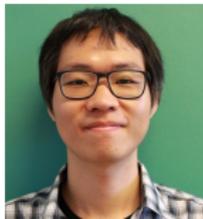


Relative Goodness-of-Fit Tests for Models with Latent Variables

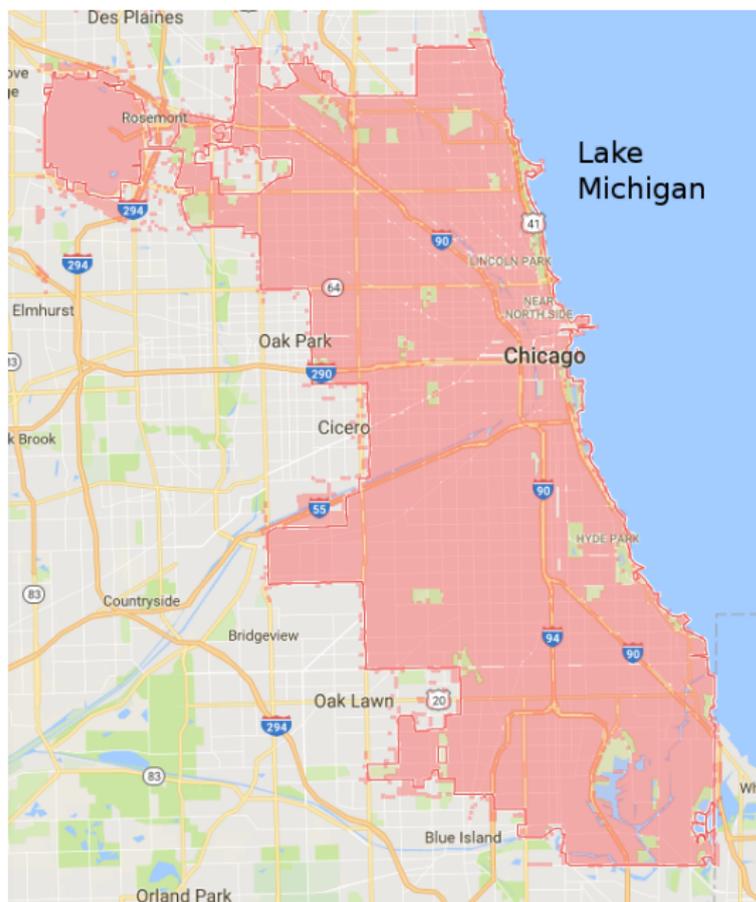
Arthur Gretton



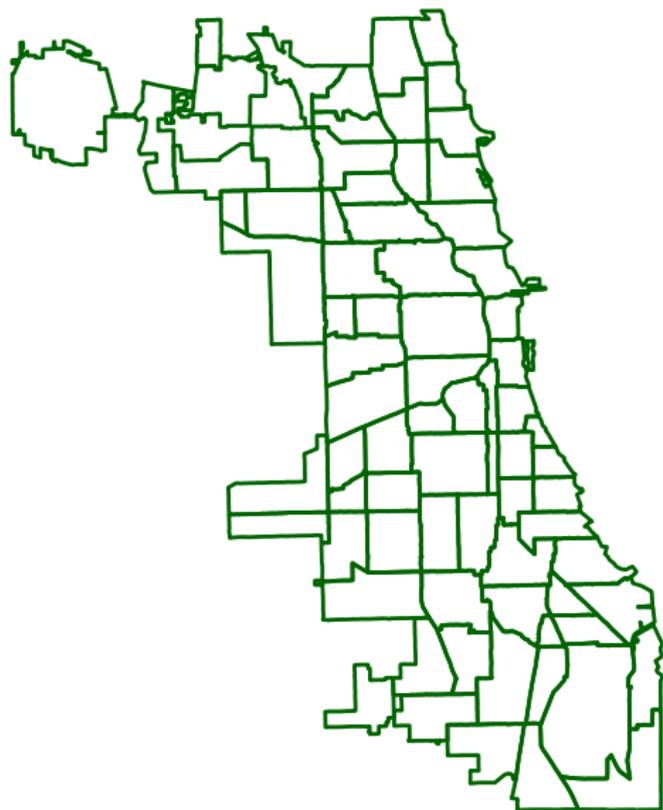
Gatsby Computational Neuroscience Unit,
University College London

