

# Kernel tests of goodness-of-fit using Stein's method

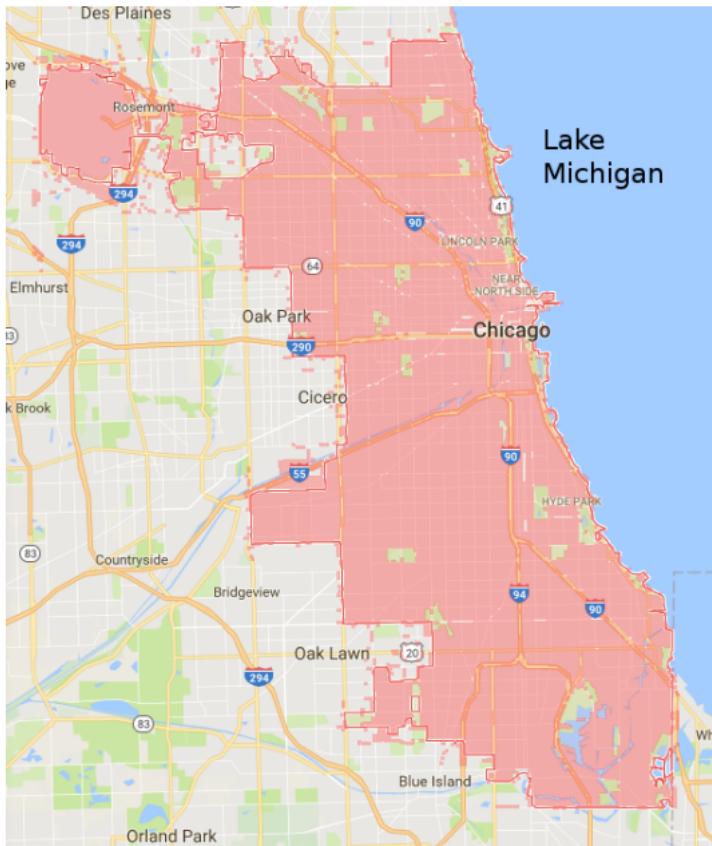


Arthur Gretton  
[gretton@gatsby.ucl.ac.uk](mailto:gretton@gatsby.ucl.ac.uk)

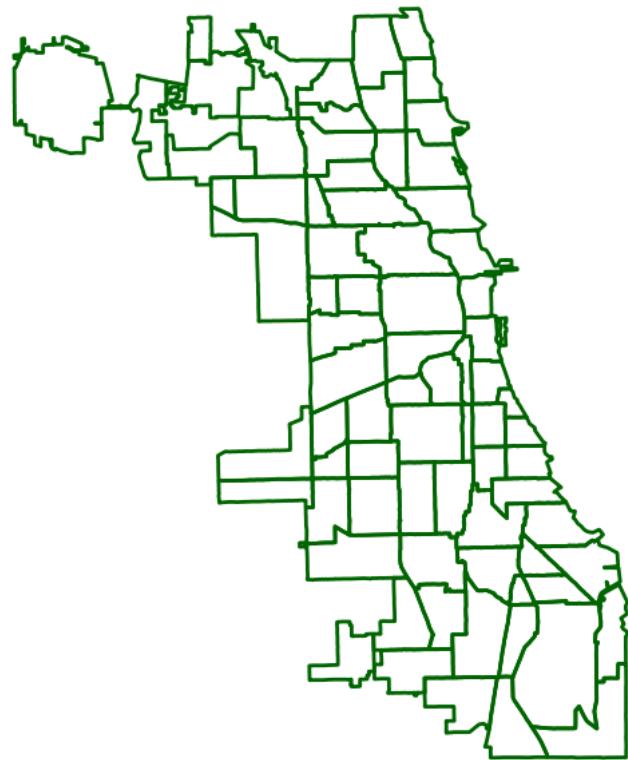
Gatsby Unit, University College London

FIMI 2020

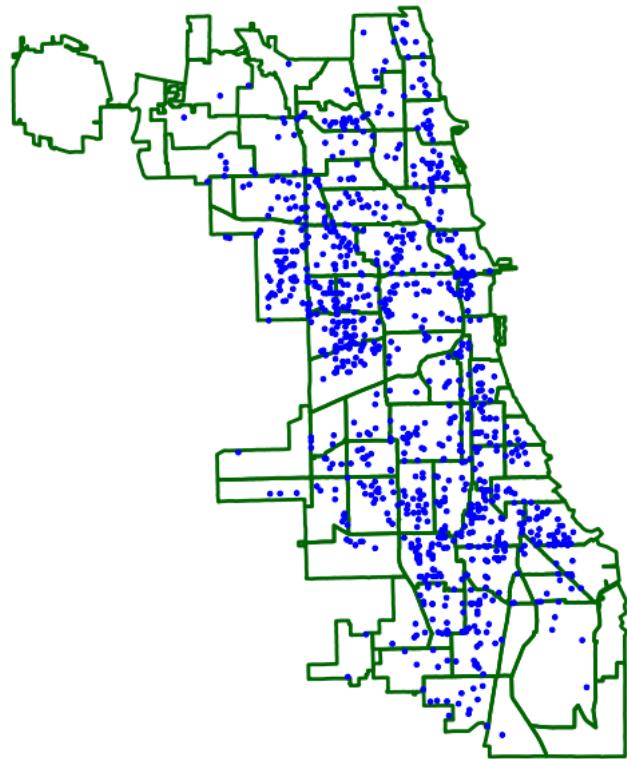
# Model Criticism



## Model Criticism

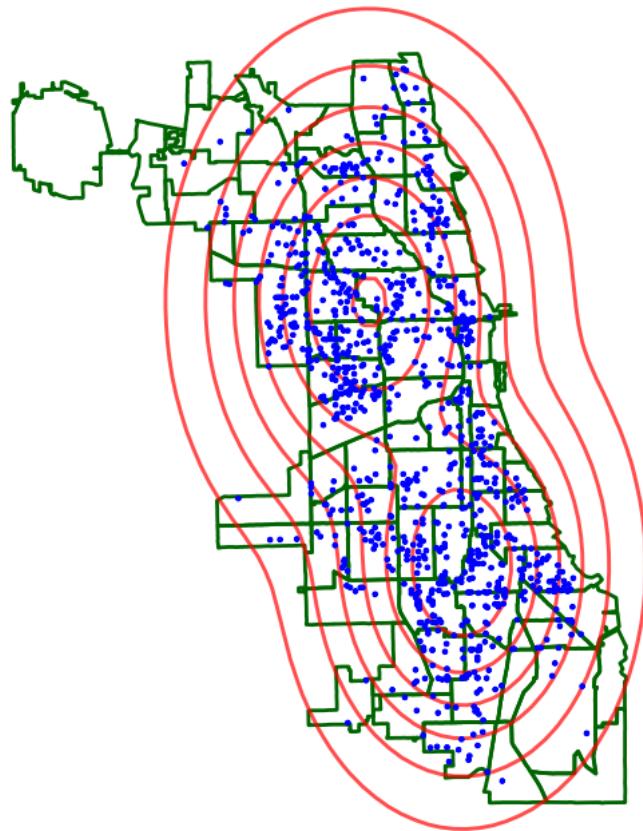


## Model Criticism



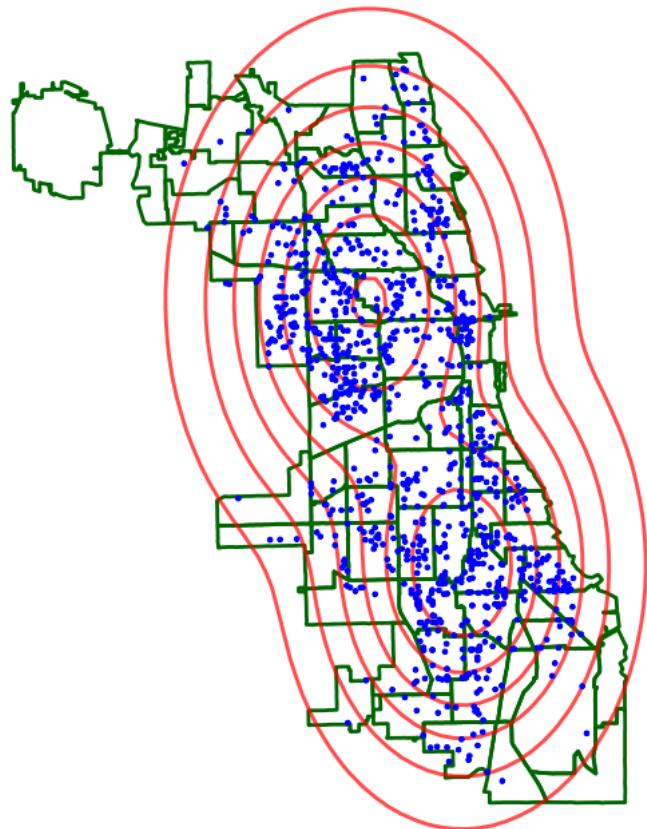
Data = robbery events in Chicago in 2016.

## Model Criticism



Is this a good **model**?

## Model Criticism



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

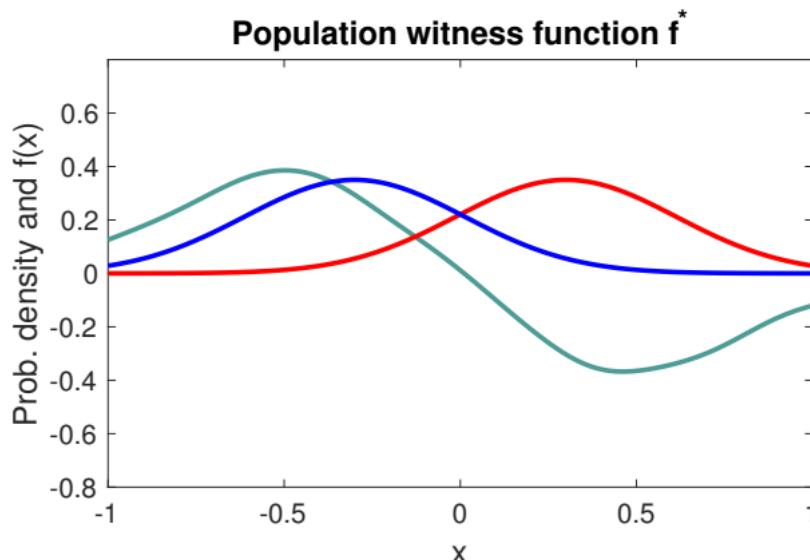
# Outline

- The kernel Stein discrepancy Chwialkowski, Strathmann, G. ICML 2016
  - Comparing two models via samples: MMD and the witness function.
  - Comparing a sample and a model: **Stein** modification of the witness class
- A Linear-Time Kernel Goodness-of-Fit Test Jitkrittum, Xu, Szabo, Fukumizu, G. NeurIPS 2017
  - Features learned to maximise (estimate of) test power
  - Better asymptotic relative efficiency vs a “naive” linear time test
- Relative hypothesis tests with latent variables Kanagawa, Jitkrittum, Mackey, Fukumizu, G. 2019

## The MMD: an integral probability metric

Maximum mean discrepancy: RKHS function for  $P$  vs  $Q$

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



## The MMD: an integral probability metric

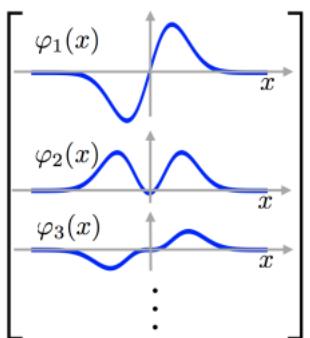
Maximum mean discrepancy: RKHS function for  $P$  vs  $Q$

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

Functions are linear combinations of features:

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

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



## The maximum mean discrepancy

The maximum mean discrepancy in terms of expected kernels:

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

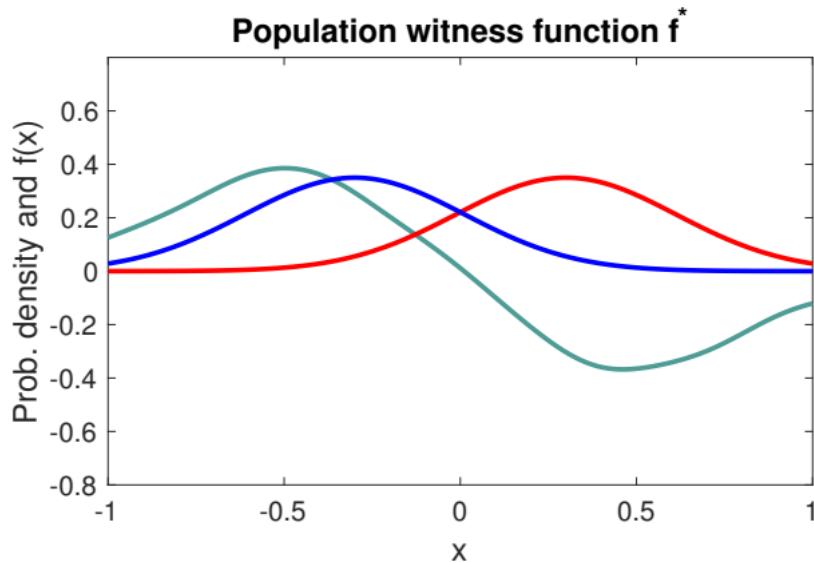
(a)= within distrib. similarity, (b)= cross-distrib. similarity.

