler influence test power?

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

Experiment with PPCA:

- $P$  : MALA with a bad step size (poor sampler)
- $Q$  : NUTS-HMC (good sampler)