BIRS Stein Workshop, April, 2022

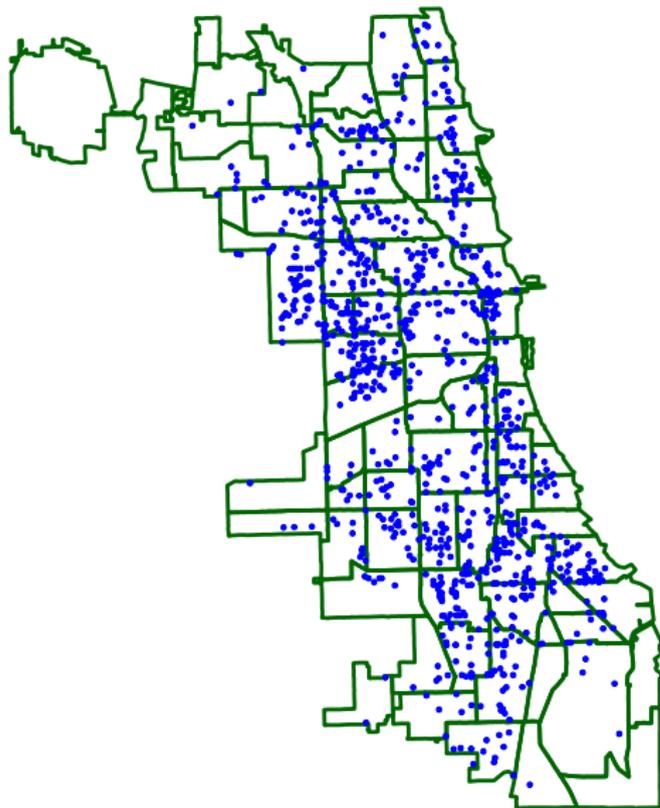
Model Criticism



Model Criticism

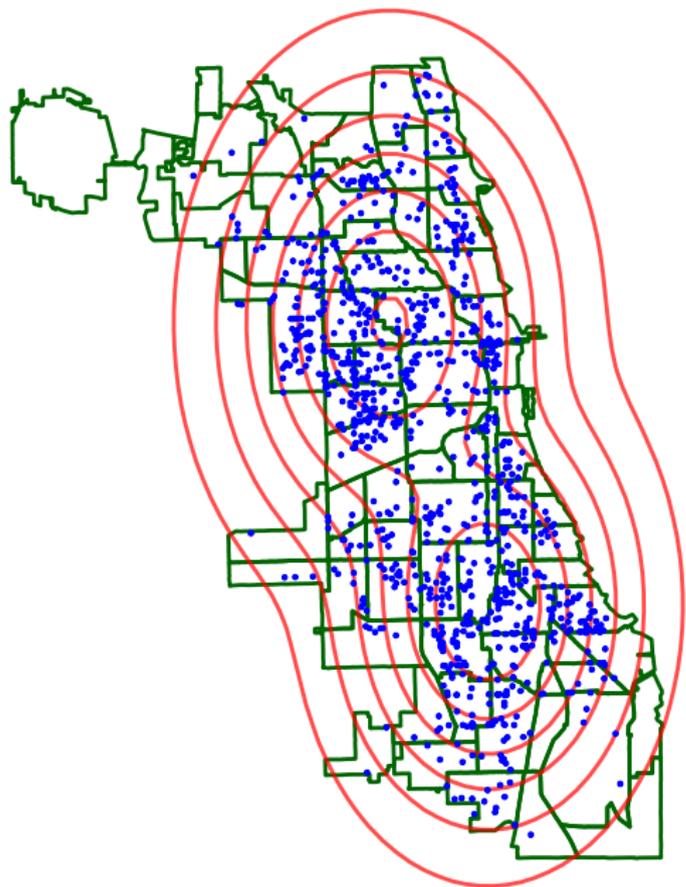


Model Criticism



Data = robbery events in
Chicago in 2016.

Model Criticism



Is this a good **model**?

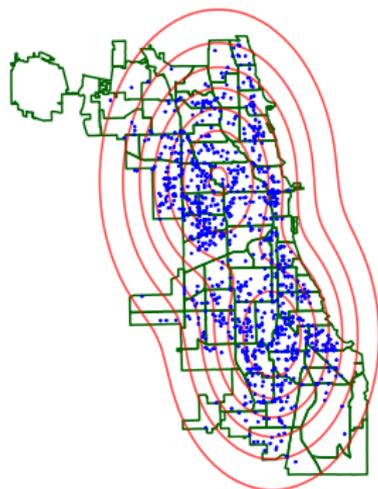
Model Criticism

"All models are wrong."

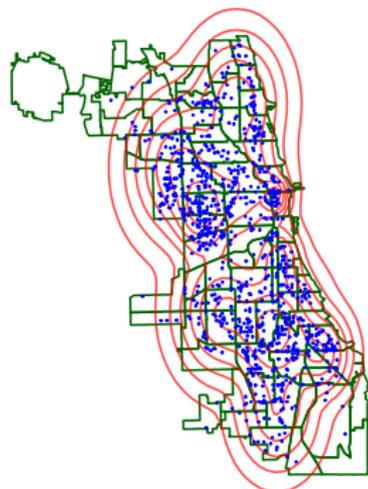
G. Box (1976)

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?

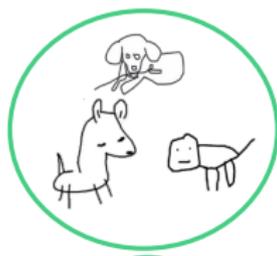
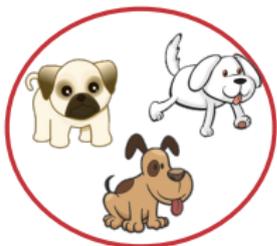


P : two components



Q : ten components

A relative test of goodness-of-fit



R

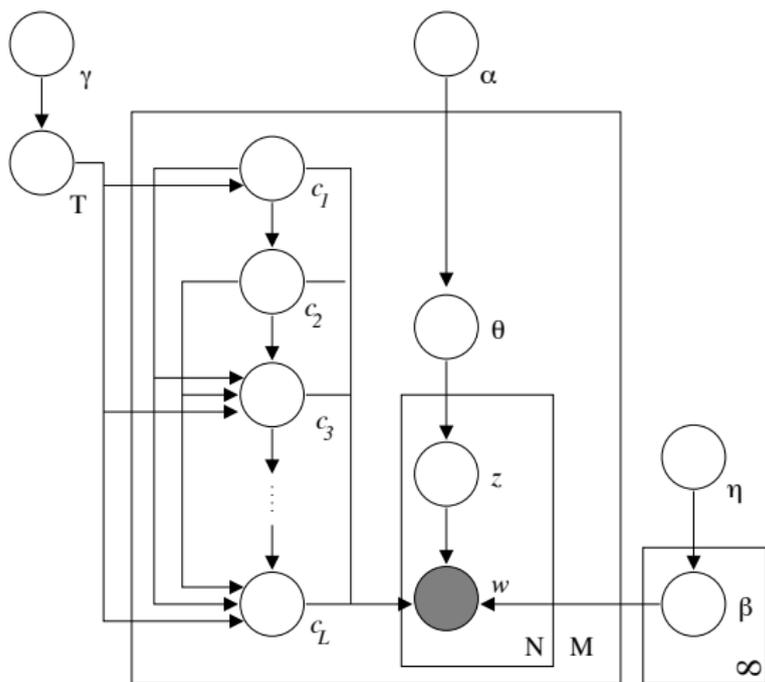
P

Q

$D(\cdot, R)$

Most interesting models have latent structure

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

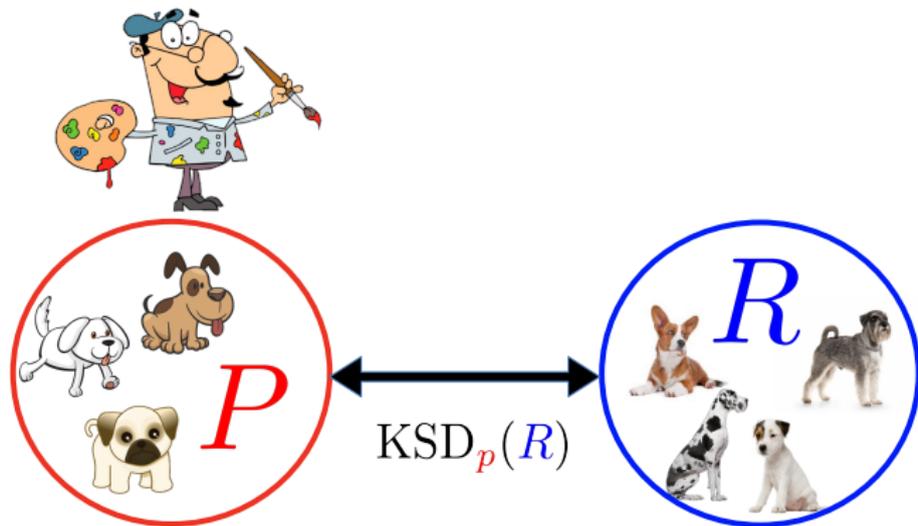


Relative goodness-of-fit tests for Models with Latent Variables

- 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 $\{\mathbf{x}_i\}_{i=1}^n \sim Q$.
- “All models are wrong” ($P \neq Q$).