$\mu_P$  are expected features,  $k$  is dot product:

$$\mu_P = \mathbf{E}_P \varphi(x) \quad k(x, x') = \langle \varphi(x), \varphi(x') \rangle$$

## Model criticism

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



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

**Problem:** usually can't compute  $\mathbf{E}_p f$  in closed form.

## Stein idea

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

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

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

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

Then

$$\mathbf{E}_p T_p f = 0$$

subject to appropriate boundary conditions.

Proof:

$$\begin{aligned} & \mathbf{E}_p [T_p f] \\ &= \int \left[ \frac{d}{dx} (f(x)p(x)) \right] dx \\ &= [f(x)p(x)]_{-\infty}^{\infty} = 0 \end{aligned}$$

## Stein idea

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

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

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

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

Then

$$\mathbf{E}_p T_p f = 0$$

subject to appropriate boundary conditions.

Proof:

$$\begin{aligned} E_p [T_p f] &= \int \left[ \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x)) \right] p(x) dx \\ &\quad \int \left[ \frac{d}{dx} (f(x)p(x)) \right] dx \\ &= [f(x)p(x)]_{-\infty}^{\infty} = 0 \end{aligned}$$

## Stein idea

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

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

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

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

Then

$$\mathbf{E}_p T_p f = 0$$

subject to appropriate boundary conditions.

Proof:

$$\begin{aligned} E_p [T_p f] & \int \left[ \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x)) \right] p(x) dx \\ & \int \left[ \frac{d}{dx} (f(x)p(x)) \right] dx \\ & = [f(x)p(x)]_{-\infty}^{\infty} = 0 \end{aligned}$$

## Stein idea

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

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

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

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

Then

$$\mathbf{E}_p T_p f = 0$$

subject to appropriate boundary conditions.

Proof:

$$\begin{aligned} E_p [T_p f] & \int \left[ \frac{1}{p(x)} \frac{d}{dx} (f(x)p(x)) \right] p(x) dx \\ & \int \left[ \frac{d}{dx} (f(x)p(x)) \right] dx \\ &= [f(x)p(x)]_{-\infty}^{\infty} = 0 \end{aligned}$$

## Kernel Stein Discrepancy

Stein operator

$$T_{\textcolor{red}{p}} \textcolor{teal}{f} = \frac{1}{p(x)} \frac{d}{dx} (\textcolor{teal}{f}(x) \textcolor{red}{p}(x))$$

Kernel Stein Discrepancy (KSD)

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

## Kernel Stein Discrepancy

Stein operator

$$T_{\textcolor{red}{p}} \textcolor{teal}{f} = \frac{1}{p(x)} \frac{d}{dx} (\textcolor{teal}{f}(x) \textcolor{red}{p}(x))$$

Kernel Stein Discrepancy (KSD)

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

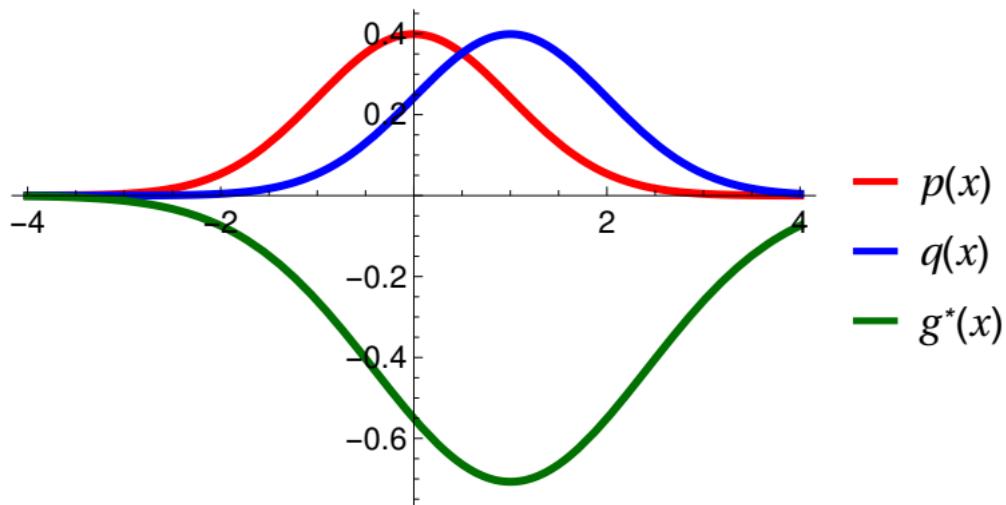
# Kernel Stein Discrepancy

Stein operator

$$T_p f = \frac{1}{p(x)} \frac{d}{dx} (f(x) p(x))$$

Kernel Stein Discrepancy (KSD)

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



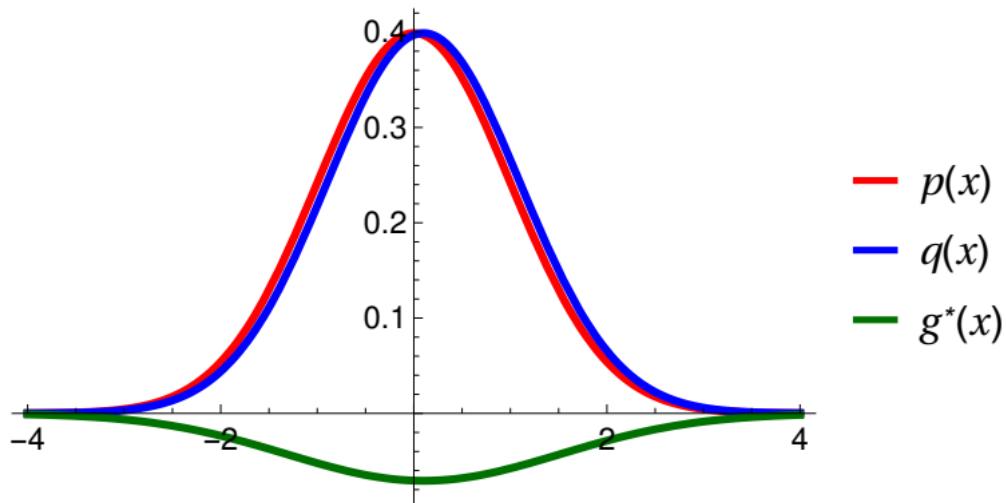
# Kernel Stein Discrepancy

Stein operator

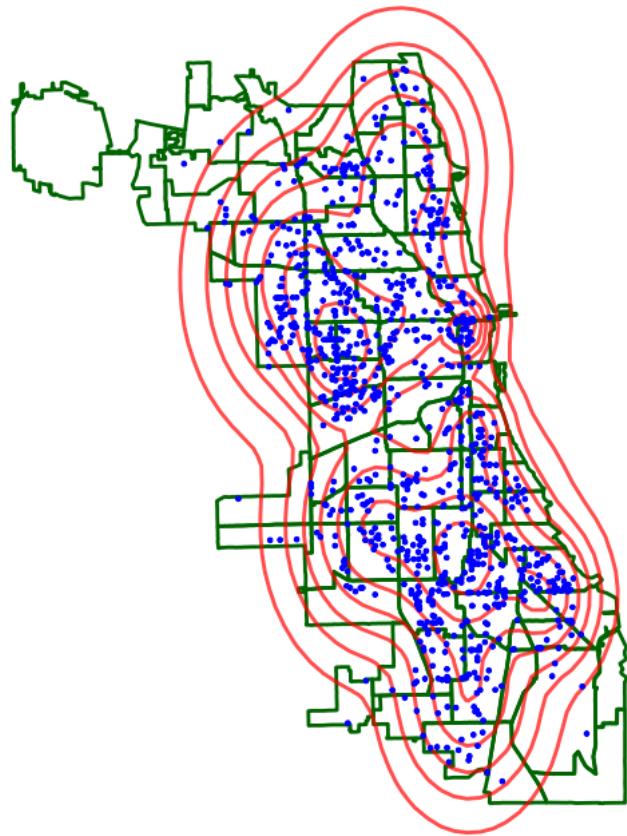
$$T_p f = \frac{1}{p(x)} \frac{d}{dx} (f(x) p(x))$$

Kernel Stein Discrepancy (KSD)

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

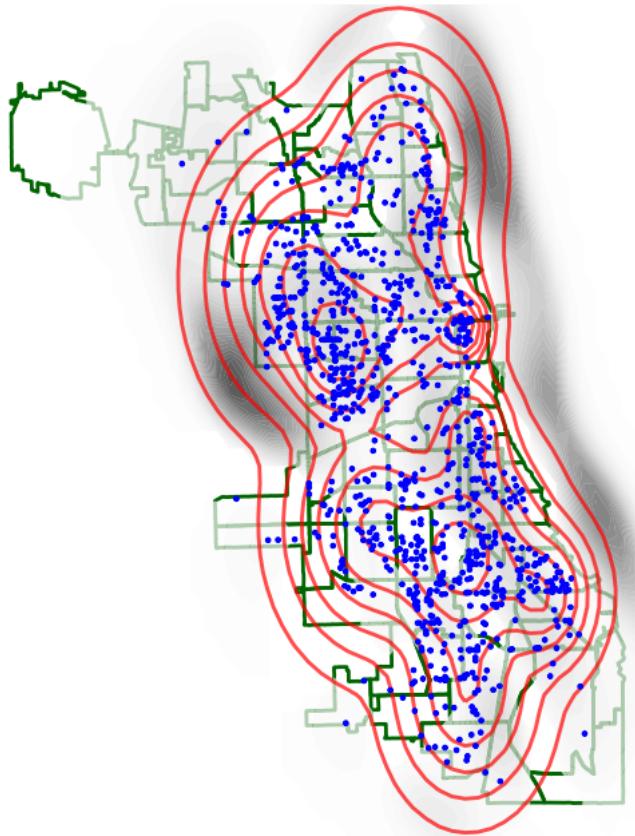


## The witness function: Chicago Crime



Model  $p = 10$ -component Gaussian mixture.

## The witness function: Chicago Crime



Witness function  $g$  shows mismatch

## Simple expression using kernels

Re-write stein operator as:

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

Can we define “Stein features” in  $\mathcal{F}$ ?

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

where  $E_{x \sim p} \xi(x) = 0$ .

## Simple expression using kernels

Re-write stein operator as:

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

Can we define “Stein features” in  $\mathcal{F}$ ?

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

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

## 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),

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

## Kernel stein discrepancy: derivation

Closed-form expression for KSD:

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

## Kernel stein discrepancy: derivation

Closed-form expression for KSD:

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

## Kernel stein discrepancy: derivation

Closed-form expression for KSD:

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

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

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

## Kernel stein discrepancy: population expression

Test statistic when  $x \in \mathbb{R}^d$ , given independent  $\mathbf{x}, \mathbf{x}' \sim q$ ,

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

where

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

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

## Kernel stein discrepancy: population expression

Test statistic when  $x \in \mathbb{R}^d$ , given independent  $\textcolor{blue}{x}, \textcolor{blue}{x}' \sim q$ ,

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

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 when  $x \in \mathbb{R}^d$ , given independent  $\mathbf{x}, \mathbf{x}' \sim q$ ,

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

where

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

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

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

## Kernel stein discrepancy: population expression

Test statistic when  $x \in \mathbb{R}^d$ , given independent  $\mathbf{x}, \mathbf{x}' \sim q$ ,

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

where

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

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

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

## 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

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

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

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

is undefined.

