

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

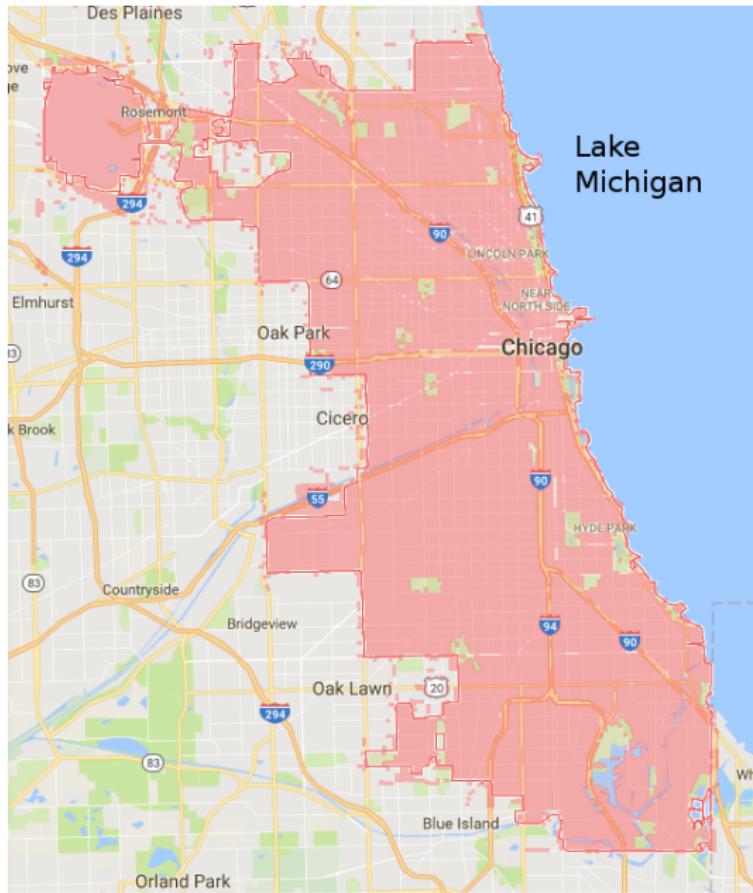
Arthur Gretton



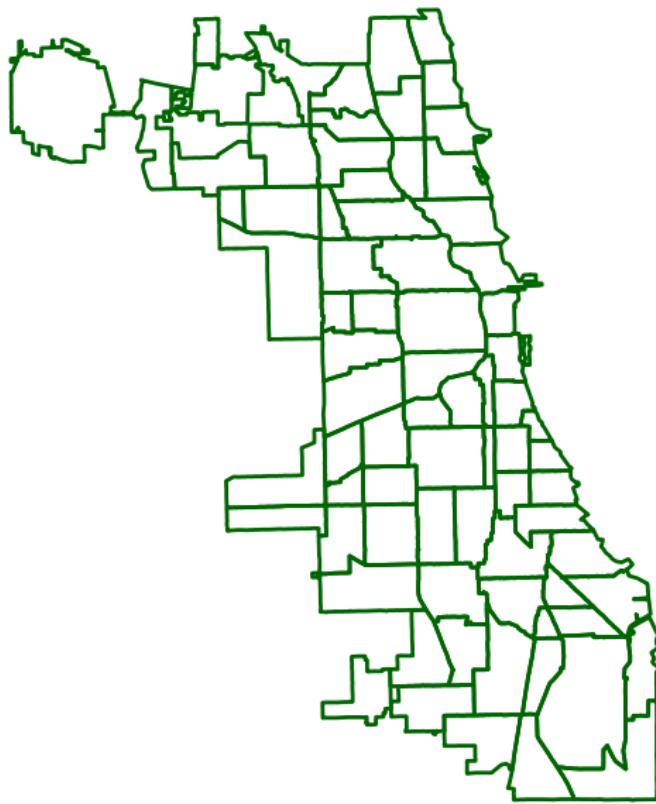
Gatsby Computational Neuroscience Unit,
University College London

June 15, 2019

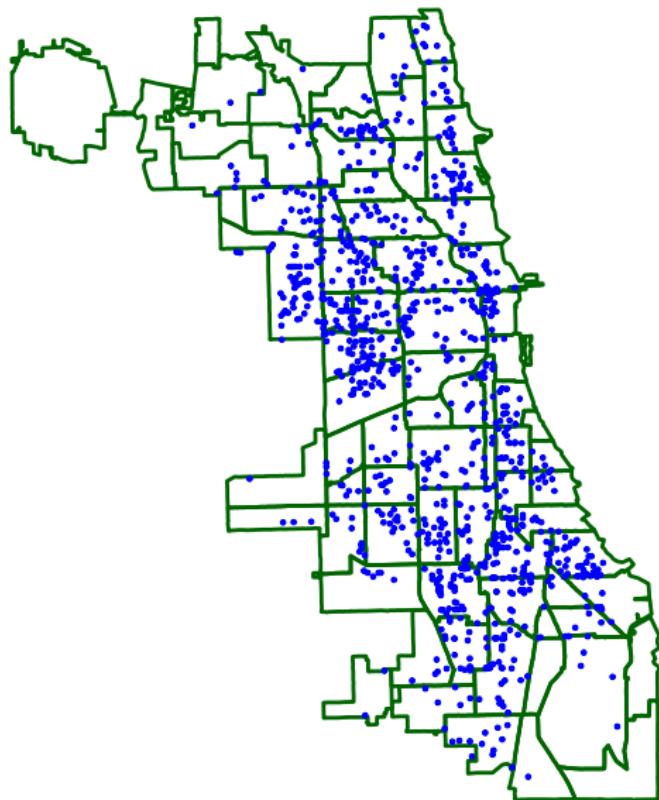
Model Criticism



Model Criticism

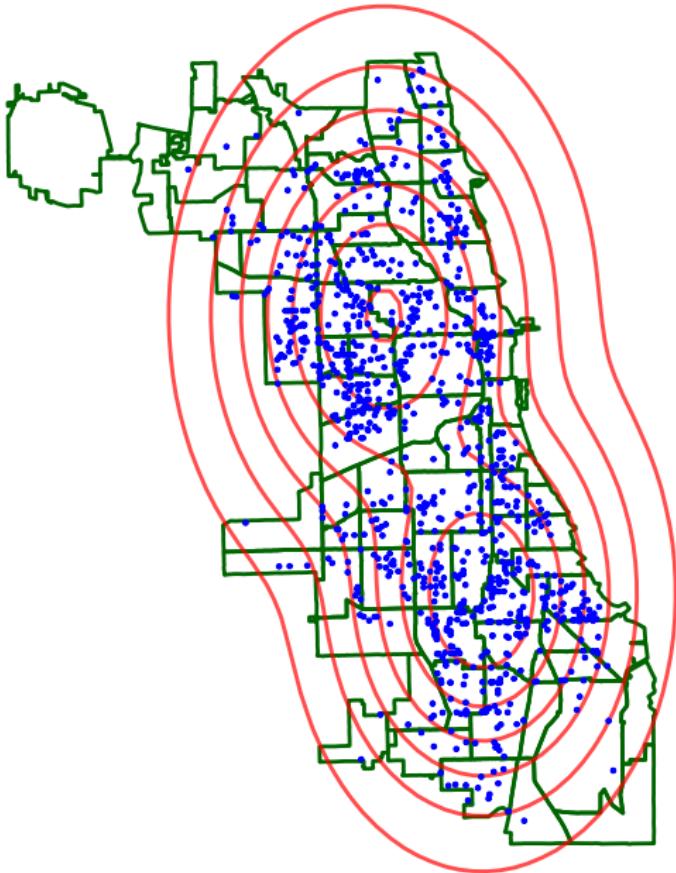


Model Criticism



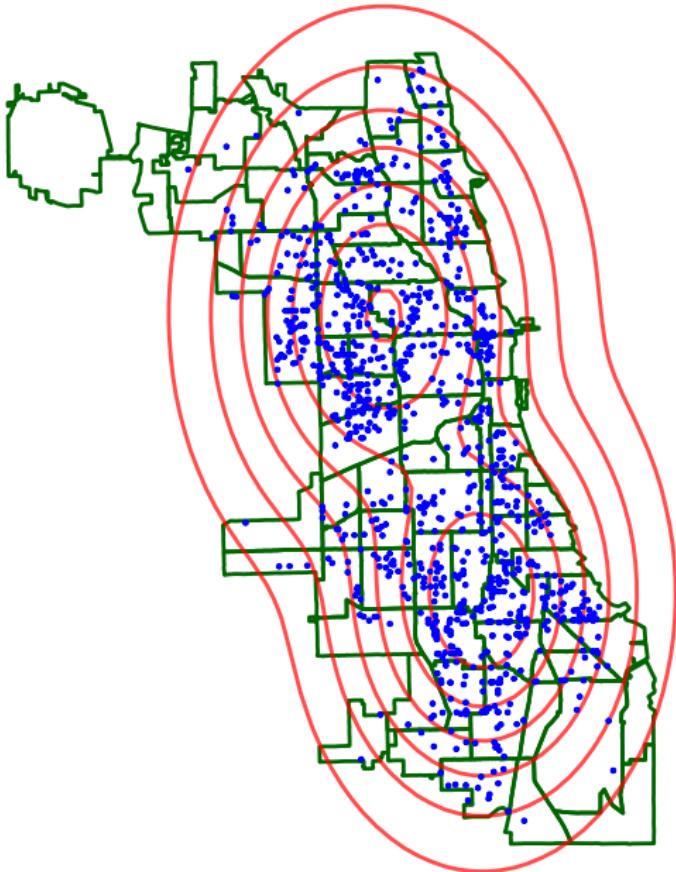
Data = robbery events in Chicago in 2016.