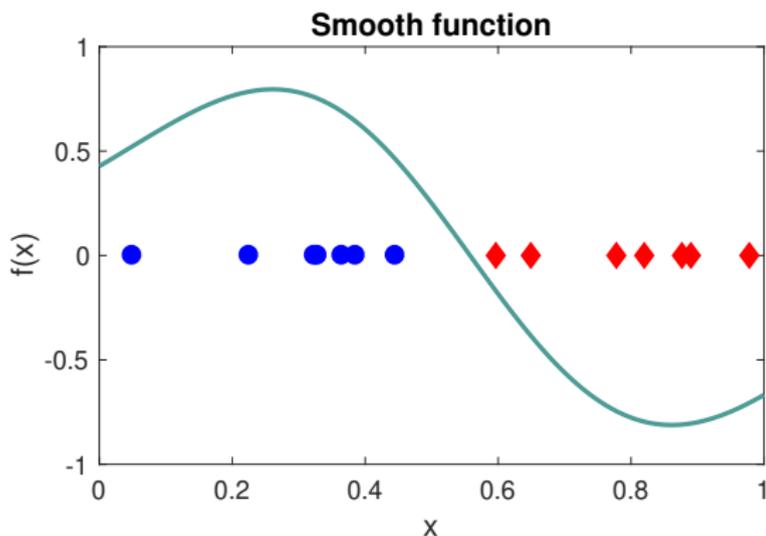


Integral probability metrics

Integral probability metric:

Find a "well behaved function" $f(x)$ to maximize

$$\mathbb{E}_Q f(Y) - \mathbb{E}_P f(X)$$

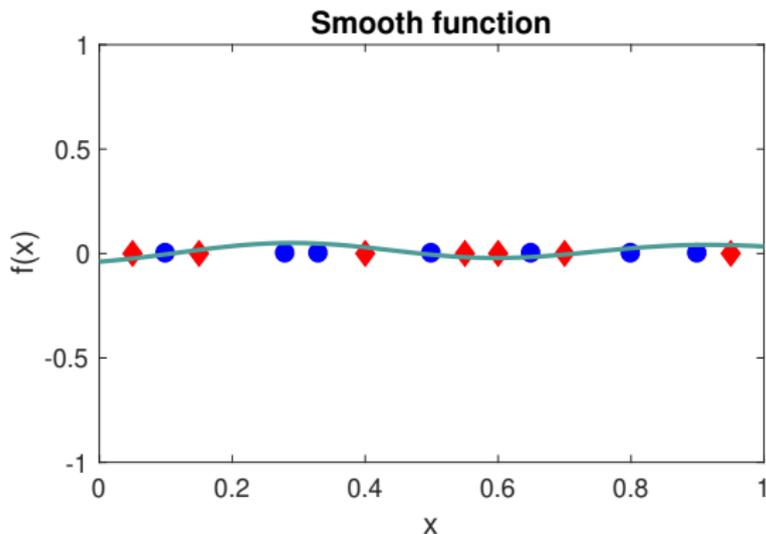


Integral probability metrics

Integral probability metric:

Find a "well behaved function" $f(x)$ to maximize

$$\mathbb{E}_Q f(Y) - \mathbb{E}_P f(X)$$



All of kernel methods

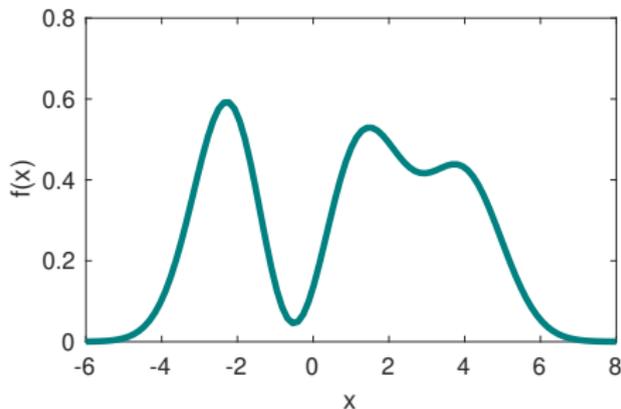
Functions are linear combinations of features:

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

All of kernel methods

“The kernel trick”

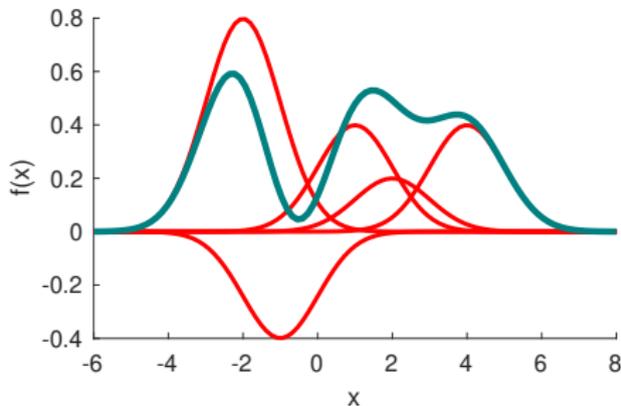
$$\begin{aligned} f(x) &= \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) \\ &= \sum_{i=1}^m \alpha_i \underbrace{k(x_i, x)}_{\langle \varphi(x_i), \varphi(x) \rangle_{\mathcal{F}}} \end{aligned}$$



All of kernel methods

“The kernel trick”

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



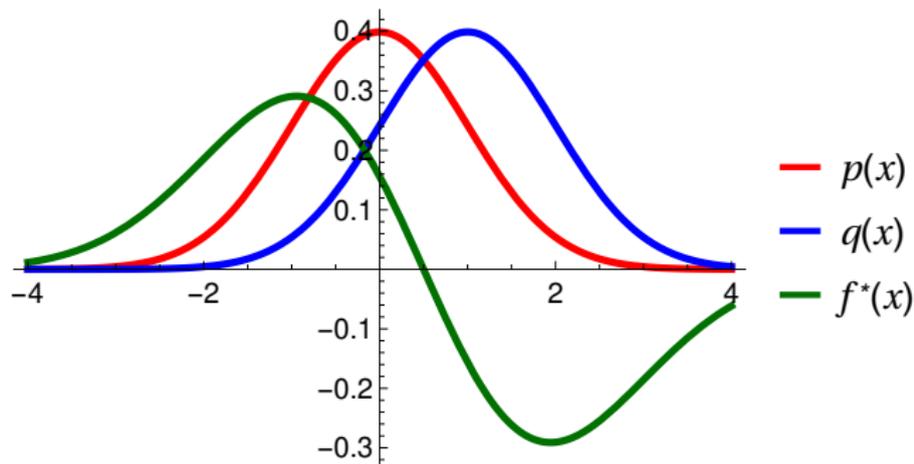
$$f_{\ell} := \sum_{i=1}^m \alpha_i \varphi_{\ell}(x_i)$$

Function of **infinitely many features** expressed using m coefficients.