## KSD for discrete-valued variables

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

The population KSD (discrete):

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

where

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

$k_1(x, x') = \Delta_x^{-1} k(x, x')$ ,  $\Delta_x^{-1}$  is cyclic backwards 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):

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

where

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

$k_1(x, x') = \Delta_x^{-1} k(x, x')$ ,  $\Delta_x^{-1}$  is cyclic backwards difference on  $x$ ,

$$\mathbf{s}_{\textcolor{red}{p}}(x) = \frac{\Delta \mathbf{p}(x)}{\mathbf{p}(x)}$$

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):

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

where

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

$k_1(x, x') = \Delta_x^{-1} k(x, x')$ ,  $\Delta_x^{-1}$  is cyclic backwards difference on  $x$ ,

$$\mathbf{s}_{\mathbf{p}}(x) = \frac{\Delta_{\mathbf{p}}(x)}{p(x)}$$

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

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

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

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

## Empirical statistic and asymptotics

The empirical statistic:

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

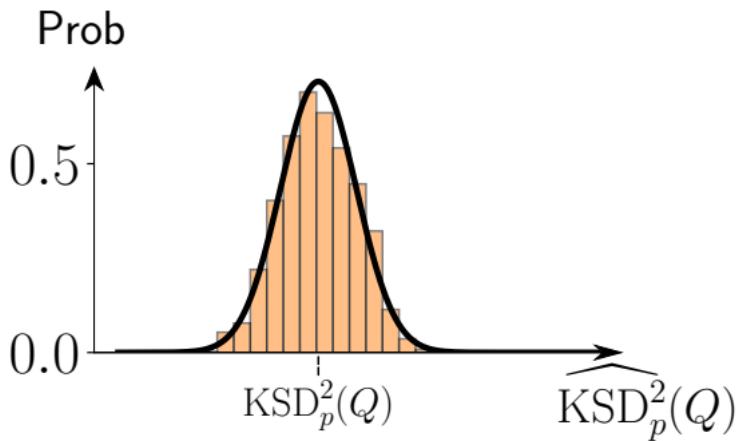
## Empirical statistic and asymptotics

The empirical statistic:

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

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

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



## Empirical statistic and asymptotics

The empirical statistic:

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

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

$$n \widehat{\text{KSD}}_{\textcolor{red}{p}}^2(\textcolor{blue}{Q}) \sim \sum_{\ell=1}^{\infty} \lambda_{\ell} Z_{\ell}^2$$

$$\lambda_i \psi_i(x') = \int_{\mathcal{X}} h_{\textcolor{red}{p}}(x, x') \psi_i(x) d\textcolor{red}{p}(x)$$

$$Z_{\ell} \sim \mathcal{N}(0, 1) \quad \text{i.i.d.}$$

Test threshold via wild bootstrap.

## A naive linear time statistic

A running average:

$$\widehat{LKS}_p^2(Q) := \frac{2}{n} \sum_{i=1}^{n/2} h_p(x_{2i-1}, x_{2i}).$$

Asymptotically normal when  $P \neq Q$  and when  $P = Q$ .

## A naive linear time statistic

A running average:

$$\widehat{LKS_p^2}(Q) := \frac{2}{n} \sum_{i=1}^{n/2} h_p(x_{2i-1}, x_{2i}).$$

Asymptotically normal when  $P \neq Q$  and when  $P = Q$ .

Can we do better? Wishlist:

- 1 still linear-time
- 2 adaptive (parameters automatically tuned)
- 3 more interpretable

# Linear-time, interpretable Goodness-of-fit Test

## Stein Witness Function at a Single Location

Idea:

$$(\text{Stein}) \text{ witness}(\mathbf{v}) = \mathbf{E}_{\mathbf{x} \sim q} [ T_p k_{\mathbf{v}}(\mathbf{x}) ] - \mathbf{E}_{\mathbf{y} \sim p} [ T_p k_{\mathbf{v}}(\mathbf{y}) ]$$

## Stein Witness Function at a Single Location

Idea:

$$(\text{Stein}) \text{ witness}(\mathbf{v}) = \mathbf{E}_{\mathbf{x} \sim q}[T_p] - \mathbf{E}_{\mathbf{y} \sim p}[T_p]$$



## Stein Witness Function at a Single Location

Idea:

$$(\text{Stein}) \text{ witness}(\mathbf{v}) = \mathbf{E}_{\mathbf{x} \sim q} [$$



## Stein Witness Function at a Single Location

Idea:

$$(\text{Stein}) \text{ witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} [$$



$$] - \mathbb{E}_{\mathbf{y} \sim p} [$$

## Stein Witness Function at a Single Location

Idea:

$$(\text{Stein}) \text{ witness}(\mathbf{v}) = \mathbf{E}_{\mathbf{x} \sim q} [$$



## Stein Witness Function at a Single Location

Idea:

$$(\text{Stein}) \text{ witness}(\mathbf{v}) = \mathbf{E}_{\mathbf{x} \sim q} [ T_p k_{\mathbf{v}}(\mathbf{x}) ]$$

## Stein Witness Function at a Single Location

Idea:

$$(\text{Stein}) \text{ witness}(\mathbf{v}) = \mathbf{E}_{\mathbf{x} \sim q} [ T_p k_{\mathbf{v}}(\mathbf{x}) ]$$

Proposal: Good  $\mathbf{v}$  should have high

$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Stein Witness Function at a Single Location

Idea:

$$(\text{Stein}) \text{ witness}(\mathbf{v}) = \mathbf{E}_{\mathbf{x} \sim q}[ T_p k_{\mathbf{v}}(\mathbf{x}) ]$$

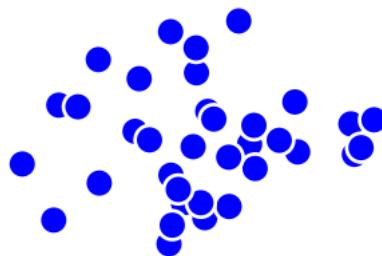
Proposal: Good  $\mathbf{v}$  should have high

signal-to-noise  
ratio

$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

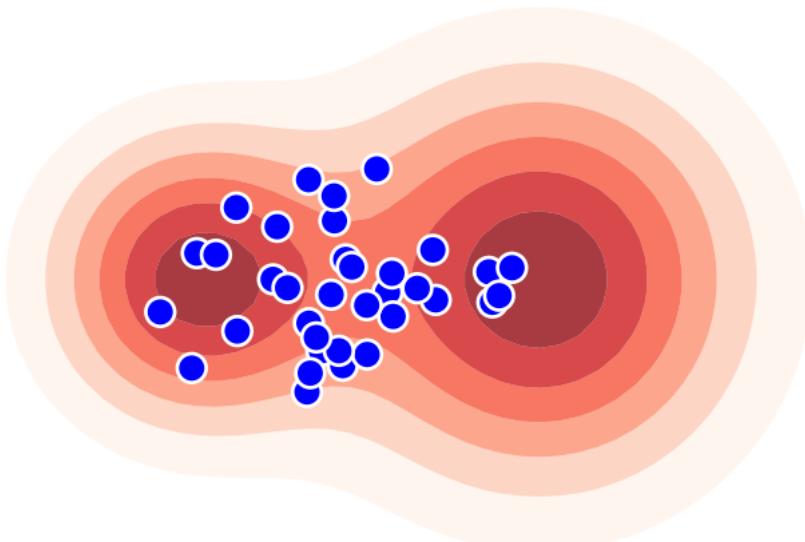
- $\text{witness}(\mathbf{v})$  and  $\text{standard deviation}(\mathbf{v})$  can be estimated in linear time.

## Proposal: Model Criticism with the Stein Witness



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

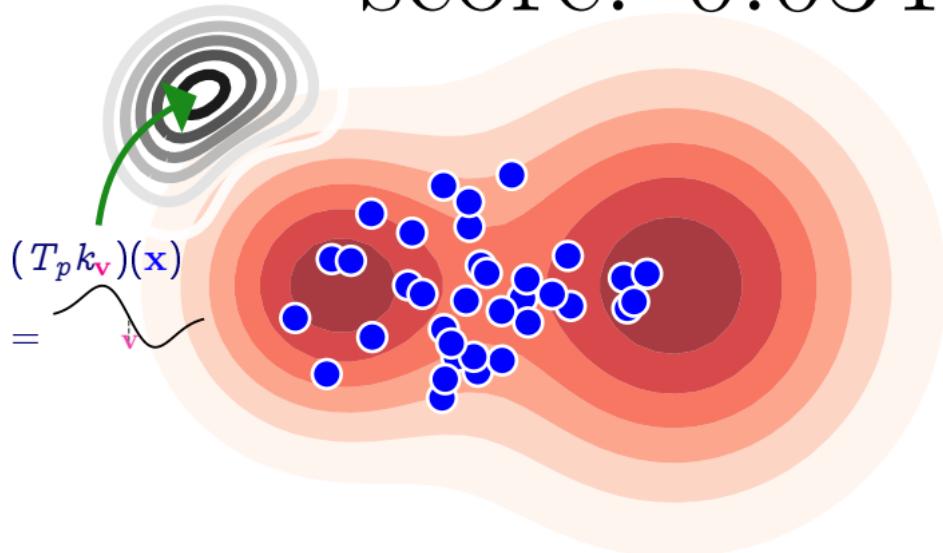
## Proposal: Model Criticism with the Stein Witness