Model Criticism



Is this a good **model**?

Model Criticism



Goals: Test if a (complicated) model fits the data.

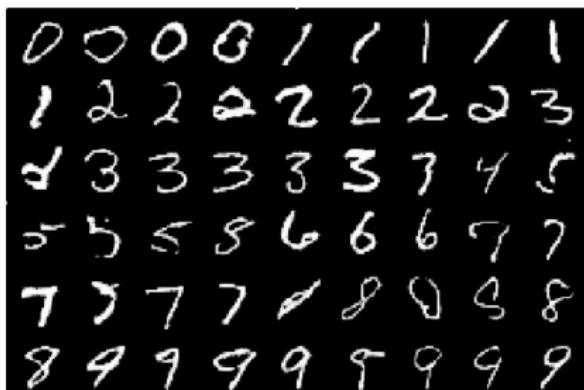
Model Criticism

"All models are wrong."

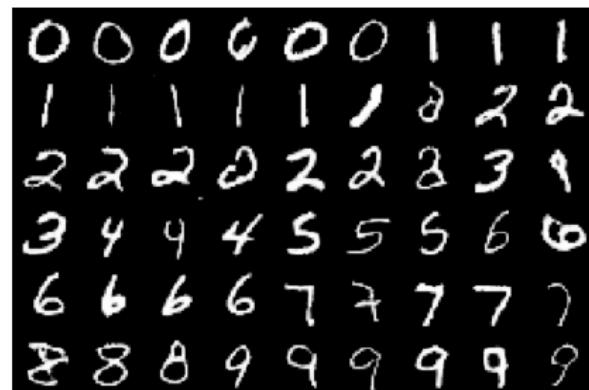
G. Box (1976)

Relative 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?



Samples from GAN,
Goodfellow et al. (2014)

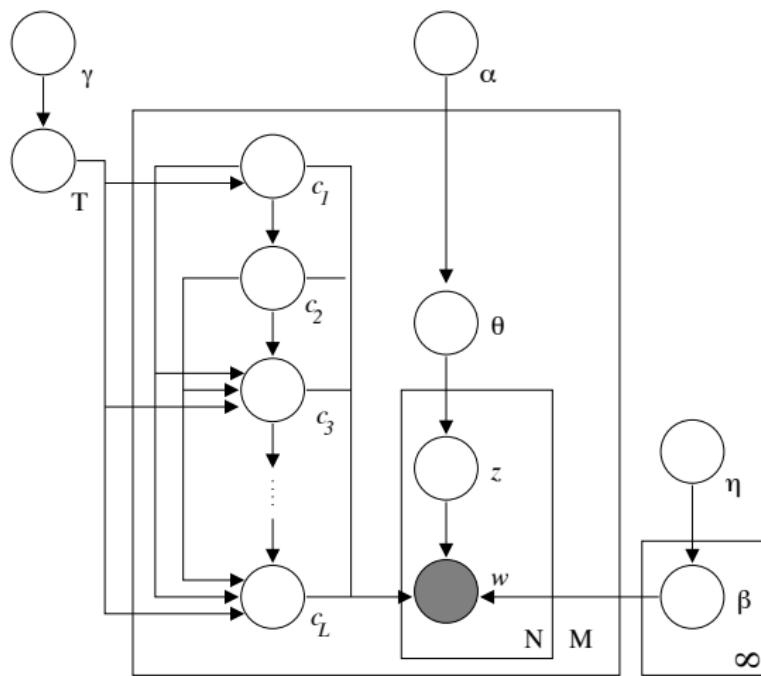


Samples from LSGAN,
Mao et al. (2017)

Which model is better?

Most interesting models have latent structure

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



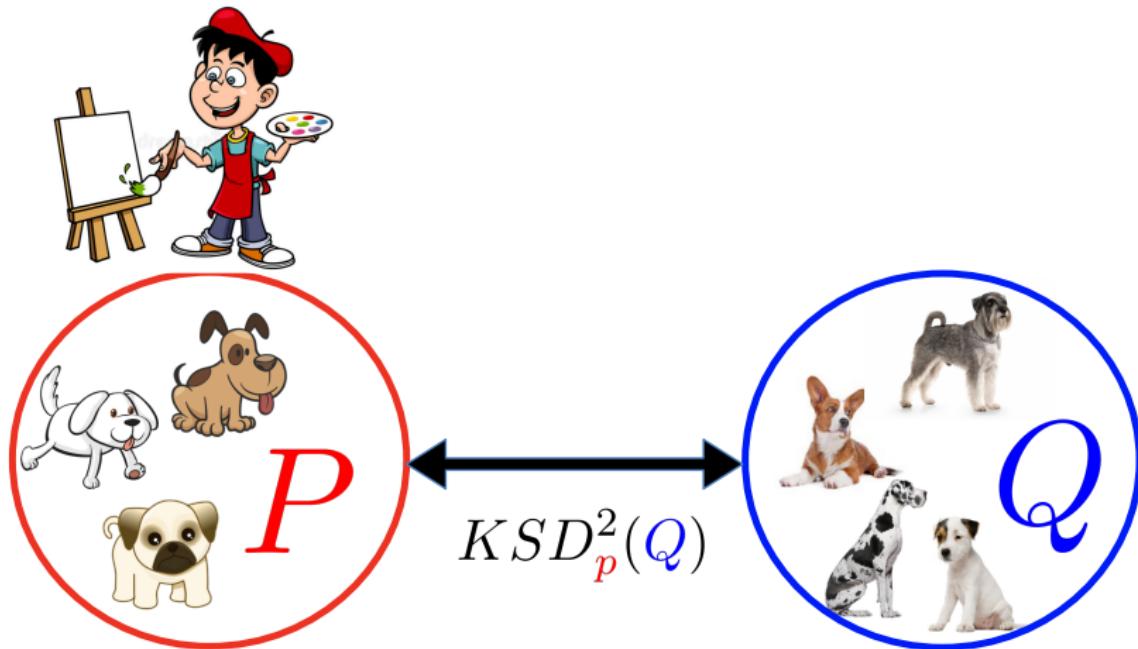
Outline

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 (new, unpublished)

Kernel Stein Discrepancy

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

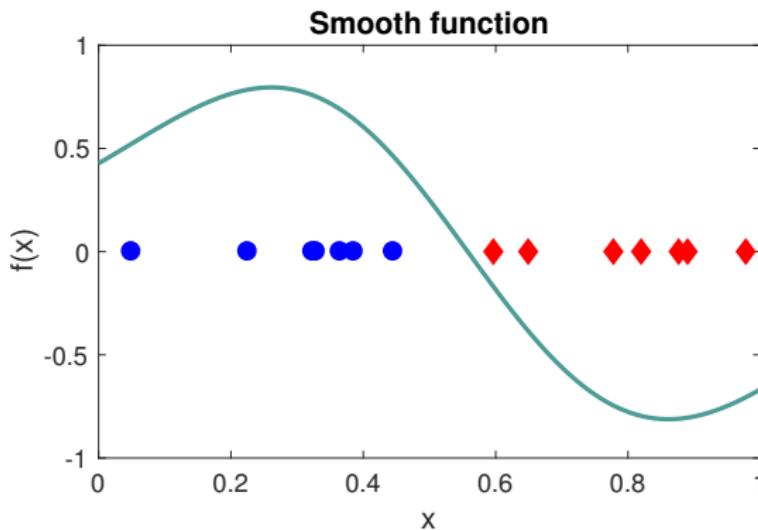


Integral probability metrics

Integral probability metric:

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

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

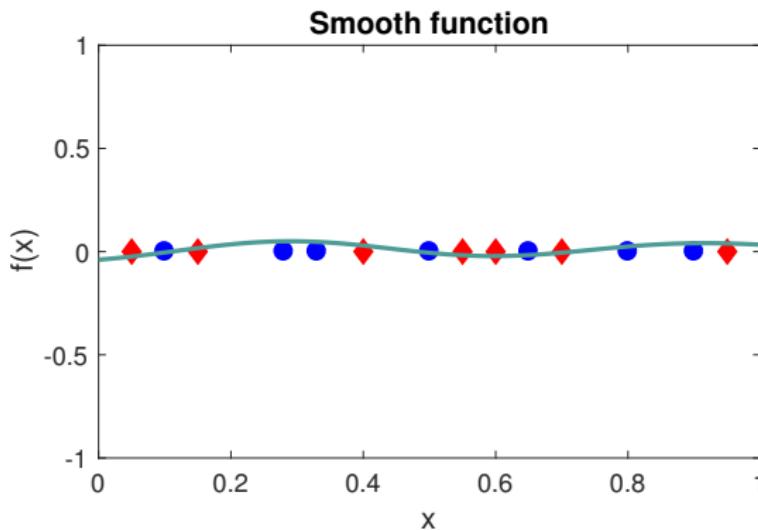


Integral probability metrics

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:

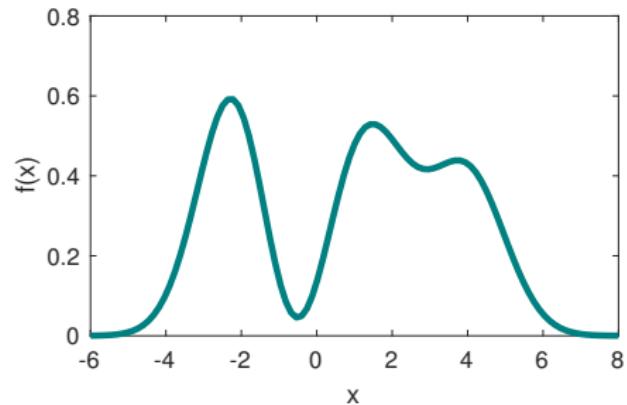
$$f(x) = \langle \mathbf{f}, \varphi(x) \rangle_{\mathcal{F}} = \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \end{bmatrix}^T \begin{bmatrix} \varphi_1(x) \\ \varphi_2(x) \\ \varphi_3(x) \\ \vdots \end{bmatrix}$$

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

All of kernel methods

“The kernel trick”

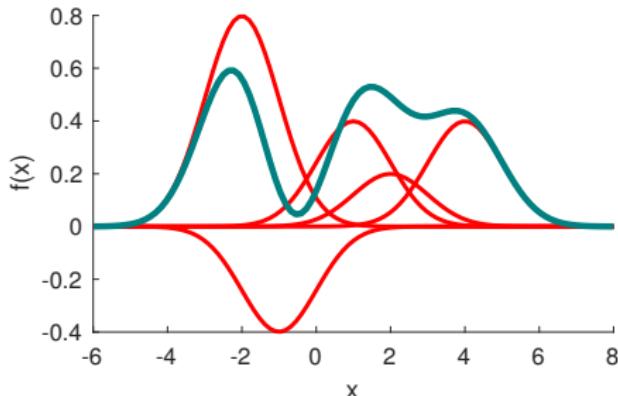
$$\begin{aligned} f(x) &= \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) \\ &= \sum_{i=1}^m \alpha_i k(x_i, x) \end{aligned}$$



All of kernel methods

“The kernel trick”

$$\begin{aligned}f(x) &= \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) \\&= \sum_{i=1}^m \alpha_i k(x_i, x)\end{aligned}$$



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

Function of infinitely many features expressed using m coefficients.

MMD as an integral probability metric

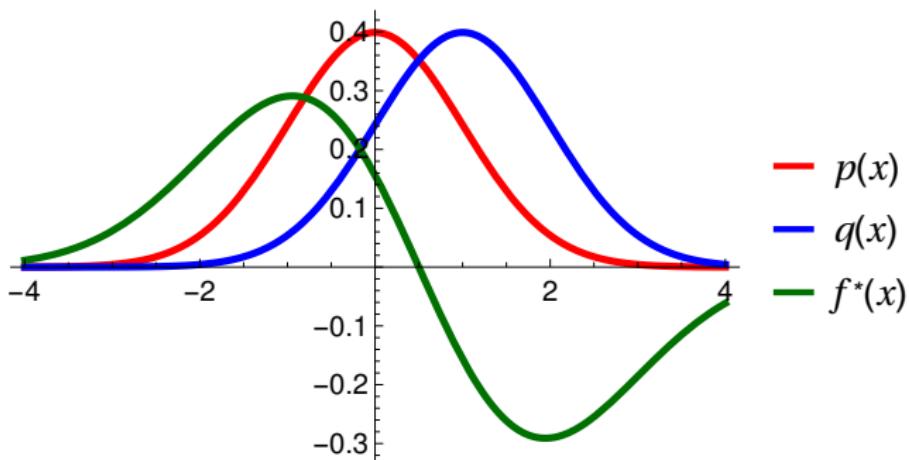
Maximum mean discrepancy: smooth function for P vs Q

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

MMD as an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

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



MMD as an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

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

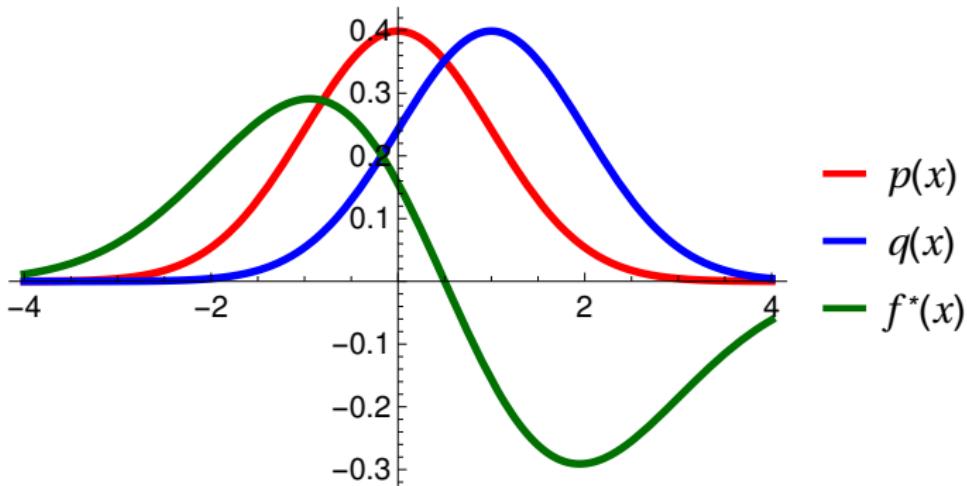
For characteristic RKHS \mathcal{F} , $\text{MMD}(P, Q; \mathcal{F}) = 0$ iff $P = Q$

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

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



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

Problem: usually can't compute \mathbf{E}_{pf} in closed form.

Stein idea

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

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

we define the (1-D) **Stein operator**

$$[\mathcal{A}_{\textcolor{red}{p}} \textcolor{teal}{f}] (x) = \frac{1}{\textcolor{red}{p}(x)} \frac{d}{dx} (\textcolor{teal}{f}(x) \textcolor{red}{p}(x))$$

Then

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

subject to appropriate boundary conditions.

Kernel Stein Discrepancy

Stein operator

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

Kernel Stein Discrepancy (KSD)

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

Kernel Stein Discrepancy

Stein operator

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

Kernel Stein Discrepancy (KSD)

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

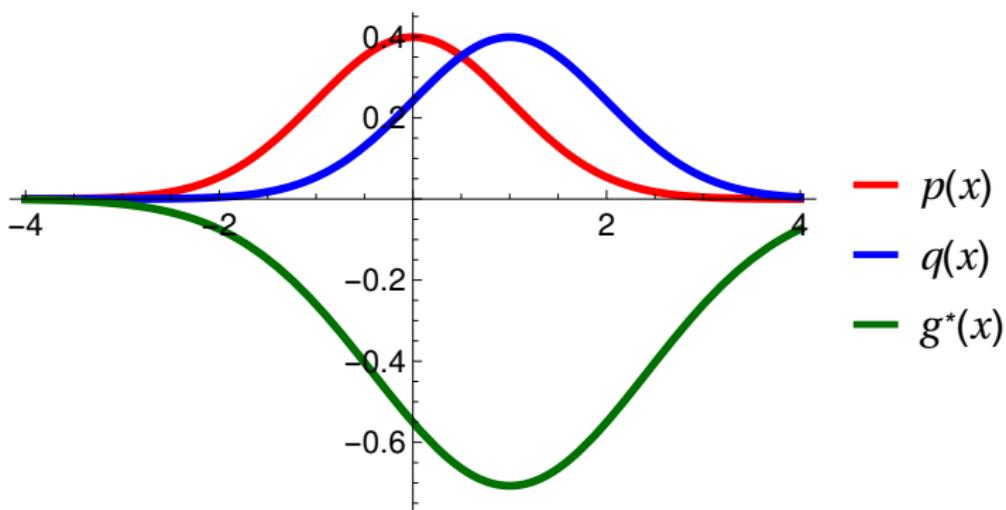
Kernel Stein Discrepancy

Stein operator

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

Kernel Stein Discrepancy (KSD)

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



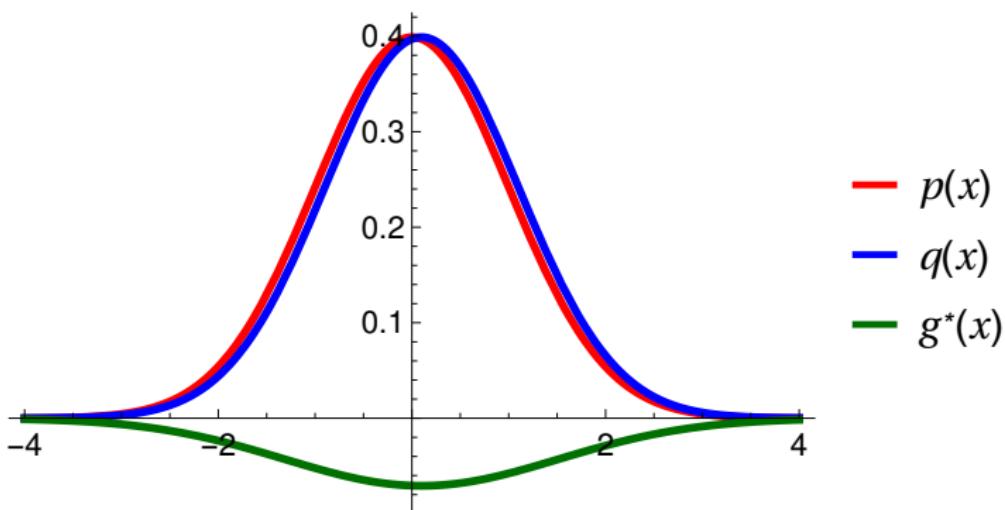
Kernel Stein Discrepancy

Stein operator

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

Kernel Stein Discrepancy (KSD)

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



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

Simple expression using kernels

Re-write stein operator as:

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

Can we define “Stein features”?