MMD: an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

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



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:

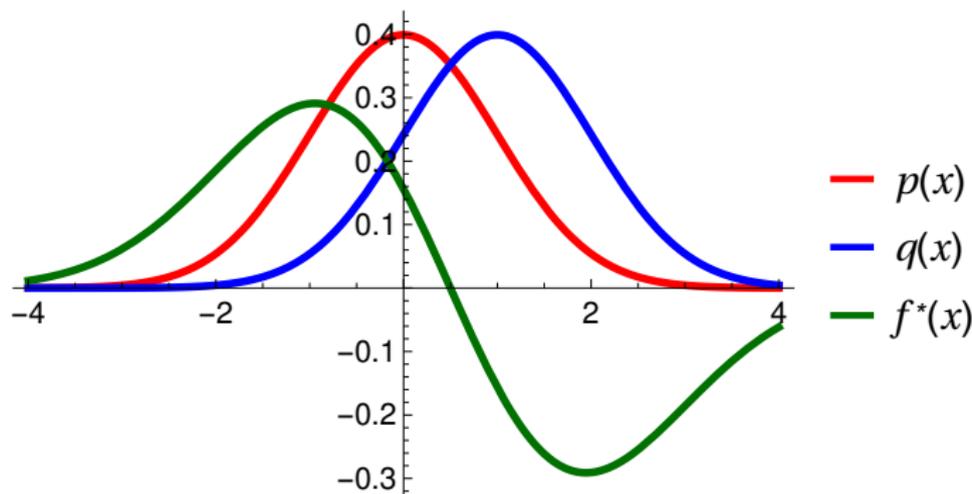
- Bounded continuous [Dudley, 2002]
- Bounded variation 1 (Kolmogorov metric) [Müller, 1997]
- 1-Lipschitz (Wasserstein distances) [Dudley, 2002]

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**

$$[\mathcal{A}_p f](x) = \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x))$$

Then

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

subject to appropriate boundary conditions.

Kernel Stein Discrepancy

Stein operator

$$\mathcal{A}_p f = \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x))$$

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

$$\mathcal{A}_p f = \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x))$$

Kernel Stein Discrepancy (KSD)

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

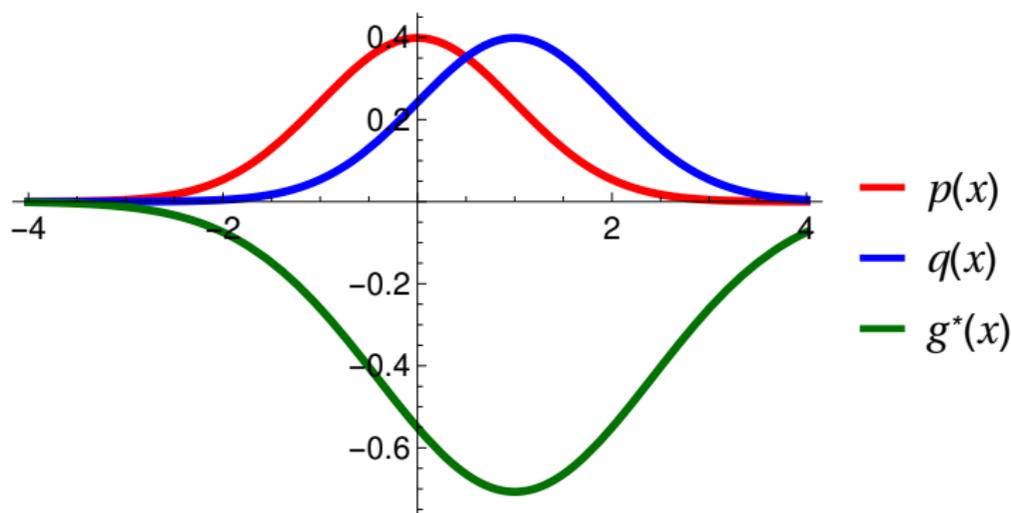
Kernel Stein Discrepancy

Stein operator

$$\mathcal{A}_p f = \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x))$$

Kernel Stein Discrepancy (KSD)

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



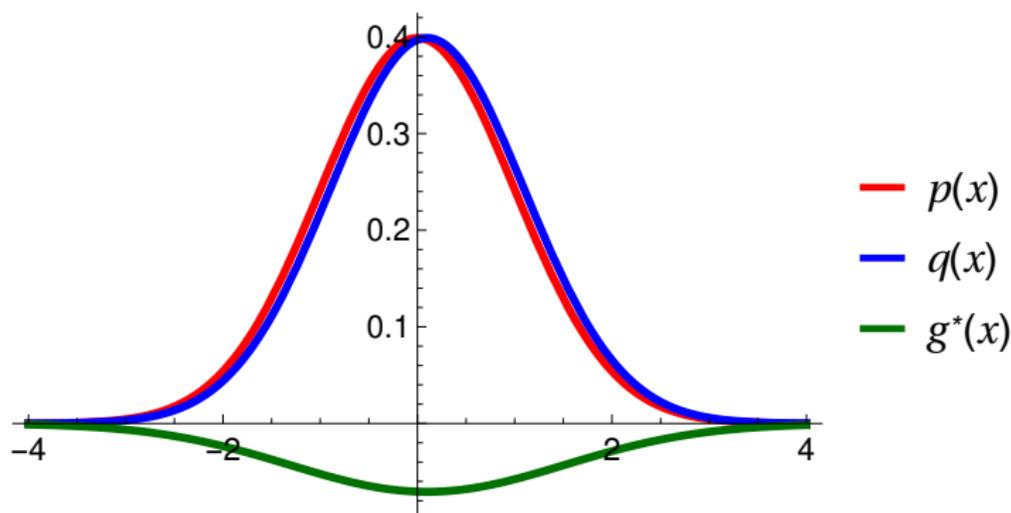
Kernel Stein Discrepancy

Stein operator

$$\mathcal{A}_p f = \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x))$$

Kernel Stein Discrepancy (KSD)

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



Simple expression using kernels

Re-write stein operator as:

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

Can we define “Stein features”?

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

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

Simple expression using kernels

Re-write stein operator as:

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

Can we define “Stein features”?

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

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

The kernel trick for derivatives

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

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

The kernel trick for derivatives

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

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

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

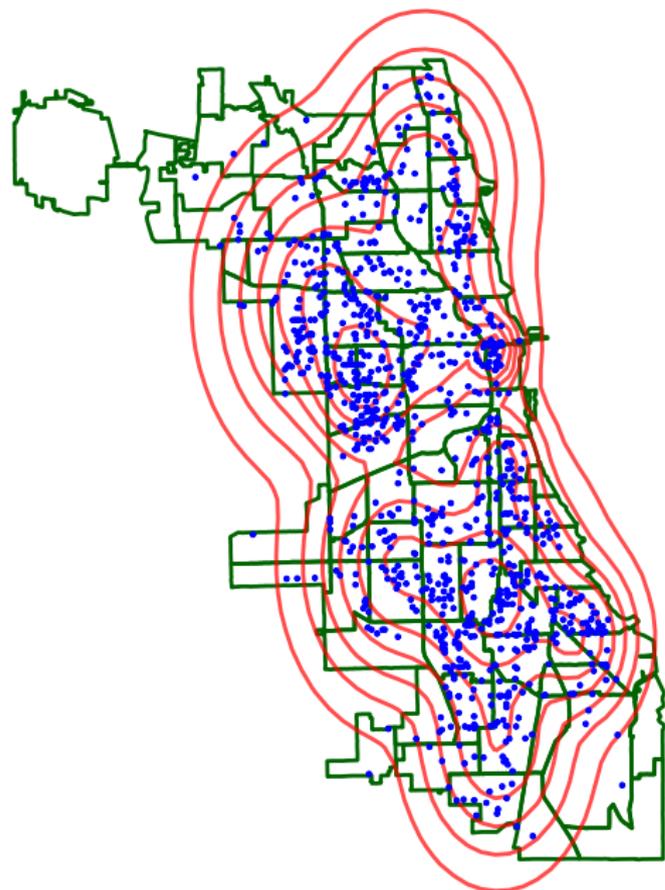
Closed-form expression for KSD: given independent $x, x' \sim Q$, then