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

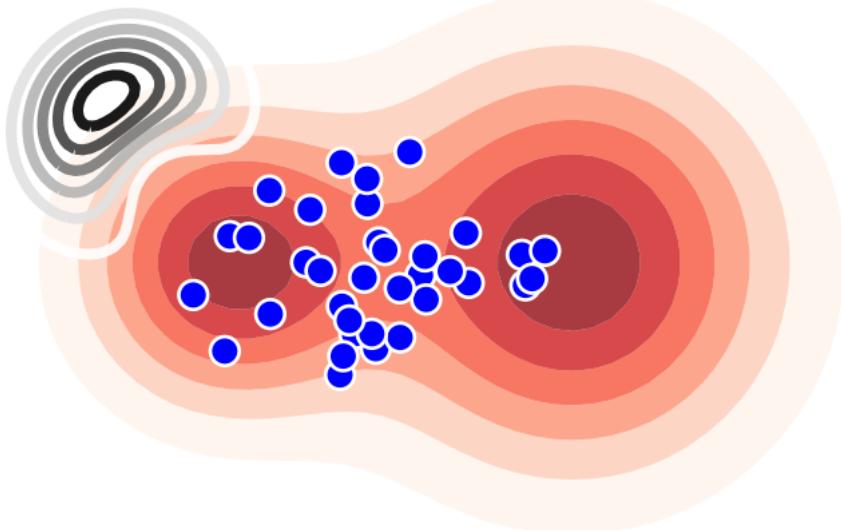
score: 0.034



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

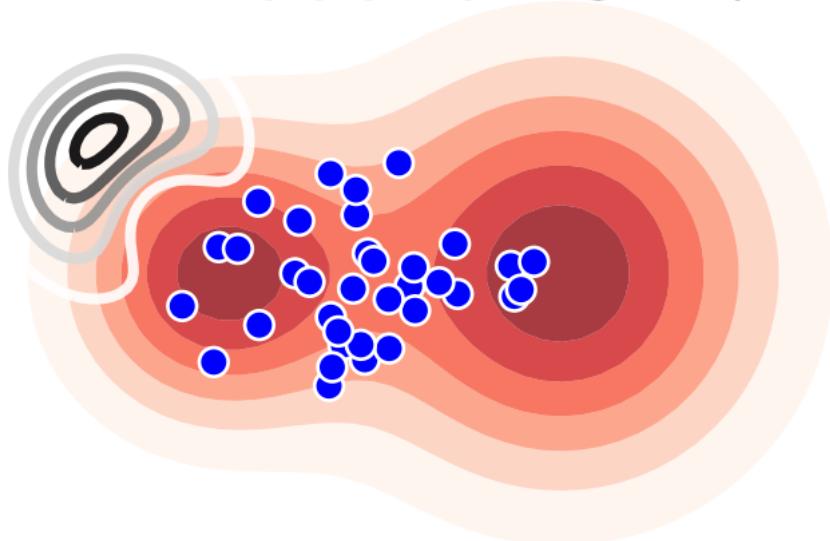
score: 0.089



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

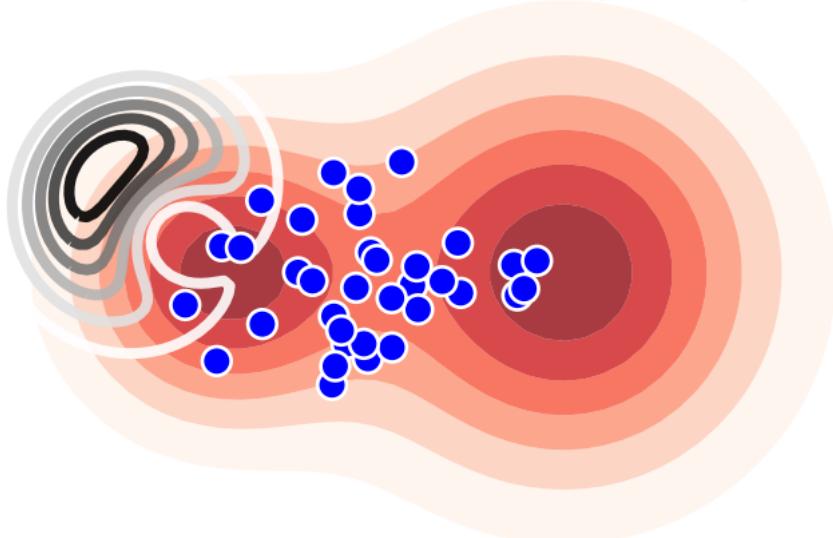
score: 0.17



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

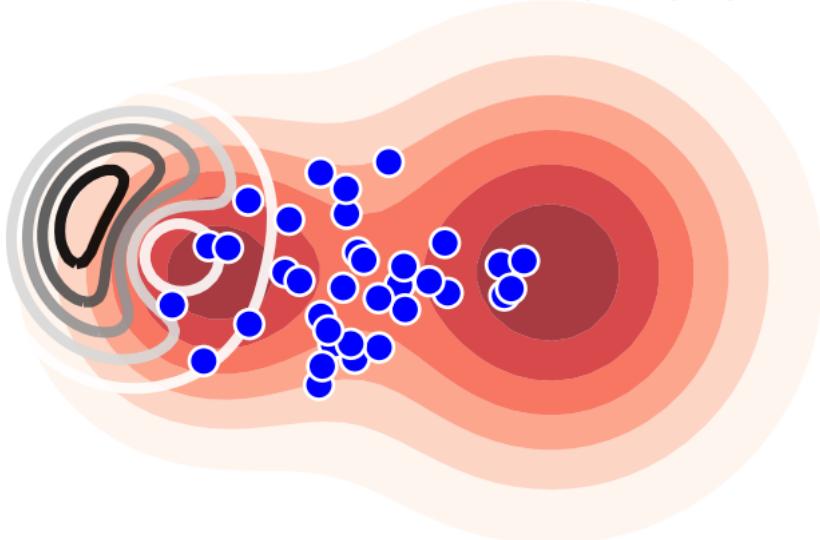
score: 0.26



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

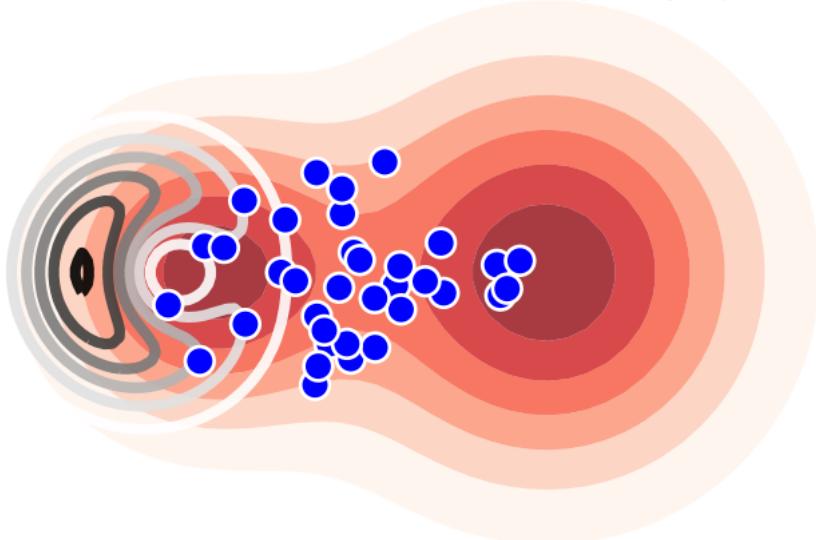
score: 0.33



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

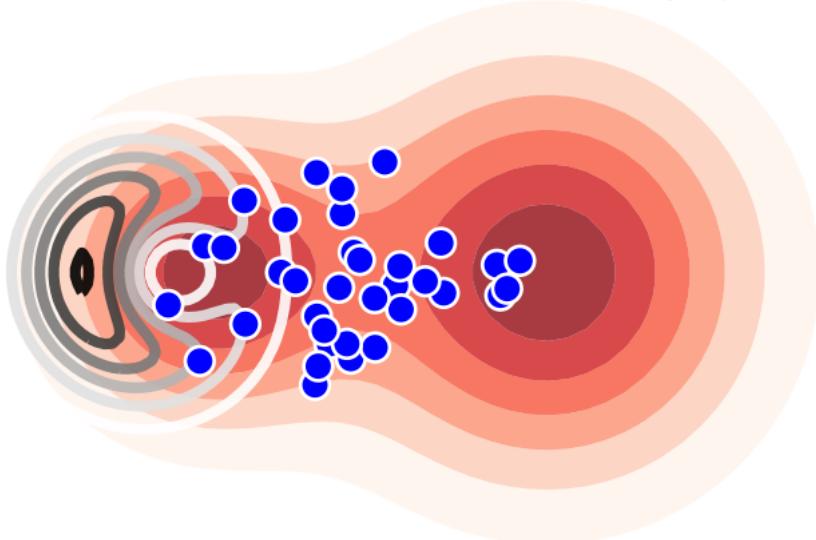
score: 0.37



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

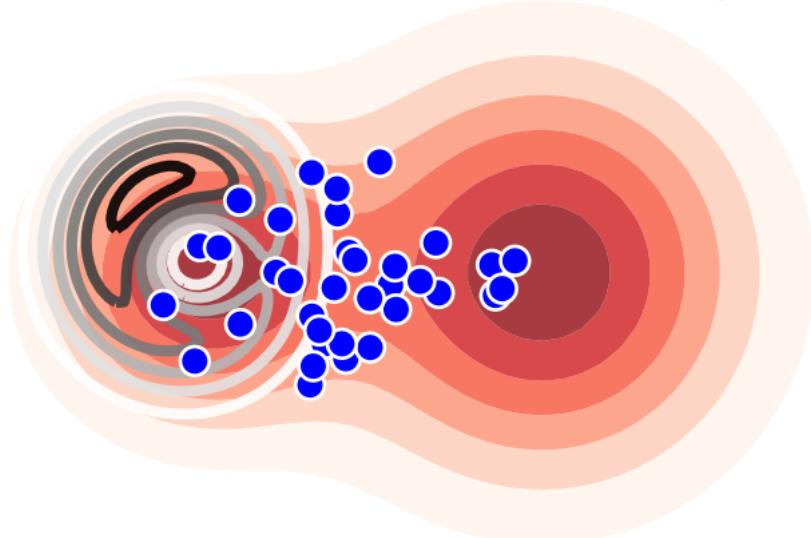
score: 0.37



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

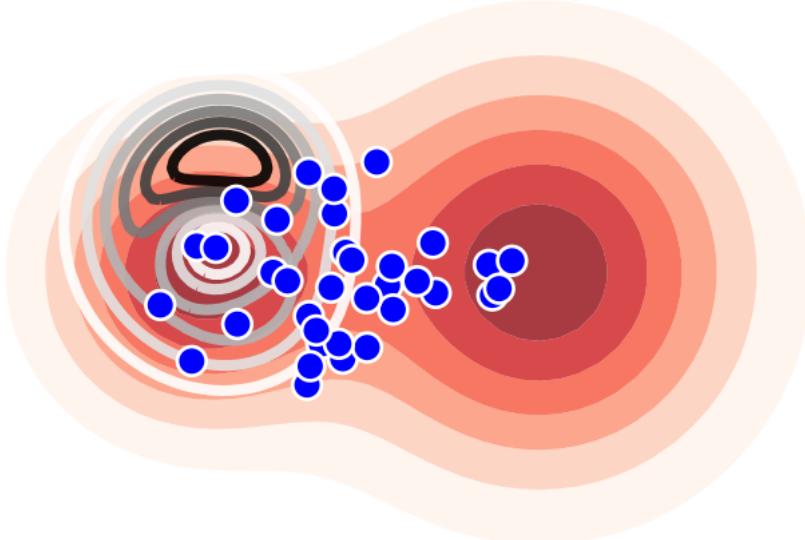
score: 0.45



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

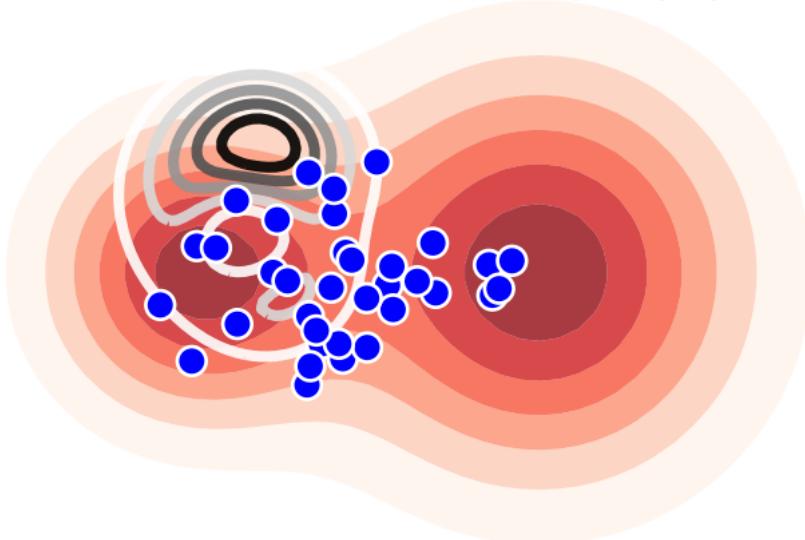
score: 0.44



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

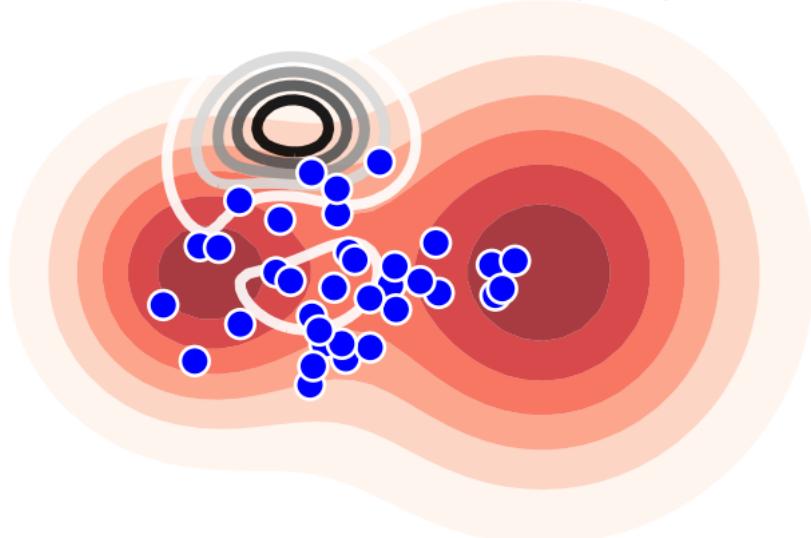
score: 0.39



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

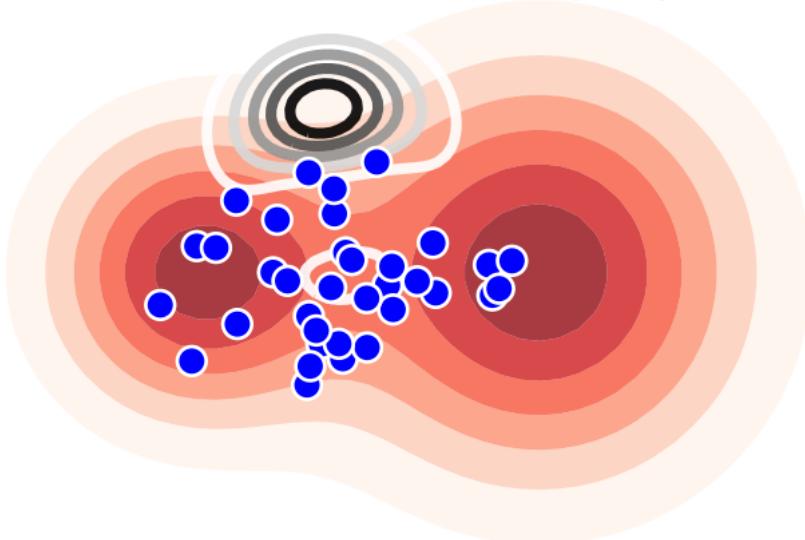
score: 0.31



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

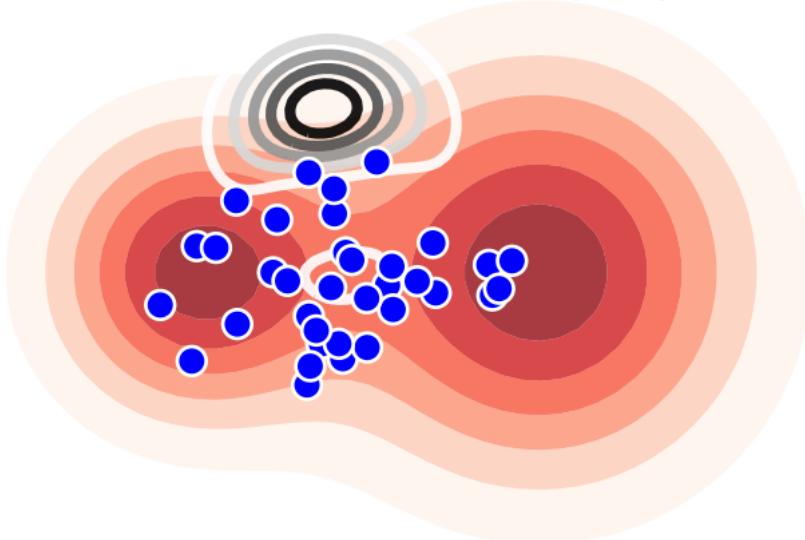
score: 0.32



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

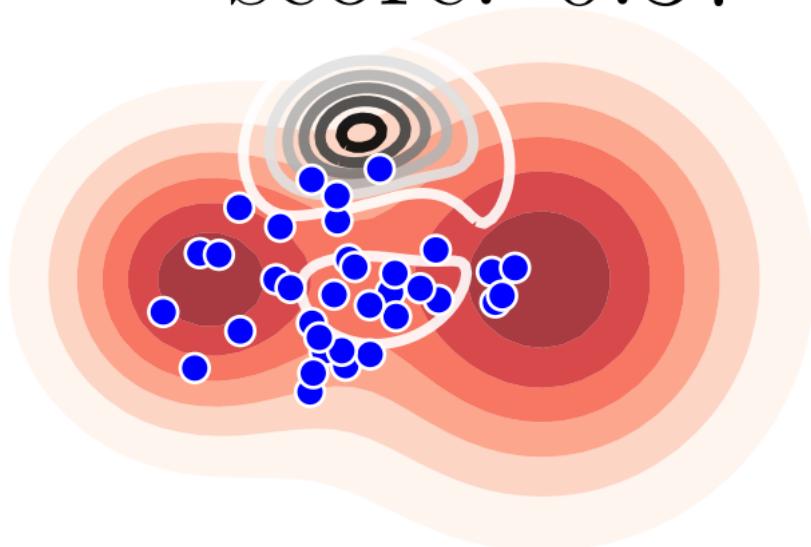
score: 0.32



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

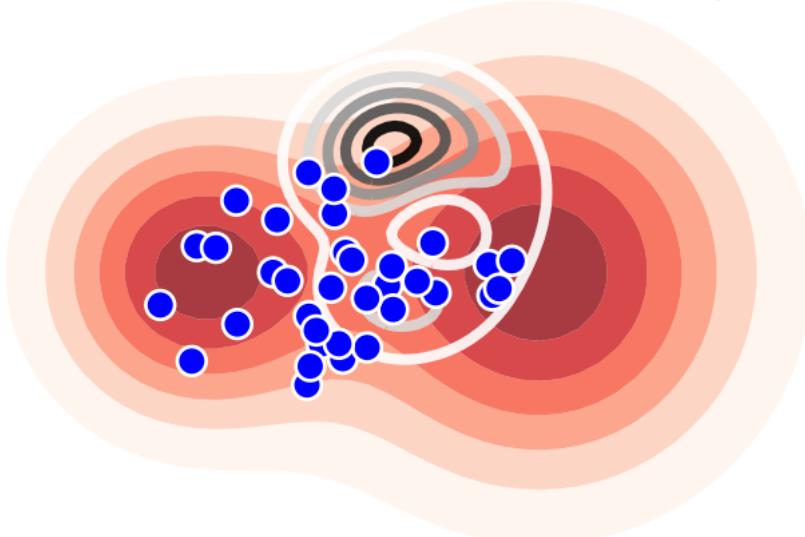
score: 0.37



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

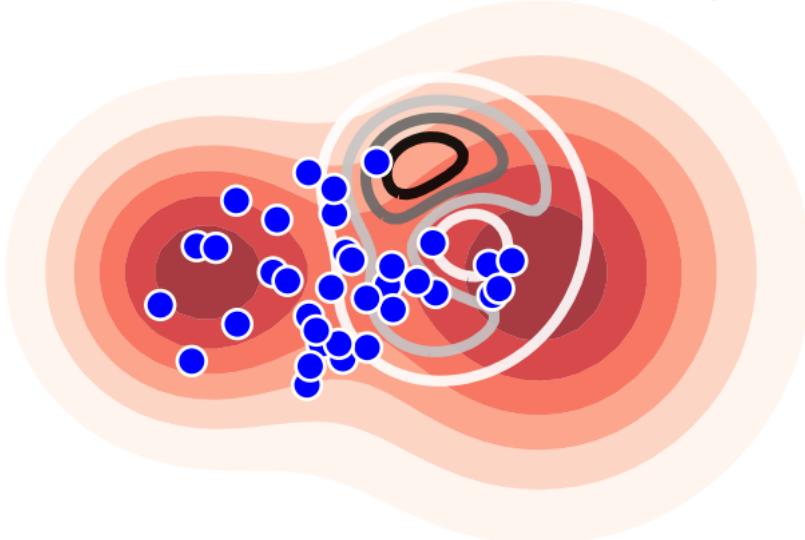
score: 0.48



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

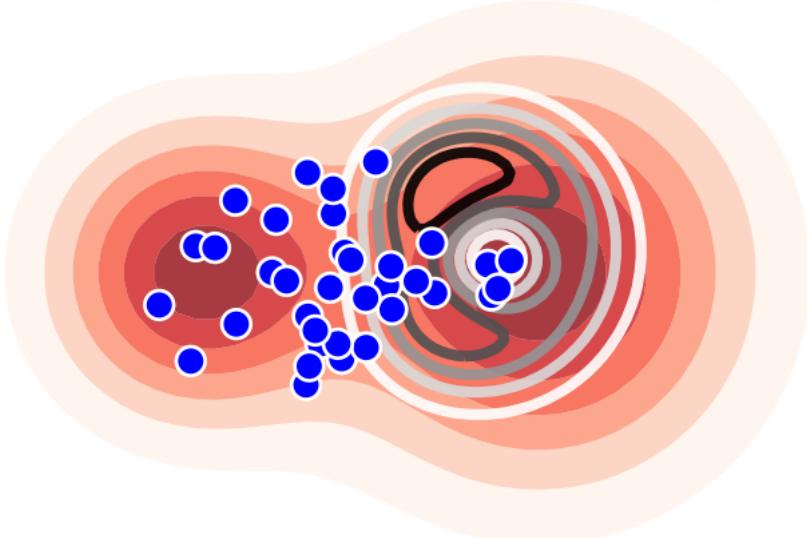
score: 0.49



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

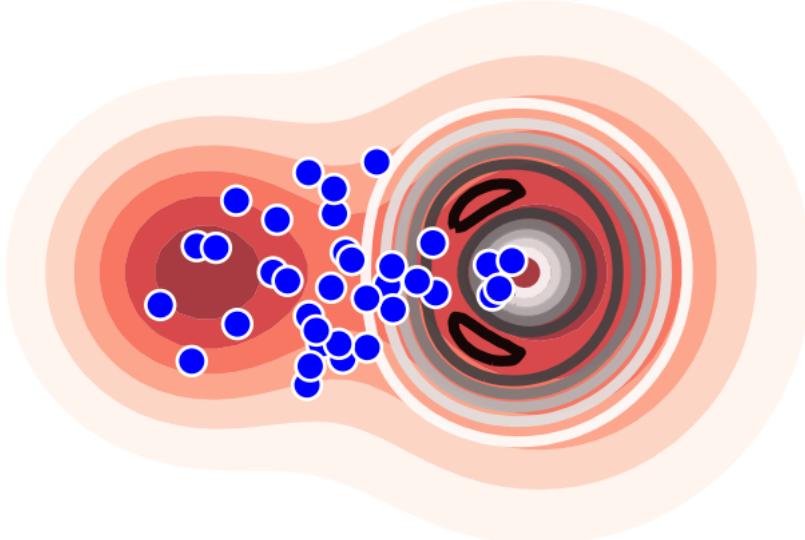
score: 0.47



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## Proposal: Model Criticism with the Stein Witness

score: 0.44



$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

## FSSD is a Discrepancy Measure

### Theorem 1.

Let  $V = \{\mathbf{v}_1, \dots, \mathbf{v}_J\} \subset \mathbb{R}^d$  be drawn i.i.d. from a distribution  $\eta$  which has a density. Let  $\mathcal{X}$  be a connected open set in  $\mathbb{R}^d$ . Assume

- 1 (*Nice RKHS*) Kernel  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is  $C_0$ -universal, and real analytic.