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

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

The kernel trick for derivatives

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

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

The kernel trick for derivatives

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

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

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 $\mathbf{x}, \mathbf{x}' \sim Q$, then

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

Kernel stein discrepancy: derivation

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

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

Kernel stein discrepancy: derivation

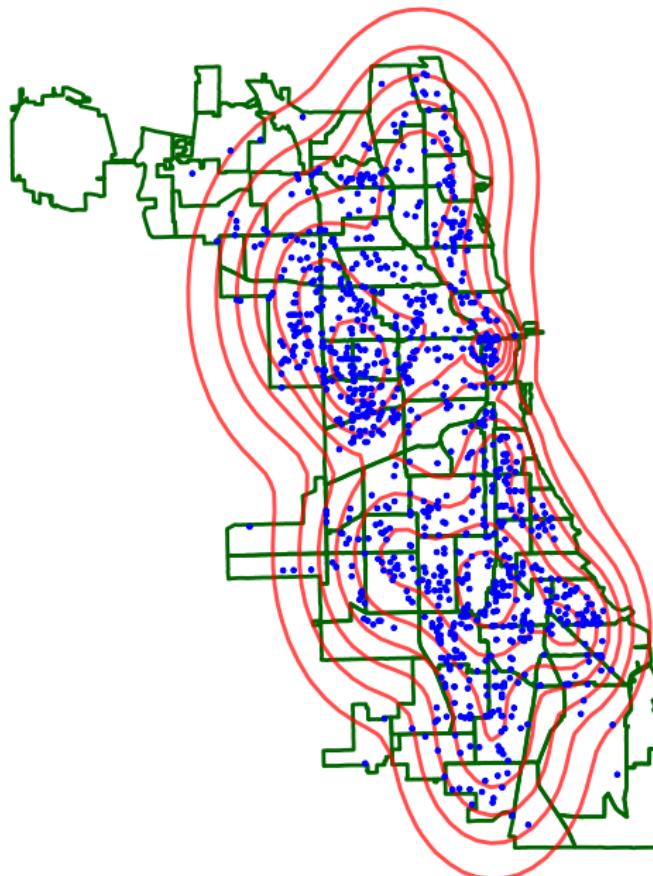
Closed-form expression for KSD: given independent $\mathbf{x}, \mathbf{x}' \sim Q$, then

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

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

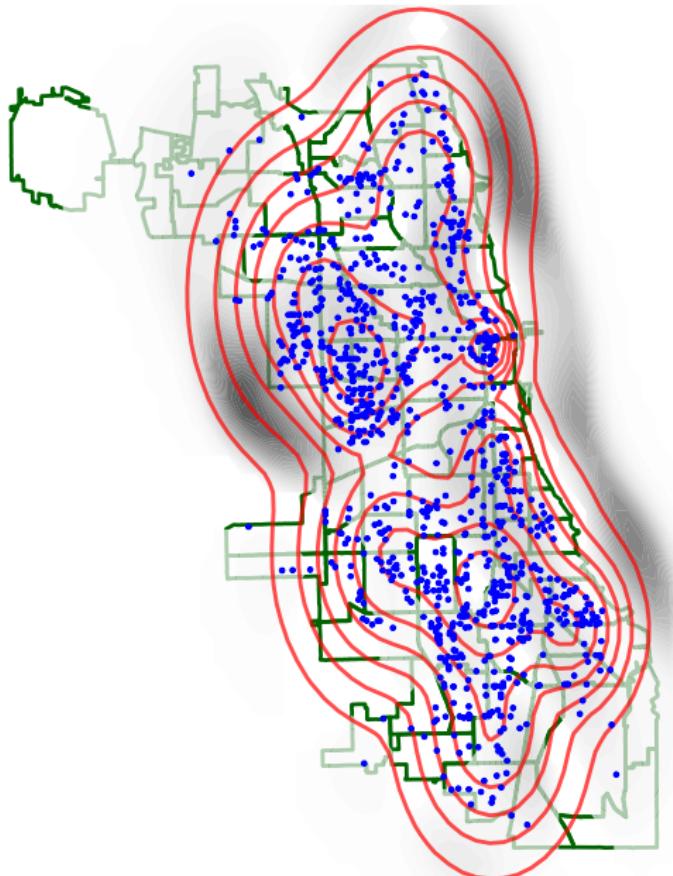
$$\mathbf{E}_{\mathbf{x} \sim q} \left(\frac{d}{dx} \log p(\mathbf{x}) \right)^2 < \infty.$$

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**,

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

Then

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

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

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

is undefined.

Does the Riesz 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

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

is undefined.

Kernel stein discrepancy: population expression

Test statistic:

$$\text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{Q}) = \|\mathbf{E}_{x \sim q} \xi_x\|_{\mathcal{F}}^2 = \mathbf{E}_{x, x' \sim Q} h_{\textcolor{red}{p}}(x, x')$$

where

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

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

Kernel stein discrepancy: population expression

Test statistic:

$$\text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{Q}) = \|\mathbf{E}_{x \sim q} \xi_x\|_{\mathcal{F}}^2 = \mathbf{E}_{x, x' \sim Q} h_{\textcolor{red}{p}}(x, x')$$

where

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

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

Kernel stein discrepancy: population expression

Test statistic:

$$\text{KSD}_{\textcolor{red}{p}}^2(\mathcal{Q}) = \|\mathbf{E}_{x \sim q} \xi_x\|_{\mathcal{F}}^2 = \mathbf{E}_{x, x' \sim Q} h_{\textcolor{red}{p}}(x, x')$$

where

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

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

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

Kernel stein discrepancy: population expression

Test statistic:

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

where

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

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

If kernel is C_0 -universal and \mathbf{Q} satisfies $\mathbf{E}_{x \sim \mathbf{Q}} \left\| \nabla \left(\log \frac{\mathbf{p}(x)}{\mathbf{q}(x)} \right) \right\|^2 < \infty$,
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):

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

where

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

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

KSD for discrete-valued variables

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

The population KSD (discrete):

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

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):

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

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:

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

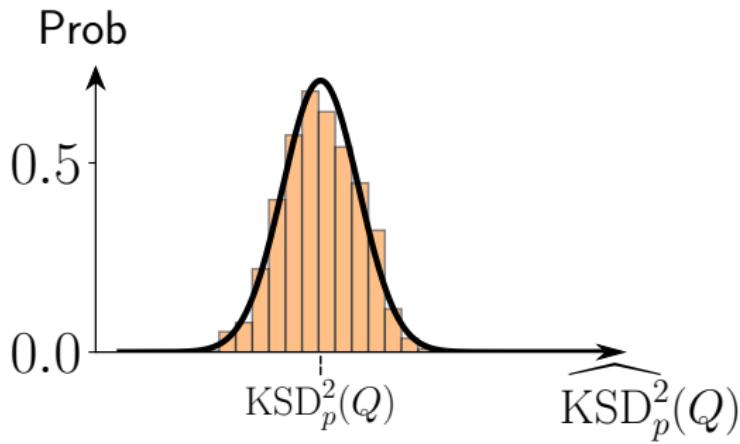
Empirical statistic, asymptotic normality for $P \neq Q$

The empirical statistic:

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

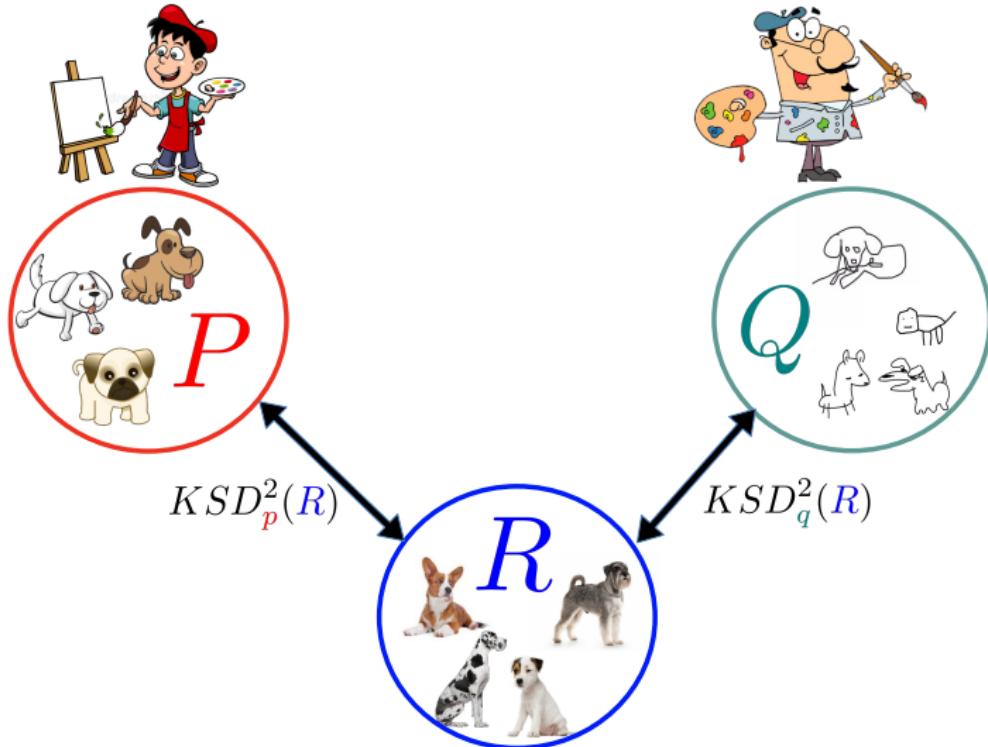
Asymptotic distribution when $P \neq Q$:

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



Relative goodness-of-fit testing

- Two generative models P and Q , data $\{x_i\}_{i=1}^n \sim R$.
- Neither model gives a perfect fit ($P \neq R$ and $Q \neq R$).

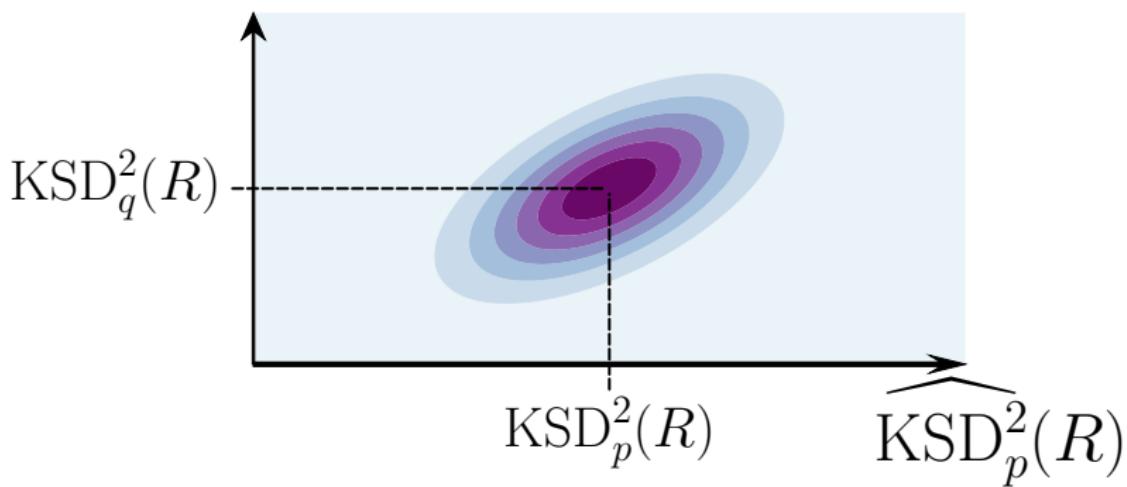


Joint asymptotic normality

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

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

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



Joint asymptotic normality

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

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

Difference in statistics is asymptotically normal:

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

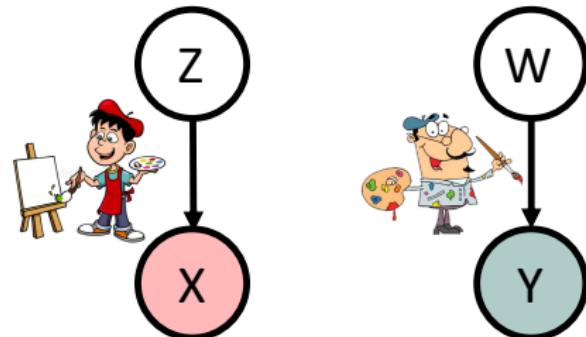
\implies a statistical test with null hypothesis $\text{KSD}_p^2(R) - \text{KSD}_q^2(R) \leq 0$ is straightforward.

Latent variable models

Can we compare latent variable models with KSD?

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

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



Recall multi-dimensional Stein operator:

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

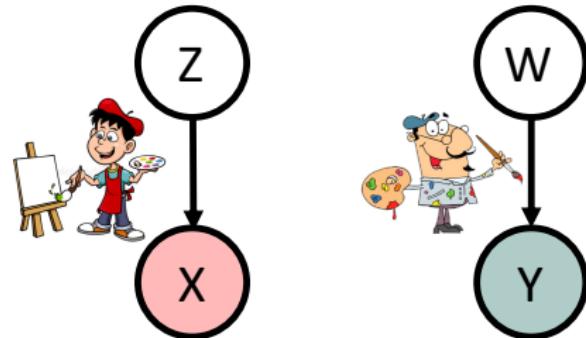
Expression (a) requires marginal $p(x)$, often intractable...

Latent variable models

Can we compare latent variable models with KSD?

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

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



Recall multi-dimensional Stein operator:

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

Expression (a) requires marginal $p(x)$, often intractable...
...but sampling can be straightforward!

Monte Carlo approximation

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 KSDs with approximate densities:

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

Monte Carlo approximation

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 KSDs with approximate densities:

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

Recall

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

→ if m is large, can we simply substitute $\textcolor{red}{p}_m$ and $\textcolor{teal}{q}_m$?

Simple proof of concept

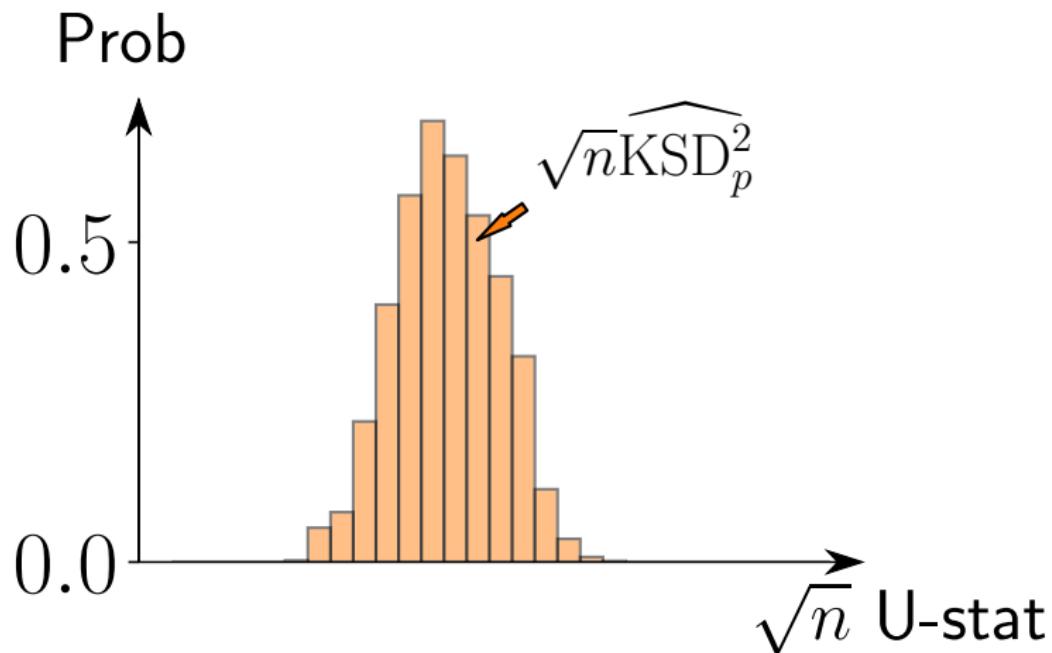
Check $\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) \approx \widehat{\text{KSD}}_{\textcolor{red}{p}_m}^2(\textcolor{blue}{R})$ with a toy model:

- Model: Beta-Binomial $\text{BetaBinom}(\alpha, \beta)$

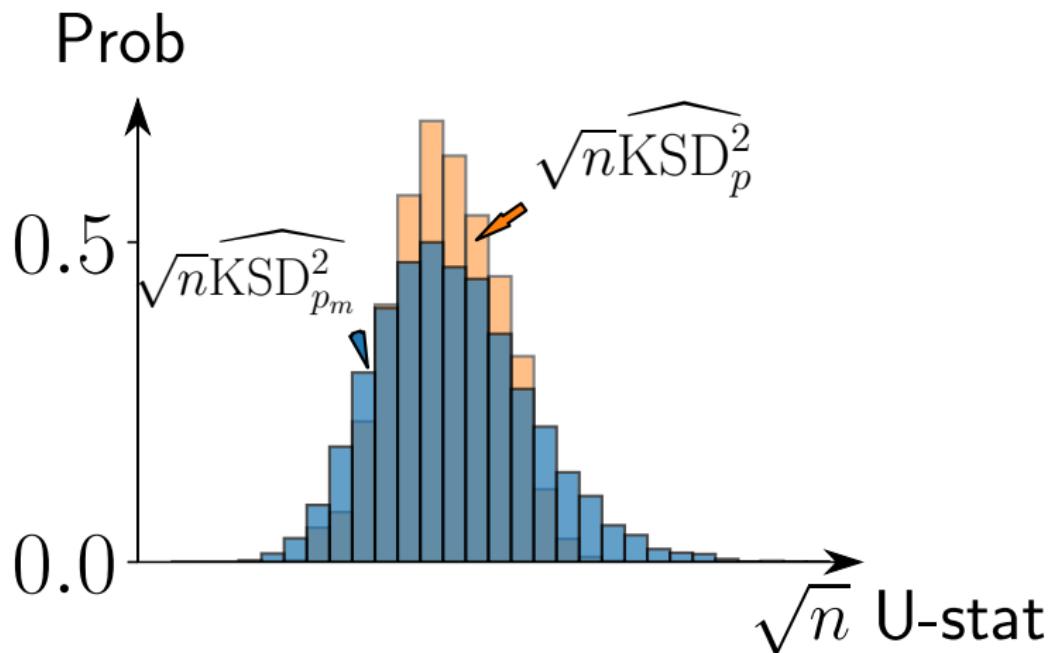
$$\textcolor{red}{p}(x|z) = \binom{N}{x} z^x (1-z)^{n-x}, \quad \textcolor{red}{p}(z) = \text{Beta}(a, b)$$

- Latent $z \in (0, 1)$: success probability for binomial likelihood
 - Marginal $\textcolor{red}{p}(x)$: tractable (given by the beta function)
- Generate $\sqrt{n}\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{R})$ and $\sqrt{n}\widehat{\text{KSD}}_{\textcolor{red}{p}_m}^2(\textcolor{blue}{R})$
→ what do their distribution look like?

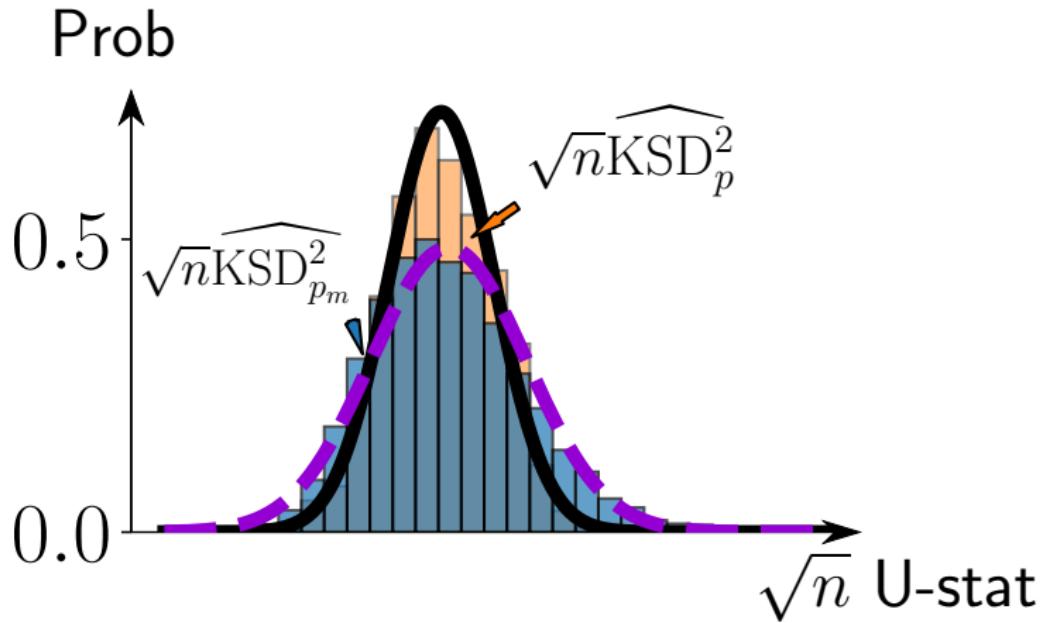
Effect of sampling the latents (Beta-binomial)



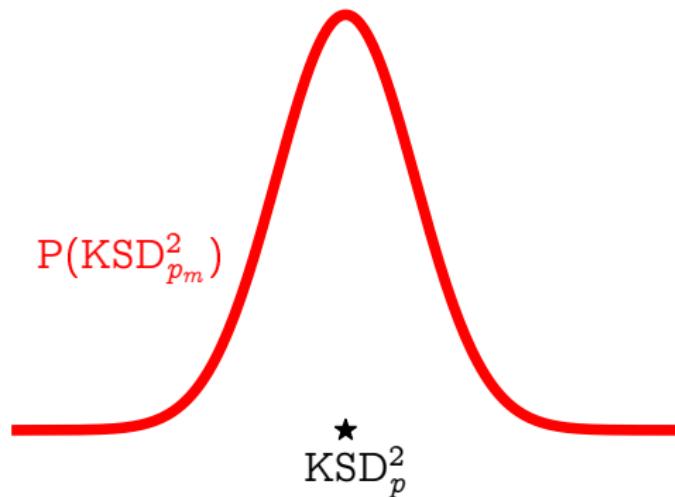
Effect of sampling the latents (Beta-binomial)



Effect of sampling the latents (Beta-binomial)

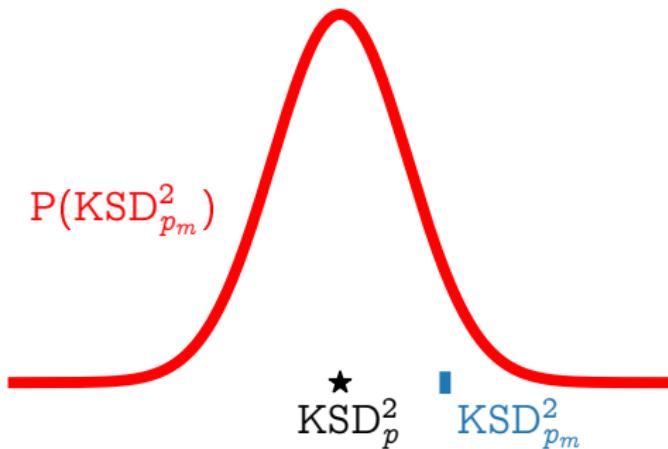


Why this happens



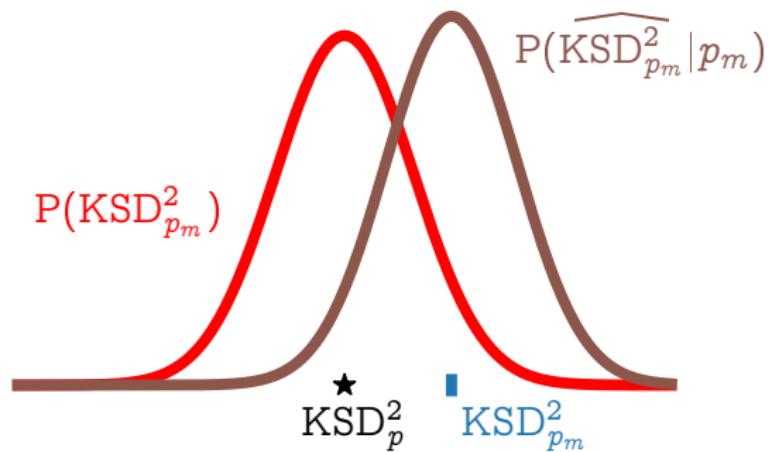
$KSD_{p_m}^2(\textcolor{blue}{R})$ is normally distributed around $KSD_p^2(\textcolor{blue}{R})$
(approximation error)

Why this happens



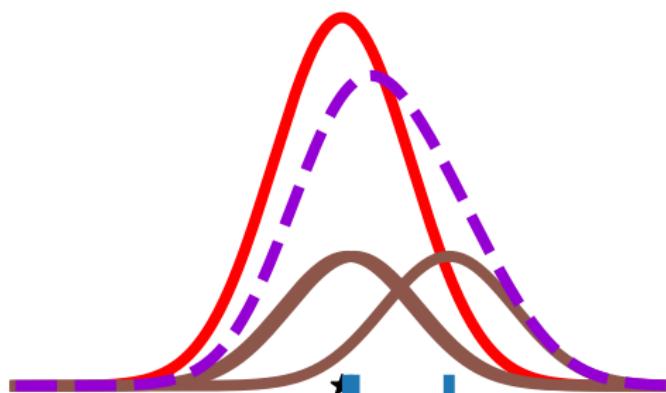
Approximation p_m gives a random draw $KSD_{p_m}^2(R)$