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

Caution: (a) requires a condition for the Riesz theorem to hold,

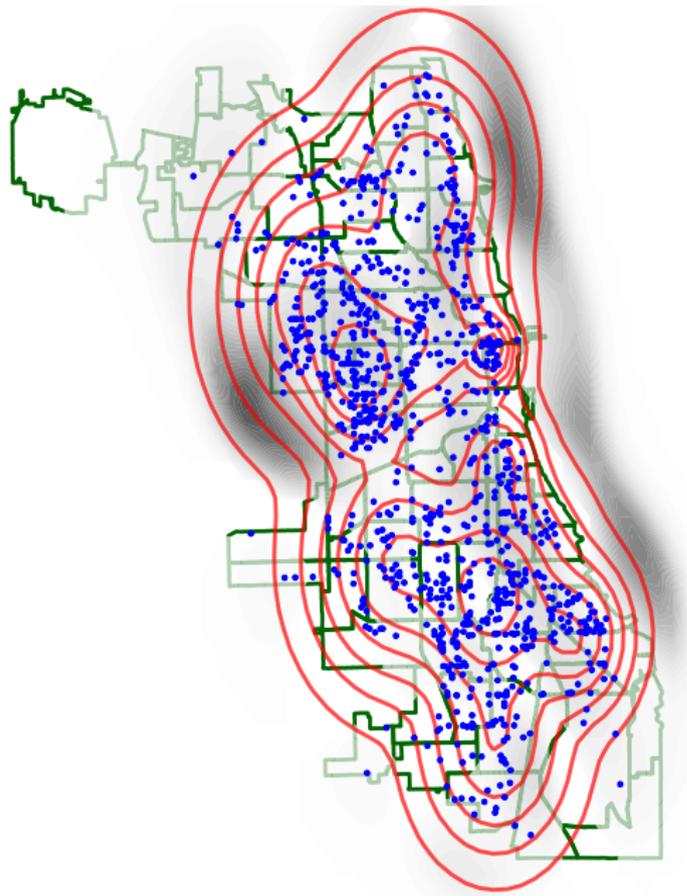
$$\mathbb{E}_{x \sim q} \left(\frac{d}{dx} \log p(x) \right)^2 < \infty.$$

The witness function: Chicago Crime



Model p = 10-component
Gaussian mixture.

The witness function: Chicago Crime



Witness function g shows mismatch

Does the Riesz 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 Riesz 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.

Kernel stein discrepancy: population expression

Test statistic:

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

where

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

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

Kernel stein discrepancy: population expression

Test statistic:

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

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

Test statistic:

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

where

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

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

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

Kernel stein discrepancy: population expression

Test statistic:

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

where

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

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

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

KSD for discrete-valued variables

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

The population KSD (discrete):

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

where

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

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

KSD for discrete-valued variables

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

The population KSD (discrete):

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

where

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

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

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

KSD for discrete-valued variables

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

The population KSD (discrete):

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

where

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

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

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

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

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

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

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

The empirical statistic:

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

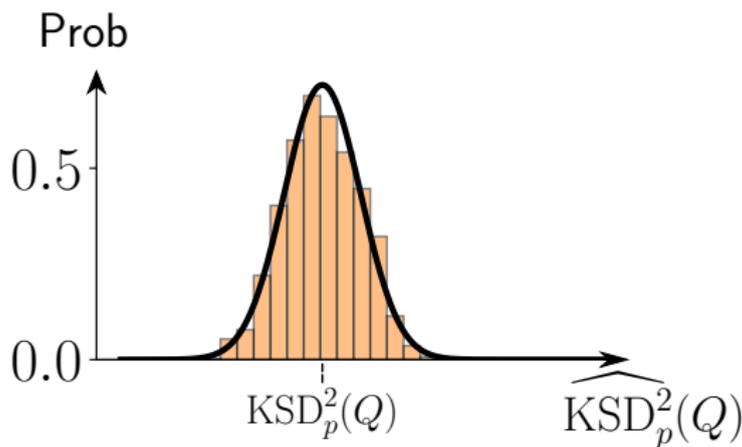
Empirical statistic, asymptotic normality for $P \neq Q$

The empirical statistic:

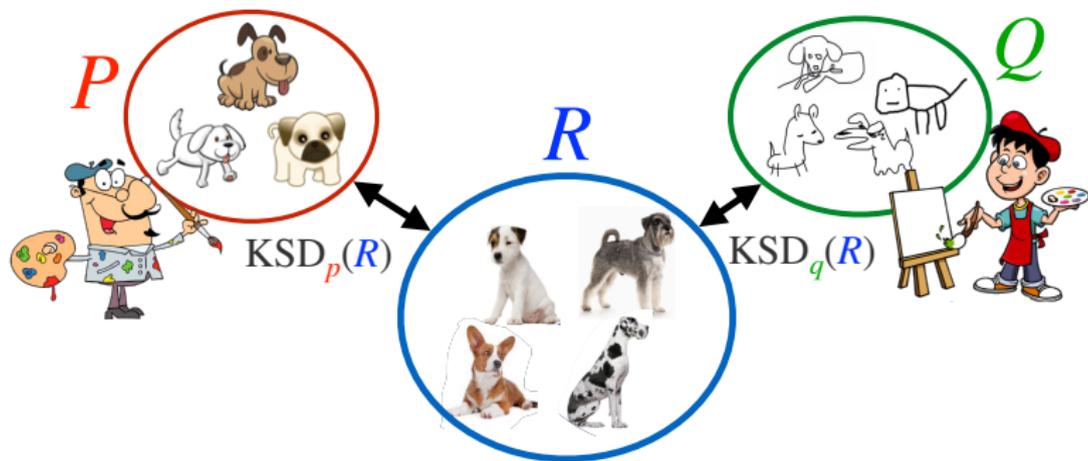
$$\widehat{\text{KSD}}_p^2(Q) := \frac{1}{n(n-1)} \sum_{i \neq j} h_p(x_i, x_j).$$

Asymptotic distribution when $P \neq Q$:

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



Relative goodness-of-fit testing



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

Hypotheses:

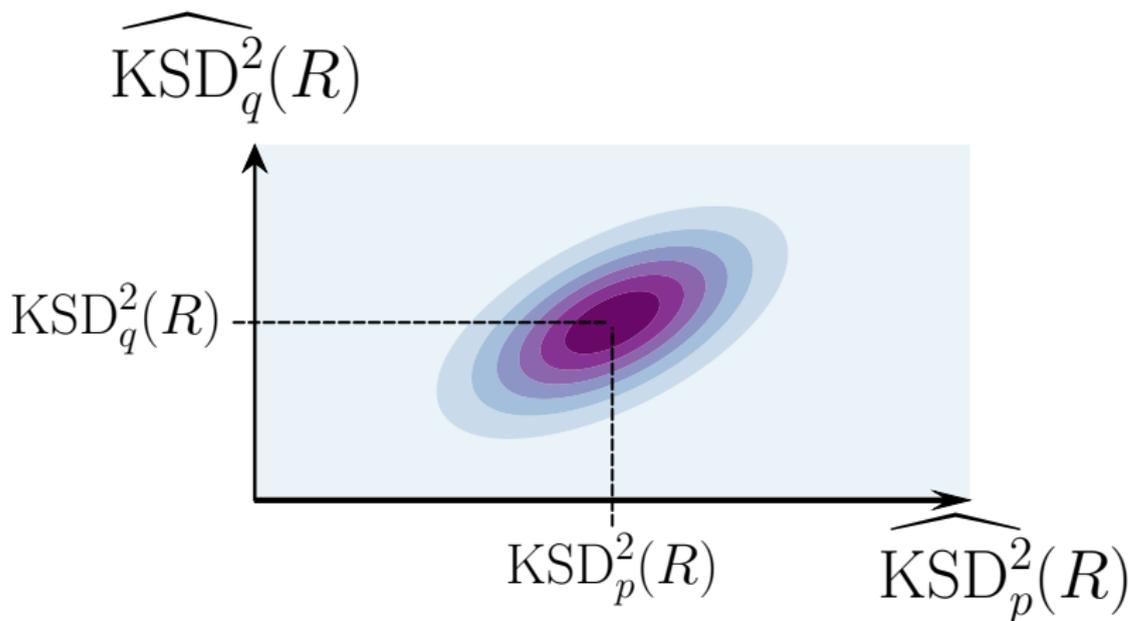
$$H_0 : KSD_p(R) \leq KSD_q(R) \text{ vs. } H_1 : KSD_p(R) > KSD_q(R)$$

(H_0 : ' P 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 $Q \neq R$

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



Relative GOF testing: joint asymptotic normality

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

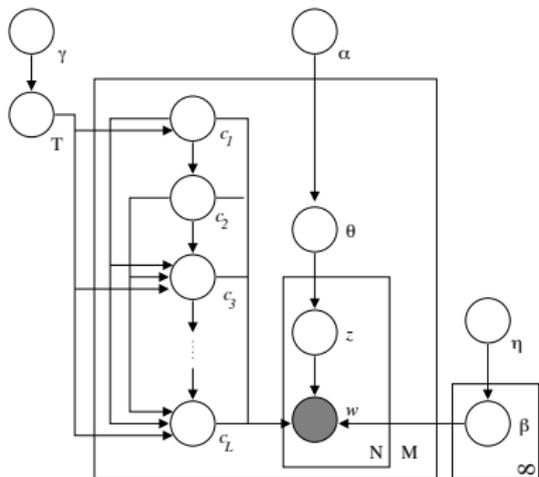
$$\sqrt{n} \begin{bmatrix} \widehat{\text{KSD}}_p^2(R) - \text{KSD}_p(R) \\ \widehat{\text{KSD}}_q^2(R) - \text{KSD}_q(R) \end{bmatrix} \xrightarrow{d} \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{h_p}^2 & \sigma_{h_p h_q} \\ \sigma_{h_p h_q} & \sigma_{h_q}^2 \end{bmatrix}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

Difference in statistics is asymptotically normal:

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

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

Latent variable models

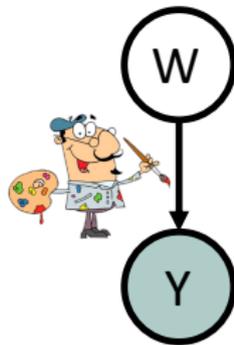
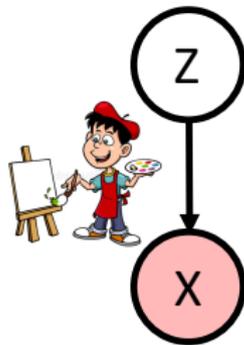


Latent variable models

Can we compare latent variable models with KSD?

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

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



Recall multi-dimensional Stein operator:

$$[T_p f](x) = f(x) \underbrace{\frac{\nabla p(x)}{p(x)}}_{(a)} + \langle \nabla, f(x) \rangle.$$

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

What not to do

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

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

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: $\frac{\nabla p_m(x)}{p_m(x)}$ very numerically unstable. Thus $\widehat{\text{KSD}}_{p_m}^2(R)$ has high variance.

MCMC approximation of score function

Result we use:

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

Proof:

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

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:

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

Proof:

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

Approximate intractable posterior $\mathbb{E}_{z|x_i}[\mathbf{s}_p(\mathbf{x}_i|z)]$

$$\bar{\mathbf{s}}_p(\mathbf{x}_i; z_i^{(t)}) := \frac{1}{m} \sum_{j=1}^m \mathbf{s}_p(\mathbf{x}_i|z_{i,j}^{(t)}) \approx \mathbf{s}_p(\mathbf{x}_i)$$

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:

$$U_n(P) = \frac{1}{n(n-1)} \sum_{i \neq j} h_p(x_i, x_j) (\approx \text{KSD}_p^2(R))$$

KSD for latent variable models

Recall earlier KSD estimate:

$$U_n(P) = \frac{1}{n(n-1)} \sum_{i \neq j} h_p(x_i, x_j) (\approx \text{KSD}_p^2(R))$$

KSD estimate for latent variable models:

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

where \bar{H}_p is the Stein kernel h_p with $s_p(x_i)$ replaced with $\bar{s}_p(x_i; z_i^{(t)})$.

Return to relative GOF test, latent variable models

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')

Return to relative GOF test, latent variable models

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')

Strategy:

- Estimate the difference $\text{KSD}_p^2(R) - \text{KSD}_q^2(R)$ by

$$D_n^{(t)}(P, Q) = U_n^{(t)}(P) - U_n^{(t)}(Q).$$

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

Asymptotic distribution for relative KSD test

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

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

where

$$\mu_{PQ} = \text{KSD}_p^2(R) - \text{KSD}_q^2(R),$$

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

Fine print:

- The double limit requires fast bias decay

$$\sqrt{n} [\mathbb{E}\{D_n^{(t)}(P, Q)\} - \mu_{PQ}] \rightarrow 0 \quad (t \rightarrow \infty).$$

- The fourth moment of $\bar{H}_p^{(t)} - \bar{H}_q^{(t)}$ has finite limit sup. ($t \rightarrow \infty$).