- 2 (*Riesz condition holds*)  $\|g\|_{\mathcal{F}}^2 < \infty$ .
- 3 (*Finite Fisher divergence*)  $\mathbb{E}_{\mathbf{x} \sim q} \|\nabla_{\mathbf{x}} \log \frac{p(\mathbf{x})}{q(\mathbf{x})}\|^2 < \infty$ .
- 4 (*vanishing boundary condition*)  $\lim_{\|\mathbf{x}\| \rightarrow \infty} p(\mathbf{x})g(\mathbf{x}) = 0$ .

Then,  $\eta$ -almost surely

$$\text{FSSD}^2 = 0 \text{ if and only if } p = q, \text{ for any } J \geq 1.$$

- Gaussian kernel  $k(\mathbf{x}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{v}\|_2^2}{2\sigma_k^2}\right)$  works.
- In practice,  $J = 1$  or  $J = 5$ .

## More on FSSD<sup>2</sup>

- When  $d > 1$ , the Stein witness  $\mathbf{g}$  has  $d$  outputs.
- Define

$$\xi(\mathbf{x}, \mathbf{v}) := \frac{1}{\mathbf{p}(\mathbf{x})} \nabla_{\mathbf{x}} [\mathbf{p}(\mathbf{x}) k(\mathbf{x}, \mathbf{v})] \in \mathbb{R}^d.$$

- $d$ -output Stein witness

$$\mathbf{g}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \xi(\mathbf{x}, \mathbf{v}) \in \mathbb{R}^d.$$

- General form:

$$\text{FSSD}^2 = \frac{1}{dJ} \sum_{j=1}^J \|\mathbf{g}(\mathbf{v}_j)\|_2^2,$$

where unbiased estimator  $\widehat{\text{FSSD}^2}$  computable in  $\mathcal{O}(d^2 J n)$ .

# Asymptotic Distributions of $\widehat{\text{FSSD}}^2$

- $\tau(\mathbf{x}) :=$  vertically stack  $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$ . Feature vector of  $\mathbf{x}$ .
- Mean feature:  $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$ .
- Equivalently,  $\text{FSSD}^2 = \frac{1}{dJ} \|\mu\|_2^2$  (mean feature).
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$  for  $r \in \{p, q\}$

## Proposition 1 (Asymptotic distributions).

Let  $Z_1, \dots, Z_{dJ} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ , and  $\{\omega_i\}_{i=1}^{dJ}$  be the eigenvalues of  $\Sigma_p$ .

- 1 Under  $H_0 : p = q$ , asymptotically  $n\widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1)\omega_i$ .
  - Easy to simulate to get p-value.
  - Simulation cost independent of  $n$ .
- 2 Under  $H_1 : p \neq q$ , we have  $\sqrt{n}(\widehat{\text{FSSD}}^2 - \text{FSSD}^2) \xrightarrow{d} \mathcal{N}(0, \sigma_{H_1}^2)$  where  $\sigma_{H_1}^2 := 4\mu^\top \Sigma_q \mu$ . Implies  $\mathbb{P}(\text{reject } H_0) \rightarrow 1$  as  $n \rightarrow \infty$ .

But, how to estimate  $\Sigma_p$ ? No sample from  $p$ !