Why this happens



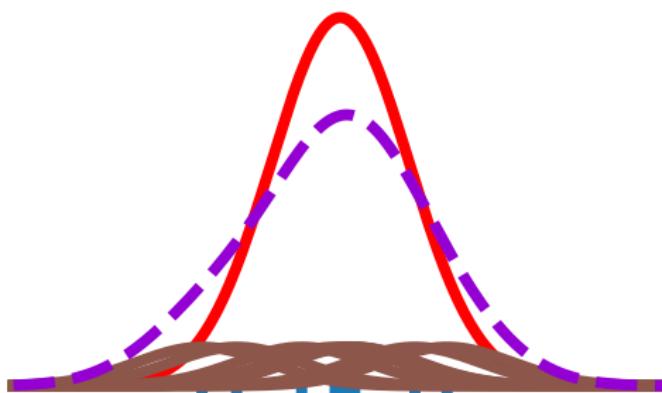
$\widehat{\text{KSD}}_{p_m}^2(R)$ is normally distributed around $\text{KSD}_{p_m}^2(R)$

Why this happens



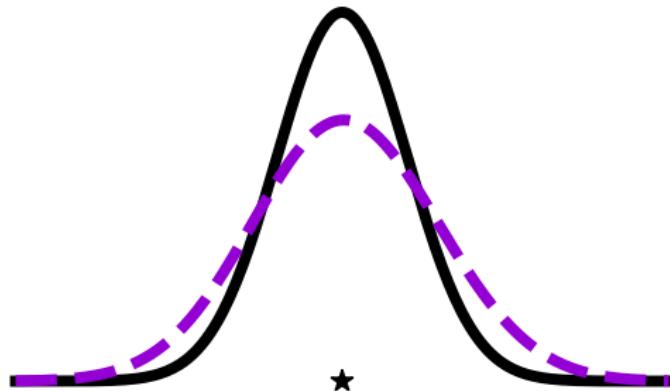
Distribution of $\widehat{\text{KSD}}_{p_m}^2(R)$ is
averaged over random draws of $\text{KSD}_{p_m}^2(R)$

Why this happens



Distribution of $\widehat{\text{KSD}}_{p_m}^2(R)$ is
averaged over random draws of $\text{KSD}_{p_m}^2(R)$

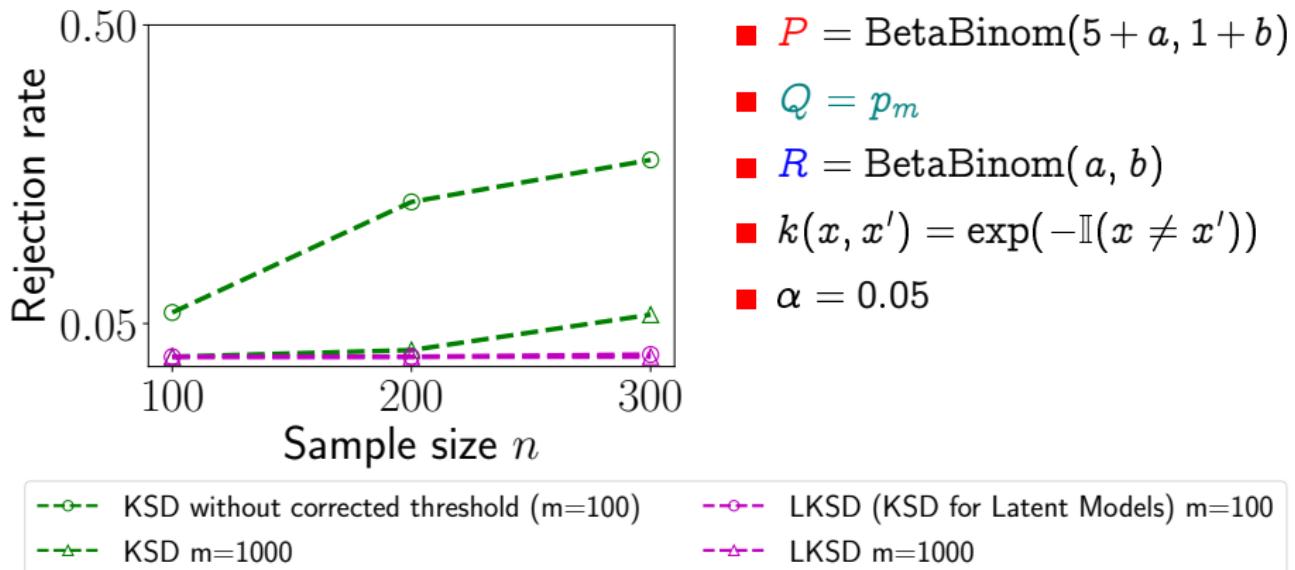
Why this happens



$\widehat{\text{KSD}}_{p_m}^2(R)$ has a higher variance than $\widehat{\text{KSD}}_p^2(R)$

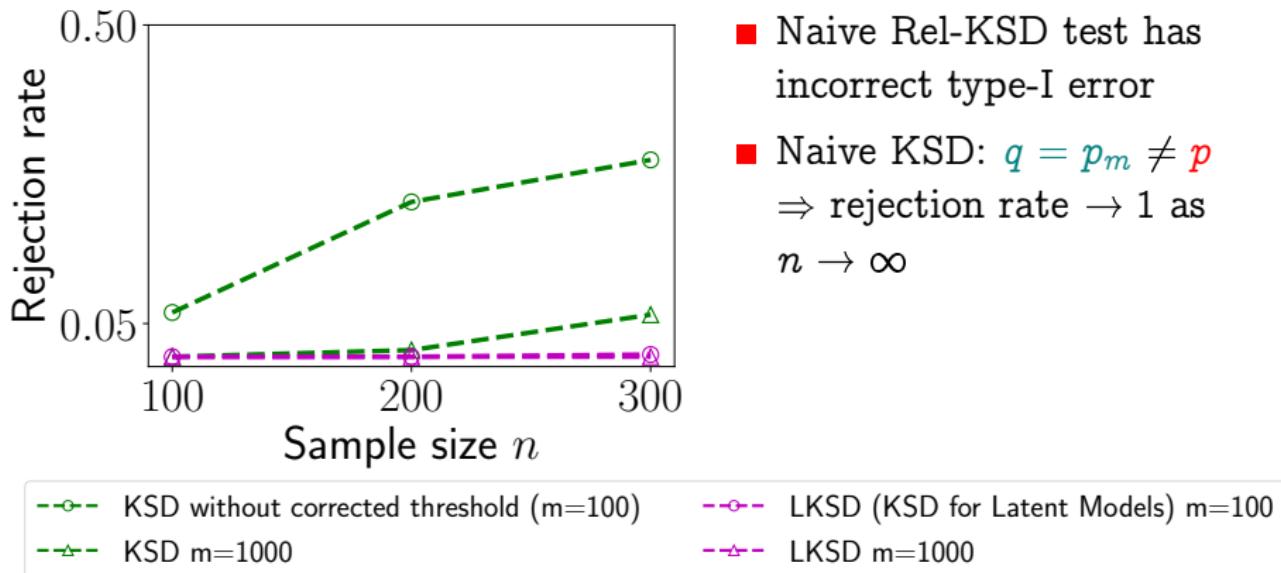
Correction for this effect

- BetaBinomial models with p vs $q = p_m$: numerical vs closed-form marginalisation.
- With correction for increased $\widehat{\text{KSD}}_{p_m}^2(R)$ variance, null accepted w.p. $1 - \alpha$.



Correction for this effect

- BetaBinomial models with p vs $q = p_m$: numerical vs closed-form marginalisation.
- With correction for increased $\widehat{\text{KSD}}_{p_m}^2(R)$ variance, null accepted w.p. $1 - \alpha$.



Asymptotics for approximate KSD

We have asymptotic normality for $\text{KSD}_{\textcolor{red}{p}_m}^2(\textcolor{blue}{R})$,

$$\sqrt{m}(\text{KSD}_{\textcolor{red}{p}_m}^2(\textcolor{blue}{R}) - \text{KSD}_{\textcolor{red}{p}}^2(\textcolor{blue}{R})) \xrightarrow{d} \mathcal{N}(0, \gamma_{\textcolor{red}{p}}^2)$$

The fine print:

- $\inf_x \textcolor{red}{p}(x) > 0$
- $\sup_x \left| \frac{d\textcolor{red}{p}(x)}{dx} \right| < \infty$
- (Uniform CLT) Likelihoods $\{\textcolor{red}{p}(x|\cdot) | x \in \mathcal{X}\}$ and derivatives $\{\frac{d}{dx} \textcolor{red}{p}(x|\cdot) | x \in \mathcal{X}\}$ are $\textcolor{red}{p}(z)$ - Donsker class

Asymptotic distribution for relative KSD test

Asymptotic distribution of approximate KSD estimate

$(n, m) \rightarrow \infty, \frac{n}{m} \rightarrow r \in [0, \infty)$:

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

where

- $c = \sigma_{\textcolor{red}{p}} \textcolor{teal}{q} \sqrt{1 + r(\gamma_{\textcolor{red}{p}} \textcolor{teal}{q} / \sigma_{\textcolor{red}{p}} \textcolor{teal}{q})^2}$
- $\gamma_{\textcolor{red}{p}}^2 = \lim_{m \rightarrow \infty} m \cdot \text{Var} [\mathbf{E}_{\mathbf{x}, \mathbf{x}'} h_{\textcolor{red}{p}}(\mathbf{x}, \mathbf{x}') - \mathbf{E}_{\mathbf{x}, \mathbf{x}'} h_{\textcolor{teal}{q}}(\mathbf{x}, \mathbf{x}')$
- $\sigma_{\textcolor{red}{p}}^2 = \lim_{n \rightarrow \infty} n \cdot \text{Var} \left[\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{R}) - \widehat{\text{KSD}}_{\textcolor{teal}{q}}^2(\textcolor{blue}{R}) \right]$

Fine print:

- $h_{\textcolor{red}{p}}(\mathbf{x}, \mathbf{x}') - h_{\textcolor{teal}{q}}(\mathbf{x}, \mathbf{x}')$ has a finite third moment
- An additional technical condition (next slide)

Main theorem

Theorem (Asymptotic distribution of random kernel U-statistic)

■ Let

- $U_{n,m}$: a U-statistic defined by a random U-statistic kernel H_m
- U_n : a U-statistic defined by a fixed U-statistic kernel h

■ Assume that

- $\sigma_{H_m}^2 \rightarrow \sigma_h^2$ in probability
- $\nu_3(H_m) \rightarrow \nu_3(h) < \infty$ in probability
where $\nu_3(H_m) = \mathbb{E}_{x,x'} |H_m(x, x') - \mathbb{E}_{x,x'} H_m(x, x')|^3$
- $Y_m := \sqrt{m} \left(\mathbb{E}_n[U_{n,m}|H_m] - \mathbb{E}_n[U_n] \right) \xrightarrow{d} Y$

■ Then, with $n/m \rightarrow r \in [0, \infty)$,

$$\lim_{n,m \rightarrow \infty} \Pr \left[\sqrt{n}(U_{n,m} - \mathbb{E}_n U_n) < t \right] = \mathbb{E}_Y \left[\Phi \left(\frac{t - \sqrt{r} Y}{\sigma_h} \right) \right]$$

Main theorem

Theorem (Asymptotic distribution of random kernel U-statistic)

■ Let

- $U_{n,m}$: a U-statistic defined by a random U-statistic kernel H_m
- U_n : a U-statistic defined by a fixed U-statistic kernel h

■ Assume that

- $\sigma_{H_m}^2 \rightarrow \sigma_h^2$ in probability
- $\nu_3(H_m) \rightarrow \nu_3(h) < \infty$ in probability
where $\nu_3(H_m) = \mathbb{E}_{x,x'} |H_m(x, x') - \mathbb{E}_{x,x'} H_m(x, x')|^3$
- $Y_m := \sqrt{m} \left(\mathbb{E}_n[U_{n,m}|H_m] - \mathbb{E}_n[U_n] \right) \xrightarrow{d} Y$

■ Then, with $n/m \rightarrow r \in [0, \infty)$,

$$\lim_{n,m \rightarrow \infty} \Pr \left[\sqrt{n}(U_{n,m} - \mathbb{E}_n U_n) < t \right] = \mathbb{E}_Y \left[\Phi \left(\frac{t - \sqrt{r} Y}{\sigma_h} \right) \right]$$

Main theorem

Theorem (Asymptotic distribution of random kernel U-statistic)

■ Let

- $U_{n,m}$: a U-statistic defined by a random U-statistic kernel H_m
- U_n : a U-statistic defined by a fixed U-statistic kernel h

■ Assume that

- $\sigma_{H_m}^2 \rightarrow \sigma_h^2$ in probability
- $\nu_3(H_m) \rightarrow \nu_3(h) < \infty$ in probability
where $\nu_3(H_m) = \mathbb{E}_{x,x'} |H_m(x, x') - \mathbb{E}_{x,x'} H_m(x, x')|^3$
- $Y_m := \sqrt{m} \left(\mathbb{E}_n[U_{n,m}|H_m] - \mathbb{E}_n[U_n] \right) \xrightarrow{d} Y$

■ Then, with $n/m \rightarrow r \in [0, \infty)$,

$$\lim_{n,m \rightarrow \infty} \Pr \left[\sqrt{n}(U_{n,m} - \mathbb{E}_n U_n) < t \right] = \mathbb{E}_Y \left[\Phi \left(\frac{t - \sqrt{r} Y}{\sigma_h} \right) \right]$$

Experiment: sensitivity to model difference

- Data $\textcolor{blue}{R}$ = Sigmoid Belief Network SBN(W):

$$\textcolor{blue}{R}(x|z) = \text{sigmoid}(\mathbf{W}z), \quad \textcolor{blue}{R}(z) = \mathcal{N}(0, I), \quad \mathbf{W} \in \mathbb{R}^{30 \times 10}$$

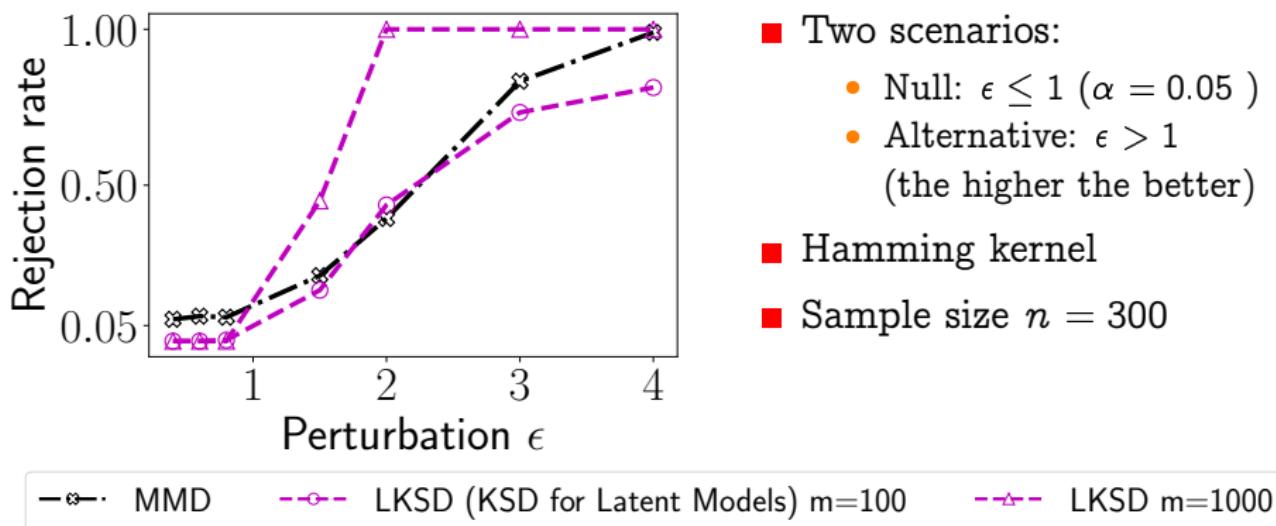
- Models: $\textcolor{red}{P}$ = SBN($W + \epsilon[1, 0, \dots, 0]$), $\textcolor{teal}{Q}$ = SBN($W + [1, 0, \dots, 0]$)
- Only the first column of weight W is perturbed by ϵ

Experiment: sensitivity to model difference

- Data $R = \text{Sigmoid Belief Network SBN}(W)$:

$$R(x|z) = \text{sigmoid}(Wz), \quad R(z) = \mathcal{N}(0, I), \quad W \in \mathbb{R}^{30 \times 10}$$

- Models: $P = \text{SBN}(W + \epsilon[1, 0, \dots, 0])$, $Q = \text{SBN}(W + [1, 0, \dots, 0])$
- Only the first column of weight W is perturbed by ϵ

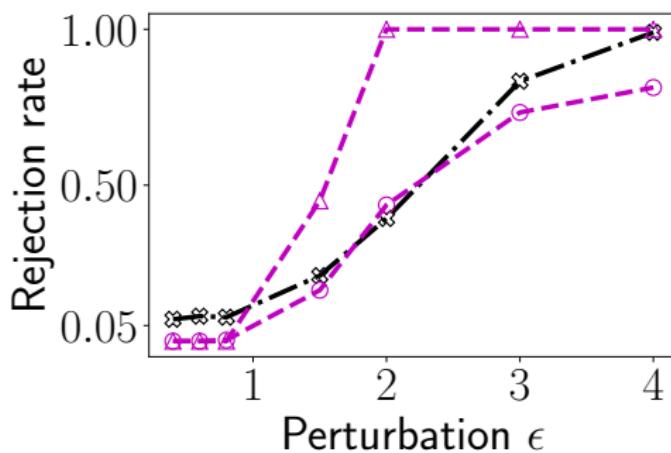


Experiment: sensitivity to model difference

- Data R = Sigmoid Belief Network SBN(W):

$$R(x|z) = \text{sigmoid}(Wz), \quad R(z) = \mathcal{N}(0, I), \quad W \in \mathbb{R}^{30 \times 10}$$

- Models: P = SBN($W + \epsilon[1, 0, \dots, 0]$), Q = SBN($W + [1, 0, \dots, 0]$)
- Only the first column of weight W is perturbed by ϵ



KSD has higher power
($\epsilon > 1$)

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

—○— MMD

-○- LKSD ($m=100$)

-△- LKSD $m=1000$

Questions?