Asymptotic distribution for relative KSD test

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

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

where

$$\mu_{PQ} = \text{KSD}_p^2(R) - \text{KSD}_q^2(R),$$

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

Fine print:

- The double limit requires fast bias decay

$$\sqrt{n} [\mathbb{E}\{D_n^{(t)}(P, Q)\} - \mu_{PQ}] \rightarrow 0 \quad (t \rightarrow \infty).$$

- The fourth moment of $\bar{H}_p^{(t)} - \bar{H}_q^{(t)}$ has finite limit sup. ($t \rightarrow \infty$).

Asymptotic distribution for relative KSD test

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

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

where

$$\mu_{PQ} = \text{KSD}_p^2(R) - \text{KSD}_q^2(R),$$

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

Level- α test:

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

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

Experiments

Experiment 1: sensitivity to model difference

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

$$\mathbf{x}_i \in \mathbb{R}^{100} \sim \mathcal{N}(A\mathbf{z}_i, I), \quad \mathbf{z}_i \in \mathbb{R}^{10} \sim \mathcal{N}(0, I_z)$$

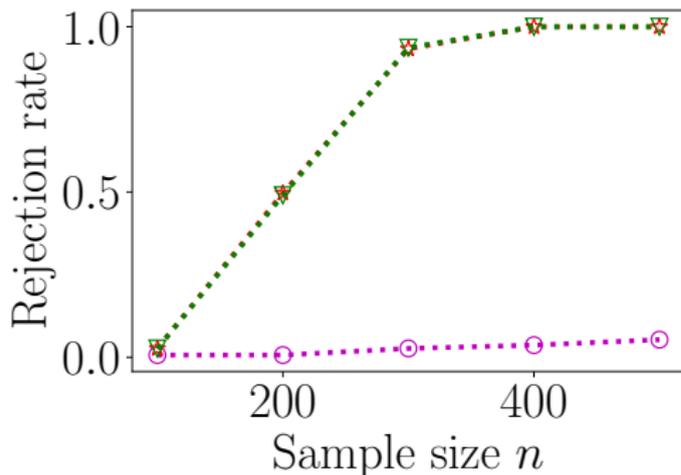
- Generate P , Q : perturb (1, 1)-entry : $A_\delta = A + \delta E_{1,1}$

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).

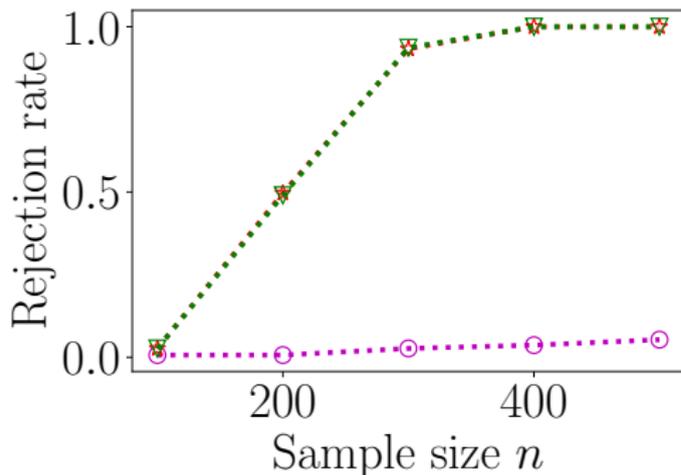
.....○..... MMD ☆..... KSD ▽..... LKSD

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 : $A_\delta = A + \delta E_{1,1}$



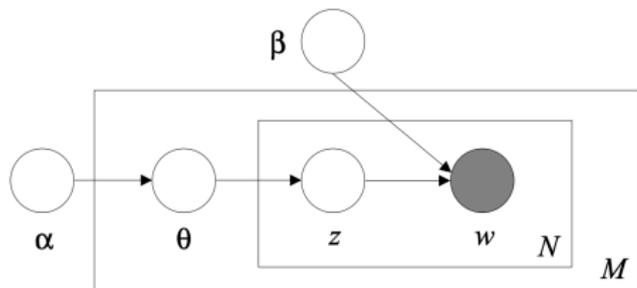
KSD = higher power

- Sample-wise difference in models = subtle (MMD fails)
- Model's information is exploited

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

Experiment 2: topic models for arXiv articles

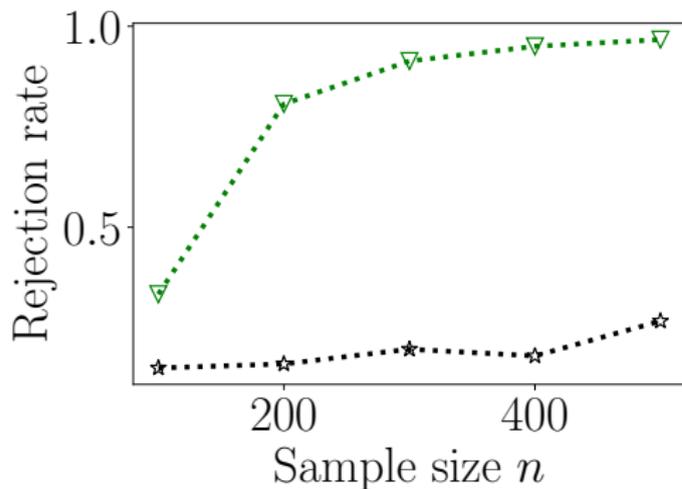
- Data R : arXiv articles from category stat.TH (stat theory) :
- Models P, Q : LDAs trained on articles from different categories
 - P : math.PR (math probability theory)
 - Q : stat.ME (stat methodology)



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)

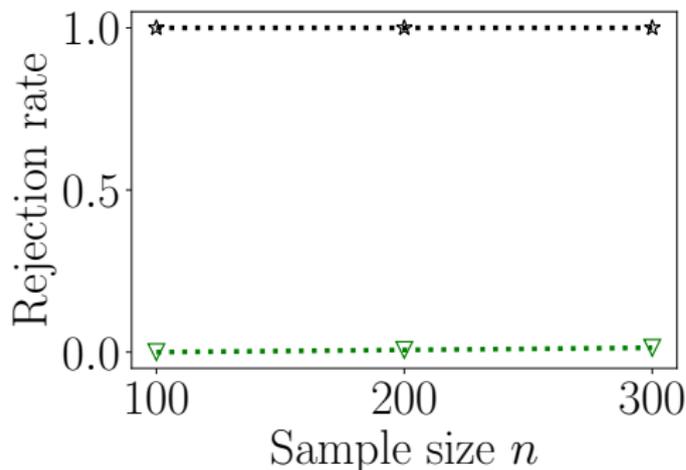


- $\mathcal{X} = \{1, \dots, L\}^D$, $D = 100$, $L = 126, 190$.
- IMQ kernel in BoW rep.:
$$k(x, x') = (1 + \|B(x) - B(x')\|_2^2)^{-1/2}$$
- MCMC size $m = 5000$ (after $t = 500$ steps).