- Theorem: Using  $\widehat{\Sigma}_q$  (computed with  $\{\mathbf{x}_i\}_{i=1}^n \sim q$ ) still leads to **23/49** asymptotically consistent test.

# Asymptotic Distributions of $\widehat{\text{FSSD}}^2$

- $\tau(\mathbf{x}) :=$  vertically stack  $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$ . Feature vector of  $\mathbf{x}$ .
- Mean feature:  $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$ .
- Equivalently,  $\text{FSSD}^2 = \frac{1}{dJ} \|\mu\|_2^2$  (mean feature).
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$  for  $r \in \{p, q\}$

**Proposition 1 (Asymptotic distributions).**

Let  $Z_1, \dots, Z_{dJ} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ , and  $\{\omega_i\}_{i=1}^{dJ}$  be the eigenvalues of  $\Sigma_p$ .

- 1 Under  $H_0 : p = q$ , asymptotically  $n\widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1)\omega_i$ .
  - Easy to simulate to get p-value.
  - Simulation cost independent of  $n$ .
- 2 Under  $H_1 : p \neq q$ , we have  $\sqrt{n}(\widehat{\text{FSSD}}^2 - \text{FSSD}^2) \xrightarrow{d} \mathcal{N}(0, \sigma_{H_1}^2)$  where  $\sigma_{H_1}^2 := 4\mu^\top \Sigma_q \mu$ . Implies  $\mathbb{P}(\text{reject } H_0) \rightarrow 1$  as  $n \rightarrow \infty$ .

**But**, how to estimate  $\Sigma_p$ ? No sample from  $p$ !

- Theorem: Using  $\widehat{\Sigma}_q$  (computed with  $\{\mathbf{x}_i\}_{i=1}^n \sim q$ ) still leads to **23/49** asymptotically consistent test.

# Asymptotic Distributions of $\widehat{\text{FSSD}}^2$

- $\tau(\mathbf{x}) :=$  vertically stack  $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$ . Feature vector of  $\mathbf{x}$ .
- Mean feature:  $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$ .
- Equivalently,  $\text{FSSD}^2 = \frac{1}{dJ} \|\mu\|_2^2$  (mean feature).
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$  for  $r \in \{p, q\}$

## Proposition 1 (Asymptotic distributions).

Let  $Z_1, \dots, Z_{dJ} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ , and  $\{\omega_i\}_{i=1}^{dJ}$  be the eigenvalues of  $\Sigma_p$ .

- 1 Under  $H_0 : p = q$ , asymptotically  $n\widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1)\omega_i$ .
  - Easy to simulate to get p-value.
  - Simulation cost independent of  $n$ .
- 2 Under  $H_1 : p \neq q$ , we have  $\sqrt{n}(\widehat{\text{FSSD}}^2 - \text{FSSD}^2) \xrightarrow{d} \mathcal{N}(0, \sigma_{H_1}^2)$  where  $\sigma_{H_1}^2 := 4\mu^\top \Sigma_q \mu$ . Implies  $\mathbb{P}(\text{reject } H_0) \rightarrow 1$  as  $n \rightarrow \infty$ .

But, how to estimate  $\Sigma_p$ ? No sample from  $p$ !

- Theorem: Using  $\hat{\Sigma}_q$  (computed with  $\{\mathbf{x}_i\}_{i=1}^n \sim q$ ) still leads to **23/49** asymptotically consistent test.

# Asymptotic Distributions of $\widehat{\text{FSSD}}^2$

- $\tau(\mathbf{x}) :=$  vertically stack  $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$ . Feature vector of  $\mathbf{x}$ .
- Mean feature:  $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$ .
- Equivalently,  $\text{FSSD}^2 = \frac{1}{dJ} \|\mu\|_2^2$  (mean feature).
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$  for  $r \in \{p, q\}$

## Proposition 1 (Asymptotic distributions).

Let  $Z_1, \dots, Z_{dJ} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ , and  $\{\omega_i\}_{i=1}^{dJ}$  be the eigenvalues of  $\Sigma_p$ .

- 1 Under  $H_0 : p = q$ , asymptotically  $n\widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1)\omega_i$ .
  - Easy to simulate to get p-value.
  - Simulation cost independent of  $n$ .
- 2 Under  $H_1 : p \neq q$ , we have  $\sqrt{n}(\widehat{\text{FSSD}}^2 - \text{FSSD}^2) \xrightarrow{d} \mathcal{N}(0, \sigma_{H_1}^2)$  where  $\sigma_{H_1}^2 := 4\mu^\top \Sigma_q \mu$ . Implies  $\mathbb{P}(\text{reject } H_0) \rightarrow 1$  as  $n \rightarrow \infty$ .

But, how to estimate  $\Sigma_p$ ? No sample from  $p$ !

- Theorem: Using  $\hat{\Sigma}_q$  (computed with  $\{\mathbf{x}_i\}_{i=1}^n \sim q$ ) still leads to **23/49** asymptotically consistent test.

# Asymptotic Distributions of $\widehat{\text{FSSD}}^2$

- $\tau(\mathbf{x}) :=$  vertically stack  $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$ . Feature vector of  $\mathbf{x}$ .
- Mean feature:  $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$ .
- Equivalently,  $\text{FSSD}^2 = \frac{1}{dJ} \|\mu\|_2^2$  (mean feature).
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$  for  $r \in \{p, q\}$

## Proposition 1 (Asymptotic distributions).

Let  $Z_1, \dots, Z_{dJ} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ , and  $\{\omega_i\}_{i=1}^{dJ}$  be the eigenvalues of  $\Sigma_p$ .

- 1 Under  $H_0 : p = q$ , asymptotically  $n\widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1)\omega_i$ .
  - Easy to simulate to get p-value.
  - Simulation cost independent of  $n$ .
- 2 Under  $H_1 : p \neq q$ , we have  $\sqrt{n}(\widehat{\text{FSSD}}^2 - \text{FSSD}^2) \xrightarrow{d} \mathcal{N}(0, \sigma_{H_1}^2)$  where  $\sigma_{H_1}^2 := 4\mu^\top \Sigma_q \mu$ . Implies  $\mathbb{P}(\text{reject } H_0) \rightarrow 1$  as  $n \rightarrow \infty$ .

But, how to estimate  $\Sigma_p$ ? No sample from  $p$ !

- Theorem: Using  $\hat{\Sigma}_q$  (computed with  $\{\mathbf{x}_i\}_{i=1}^n \sim q$ ) still leads to  $_{23/49}$  asymptotically consistent test.

## Parameter Tuning

- Jointly optimise locations  $V = \{\mathbf{v}_1, \dots, \mathbf{v}_J\}$  for more test power

**Proposition 2** (Approx. power for large  $n$ ).

Under  $H_1$ , for large  $n$  and fixed threshold  $r$ , the test power  
 $\mathbb{P}(\text{reject } H_0 \mid H_1 \text{ true})$

$$\mathbb{P}_{H_1}(n\widehat{\text{FSSD}}^2 > r) \approx 1 - \Phi\left(\frac{r}{\sqrt{n}\sigma_{H_1}} - \sqrt{n}\frac{\text{FSSD}^2}{\sigma_{H_1}}\right),$$

where  $\Phi = \text{CDF of } \mathcal{N}(0, 1)$ .

- For large  $n$ , second term dominates. So

$$\arg \max_{V, \sigma_k^2} (\text{power}) \approx \arg \max_{V, \sigma_k^2} \frac{\widehat{\text{FSSD}}^2}{\widehat{\sigma}_{H_1}^2}.$$

- Split  $\{\mathbf{x}_i\}_{i=1}^n$  into independent training/test sets. Optimize  $V$  on tr. Test on te.

## Parameter Tuning

- Jointly optimise locations  $V = \{\mathbf{v}_1, \dots, \mathbf{v}_J\}$  for more test power

### Proposition 2 (Approx. power for large $n$ ).

Under  $H_1$ , for large  $n$  and fixed threshold  $r$ , the test power  
 $\mathbb{P}(\text{reject } H_0 \mid H_1 \text{ true})$

$$\mathbb{P}_{H_1}(n\widehat{\text{FSSD}}^2 > r) \approx 1 - \Phi\left(\frac{r}{\sqrt{n}\sigma_{H_1}} - \sqrt{n}\frac{\text{FSSD}^2}{\sigma_{H_1}}\right),$$

where  $\Phi = \text{CDF of } \mathcal{N}(0, 1)$ .

- For large  $n$ , second term dominates. So

$$\arg \max_{V, \sigma_k^2} (\text{power}) \approx \arg \max_{V, \sigma_k^2} \frac{\widehat{\text{FSSD}}^2}{\widehat{\sigma}_{H_1}^2}.$$

- Split  $\{\mathbf{x}_i\}_{i=1}^n$  into independent training/test sets. Optimize  $V$  on tr. Test on te.

## Parameter Tuning

- Jointly optimise locations  $V = \{\mathbf{v}_1, \dots, \mathbf{v}_J\}$  for more test power

### Proposition 2 (Approx. power for large $n$ ).

Under  $H_1$ , for large  $n$  and fixed threshold  $r$ , the test power  
 $\mathbb{P}(\text{reject } H_0 \mid H_1 \text{ true})$

$$\mathbb{P}_{H_1}(n\widehat{\text{FSSD}}^2 > r) \approx 1 - \Phi\left(\frac{r}{\sqrt{n}\sigma_{H_1}} - \sqrt{n}\frac{\text{FSSD}^2}{\sigma_{H_1}}\right),$$

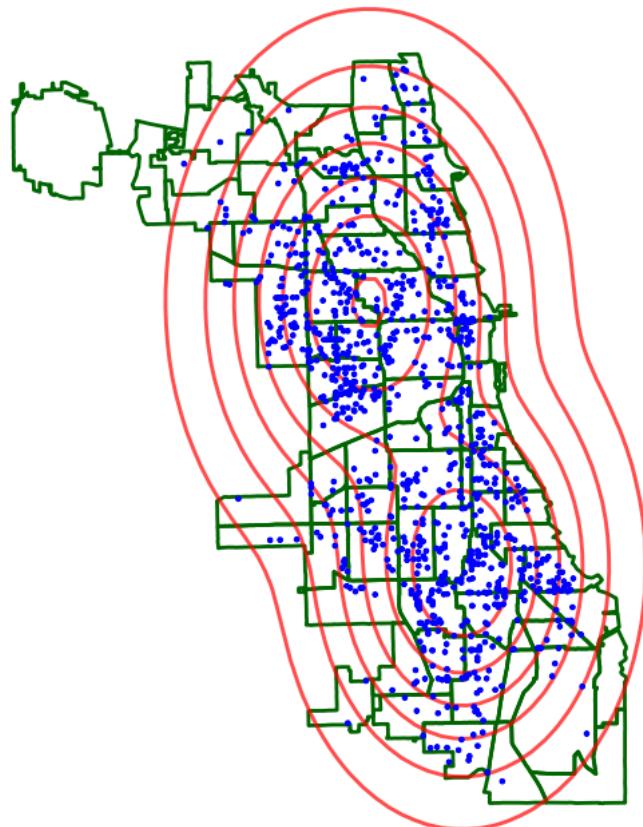
where  $\Phi = \text{CDF of } \mathcal{N}(0, 1)$ .

- For large  $n$ , second term dominates. So

$$\arg \max_{V, \sigma_k^2} (\text{power}) \approx \arg \max_{V, \sigma_k^2} \frac{\widehat{\text{FSSD}}^2}{\widehat{\sigma}_{H_1}^2}.$$

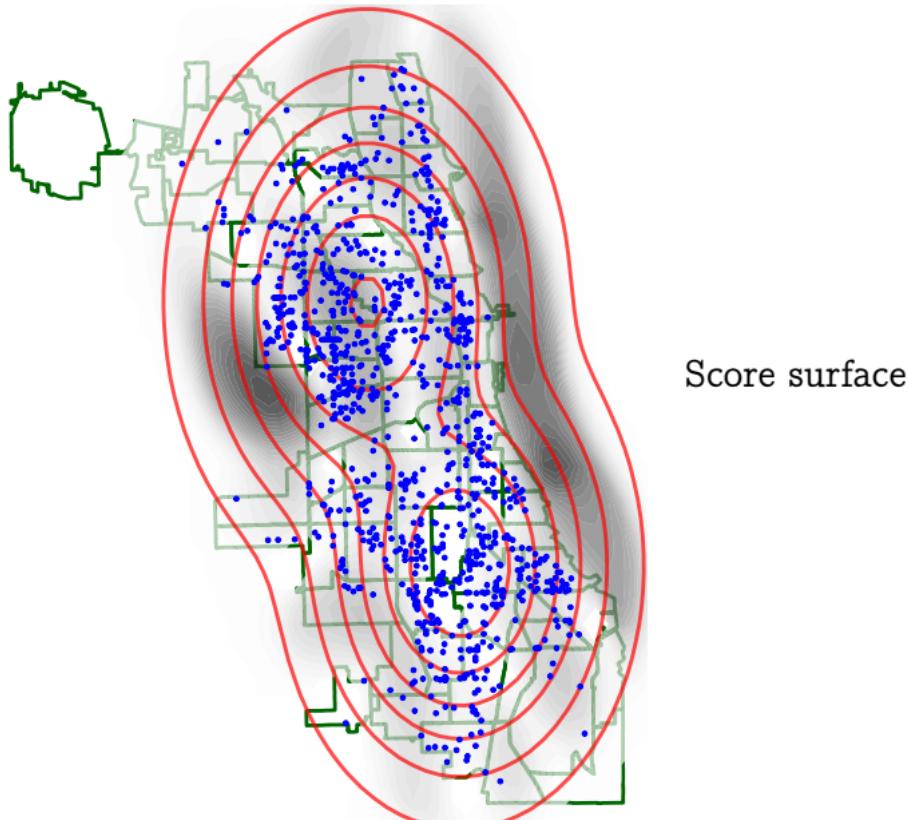
- Split  $\{\mathbf{x}_i\}_{i=1}^n$  into independent training/test sets. Optimize  $V$  on tr.  
Test on te.

## Interpretable Features: Chicago Crime



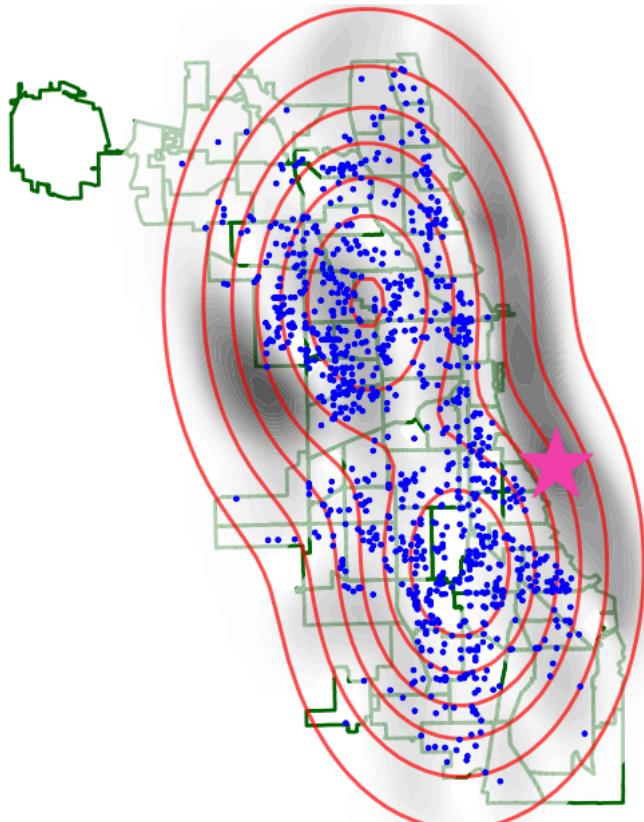
Model  $p$  = 2-component Gaussian mixture.

## Interpretable Features: Chicago Crime



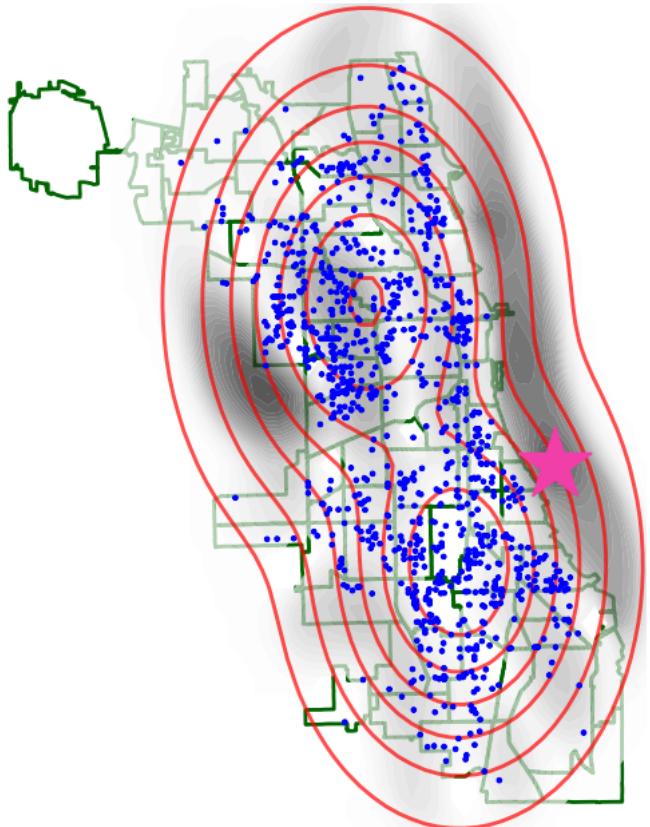
Score surface

## Interpretable Features: Chicago Crime



★ = optimized  $v$ .

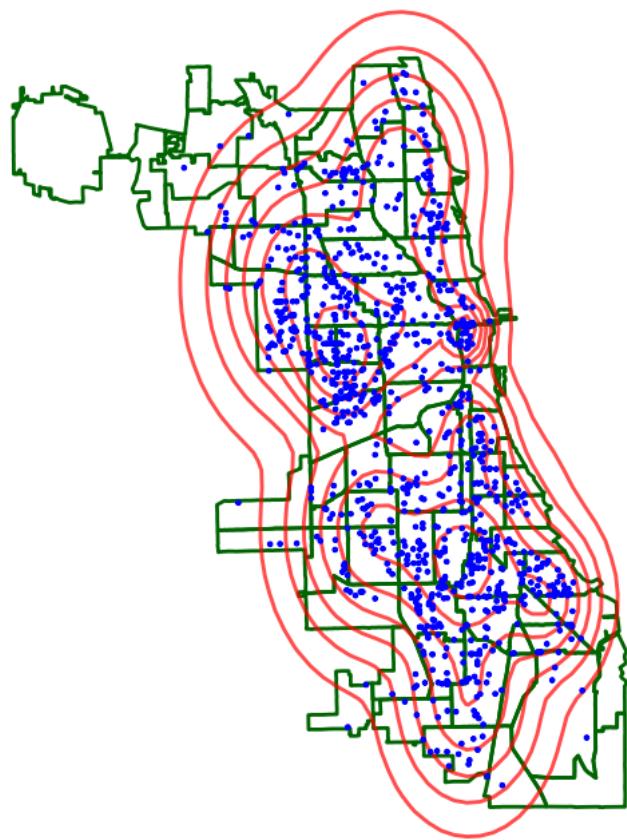
## Interpretable Features: Chicago Crime



$\star$  = optimized  $v$ .  
No robbery in Lake Michigan.

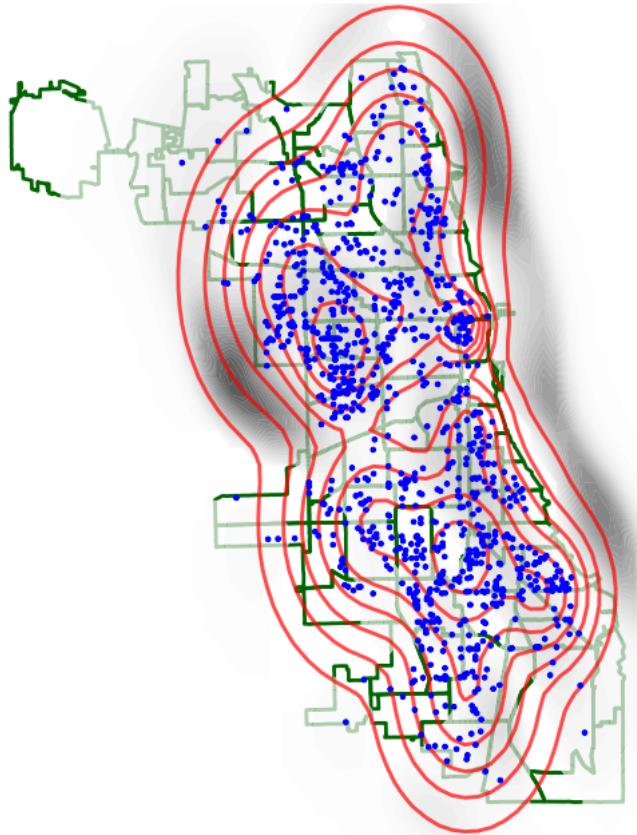


## Interpretable Features: Chicago Crime



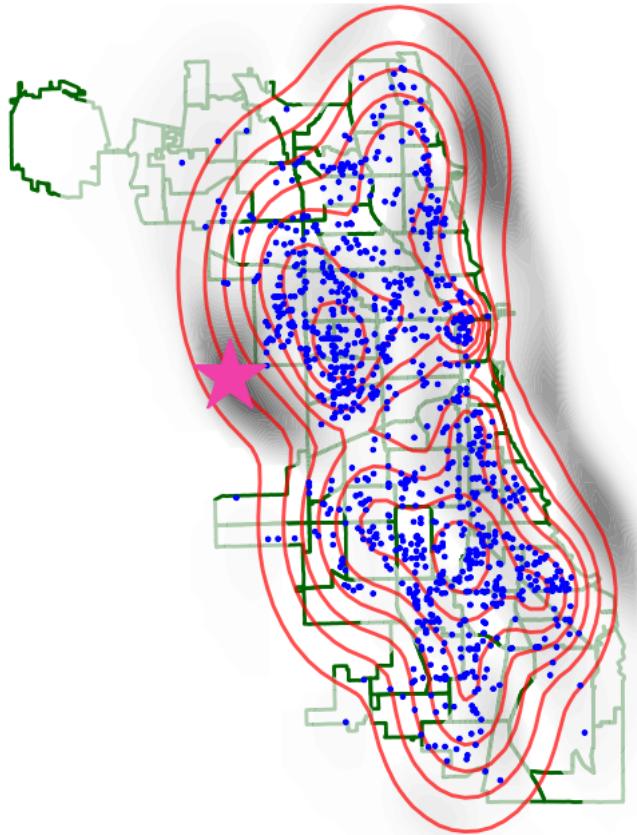
Model  $p = 10$ -component Gaussian mixture.

## Interpretable Features: Chicago Crime



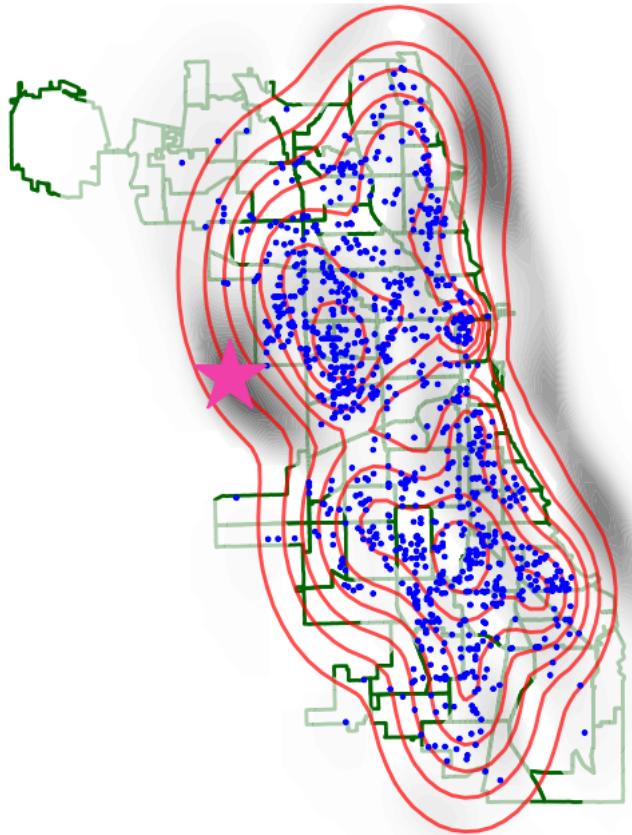
Capture the right tail better.

## Interpretable Features: Chicago Crime



Still, does not capture the left tail.