...☆... MMD (IMQBoW) ...▽... LKSD (IMQ BoW)

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)



.....*..... MMD (IMQBoW) ▽..... LKSD (IMQ BoW)

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

LDA's score = concatenation of 1d-score functions (by conditional independence)

$$s_p(x) = (s_{p,1}(x), \dots, s_{p,d}(x), \dots, s_{p,D}(x))$$

$$\text{where } s_{p,d}(x) = \mathbb{E}_{z^d|x}[s_p(x^d|z^d)] = \mathbb{E}_{z^d|x} \left[\frac{p(x^d+1|z^d, \beta)}{p(x^d|z^d, \beta)} \right] - 1$$

⇒ Higher chance of instability

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

LDA's score = concatenation of 1d-score functions (by conditional independence)

$$s_p(x) = (s_{p,1}(x), \dots, s_{p,d}(x), \dots, s_{p,D}(x))$$

$$\text{where } s_{p,d}(x) = \mathbb{E}_{z^d|x}[s_p(x^d|z^d)] = \mathbb{E}_{z^d|x} \left[\frac{p(x^d+1|z^d, \beta)}{p(x^d|z^d, \beta)} \right] - 1$$

⇒ Higher chance of instability

Observations on the sampler

Requirements on n and t

The KSD difference estimate $D_n^{(t)}(P, Q)$ is biased for finite t :

$$\mathbb{E}[D_n^{(t)}(P, Q)] \neq \mu_{P, Q} := \text{KSD}_P^2(R) - \text{KSD}_Q^2(R)$$

If the bias decay is slower than \sqrt{n} , i.e.,

$$\underbrace{\sqrt{n} \left(\mathbb{E}[D_n^{(t)}(P, Q)] - \mu_{P, Q} \right)}_{\text{bias}(t) \downarrow 0} \not\rightarrow 0,$$

then, the asymptotic normality around $\mu_{P, Q}$ does not hold:

$$\sqrt{n} \left[D_n^{(t)}(P, Q) - \mu_{P, Q} \right] \not\xrightarrow{d} \mathcal{N}(0, \sigma_{P, Q}^2).$$

Poor MCMC hurts the test

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

Experiment with PPCA:

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

Expectation:

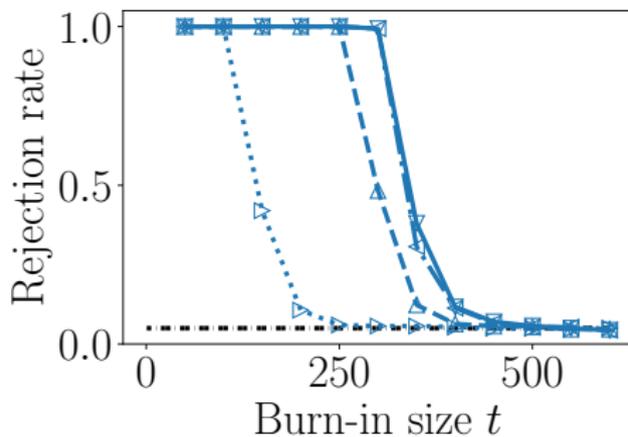
If poor, the test would reject even if P and Q are equally good

Poor MCMC hurts the test

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

Experiment with PPCA:

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



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

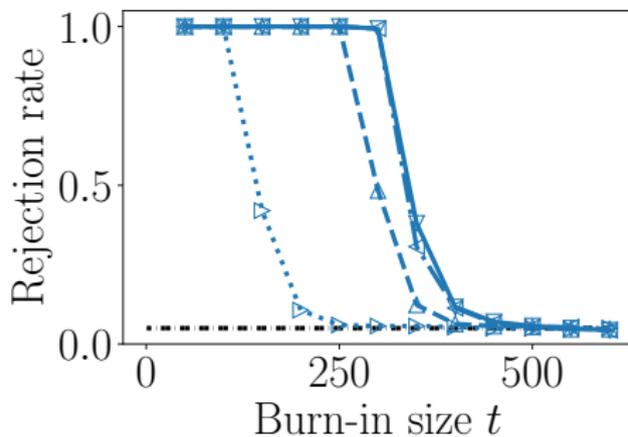
— ∇ — $m = 1$ - \triangleleft - $m = 10$ - \triangle - $m = 100$ \cdots \triangleright \cdots $m = 1000$

Poor MCMC hurts the test

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

Experiment with PPCA:

- P : MALA with a bad step size (poor sampler)
- 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

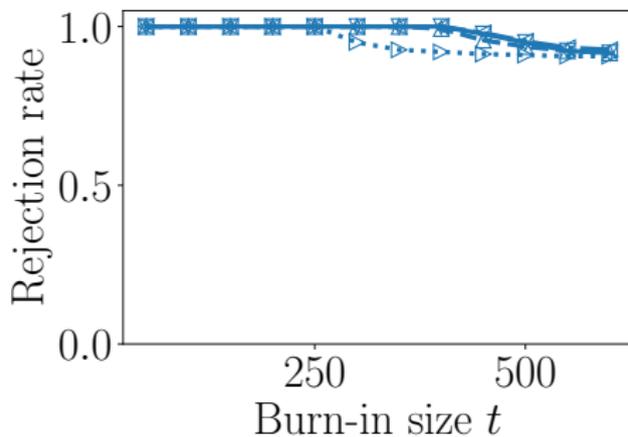
—▽— $m = 1$ -◀- $m = 10$ -△- $m = 100$ ···▷··· $m = 1000$

Poor MCMC hurts the test

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

Experiment with PPCA:

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



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

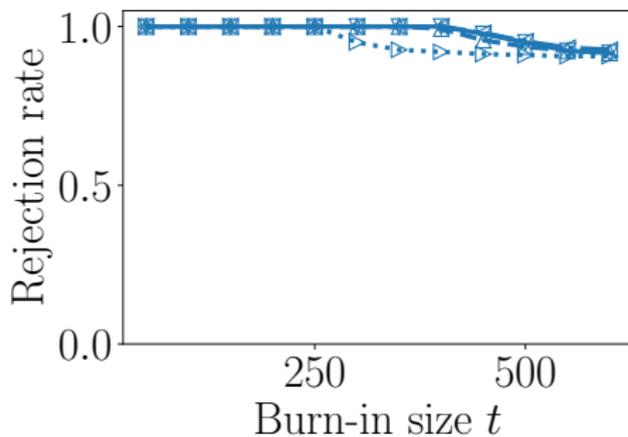
— ∇ — $m = 1$ - \leftarrow - $m = 10$ - \triangle - $m = 100$ \cdots \triangleright \cdots $m = 1000$

Poor MCMC hurts the test

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)



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

Larger $n \implies$ more sensitive to mismatch

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

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?