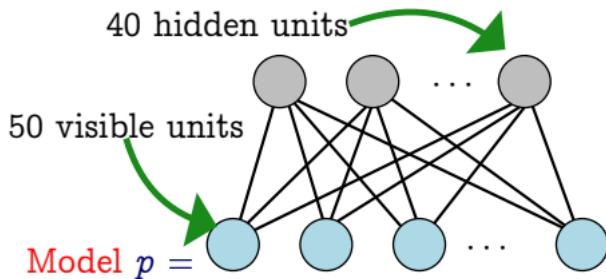
## Interpretable Features: Chicago Crime



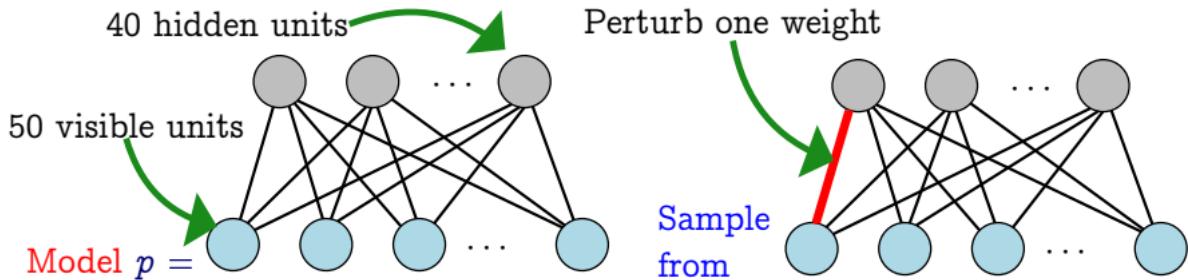
Still, does not capture the left tail.

Learned test locations are interpretable.

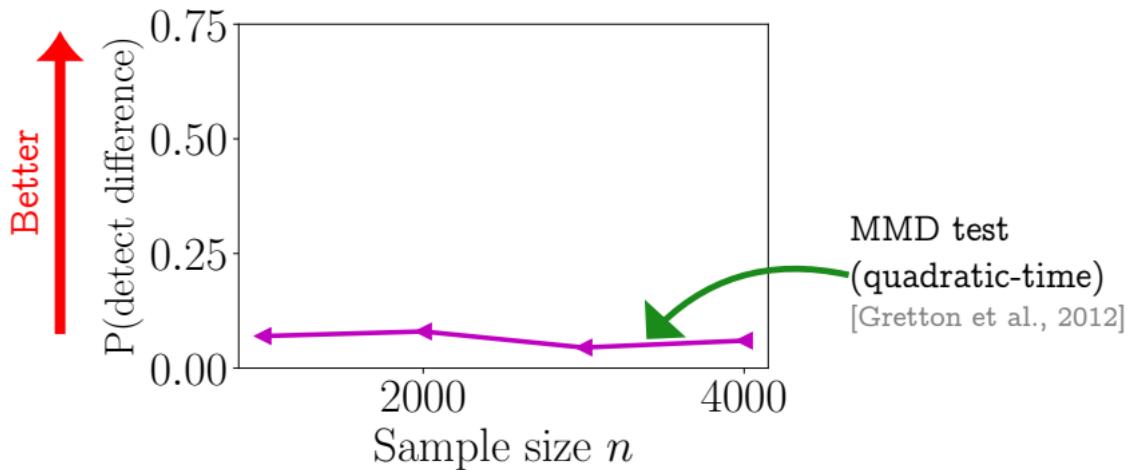
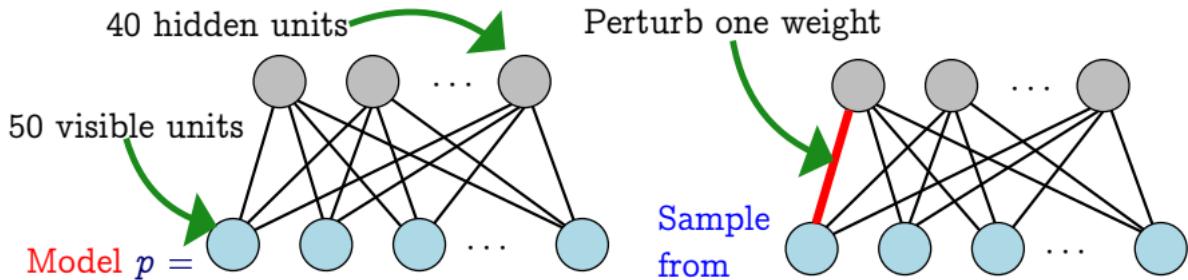
## Experiment: Restricted Boltzmann Machine (RBM)



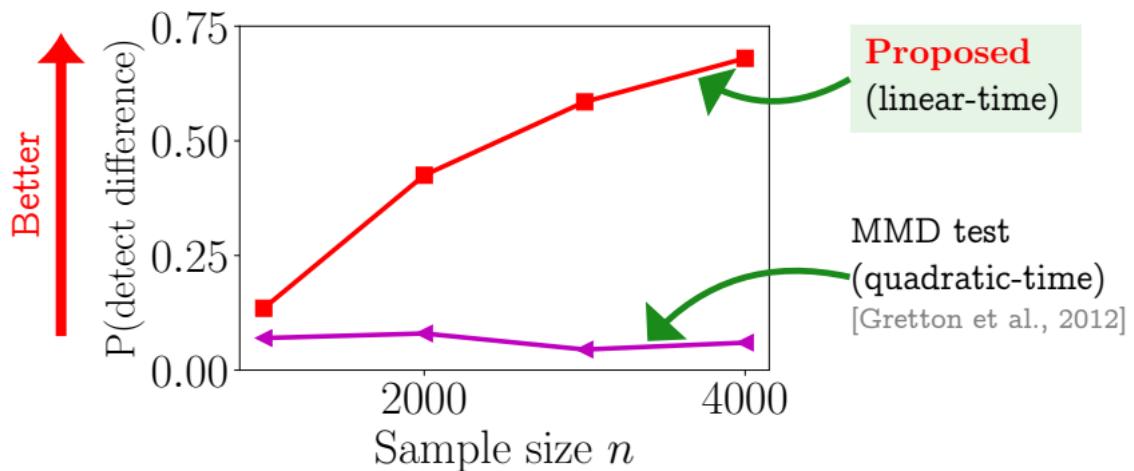
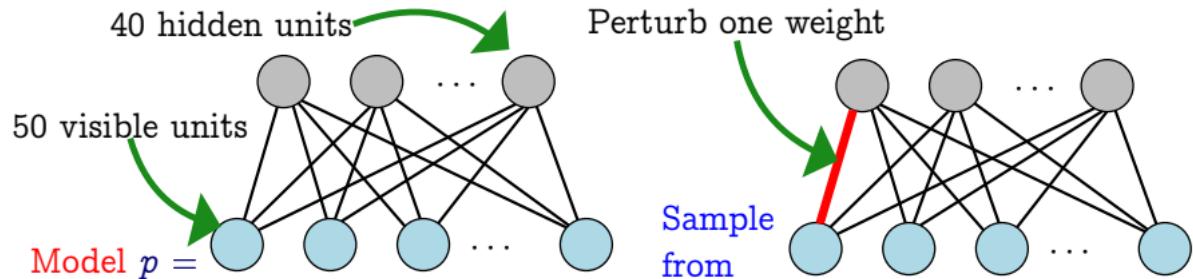
## Experiment: Restricted Boltzmann Machine (RBM)



## Experiment: Restricted Boltzmann Machine (RBM)



## Experiment: Restricted Boltzmann Machine (RBM)



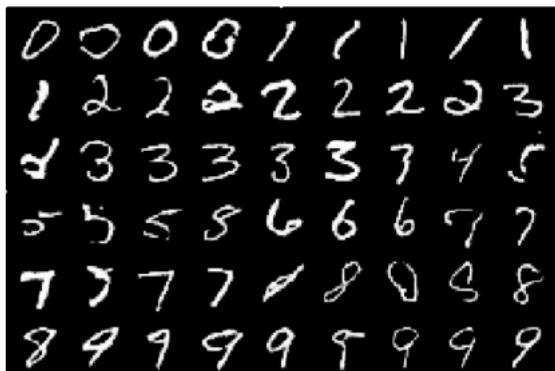
## Model Criticism

*"All models are wrong."*

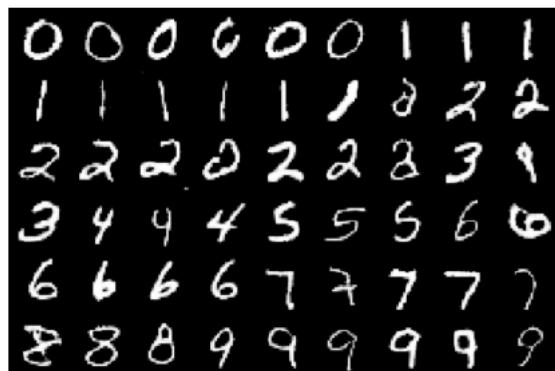
G. Box (1976)

## 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)

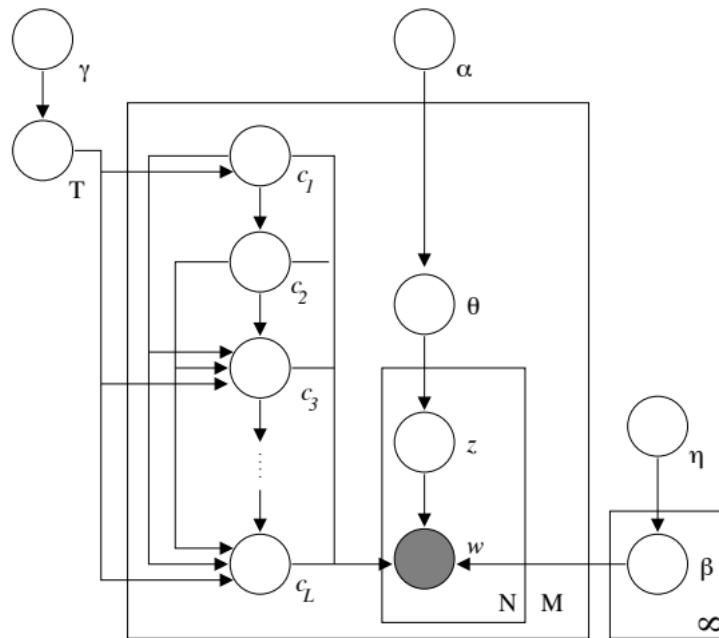


Samples from LSGAN, Mao et al. (2017)

**Which model is better?**

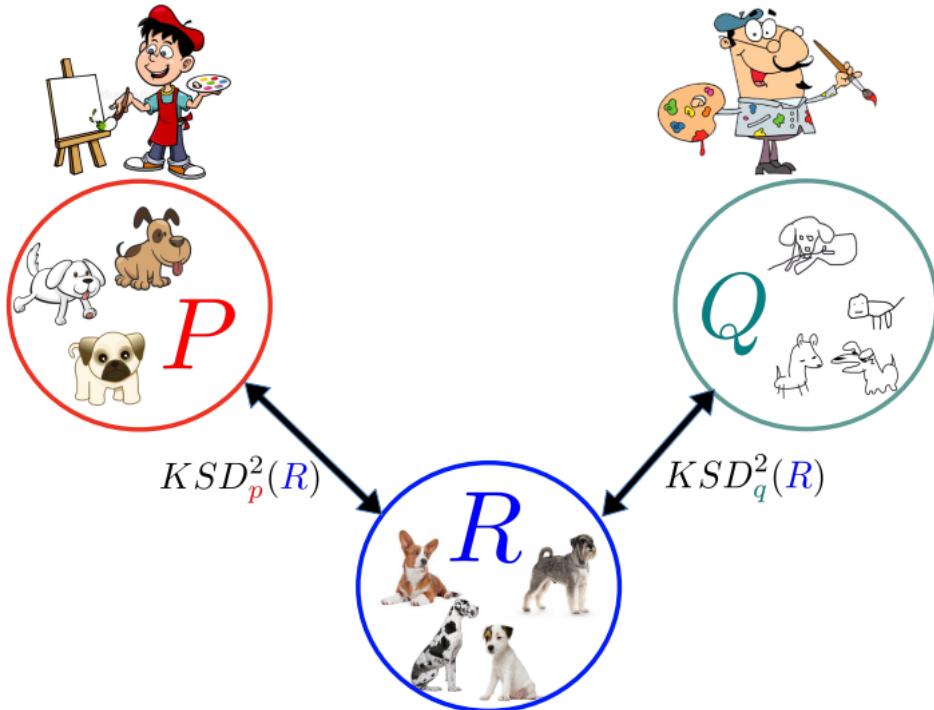
## Most interesting models have latent structure

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



## Relative goodness-of-fit testing

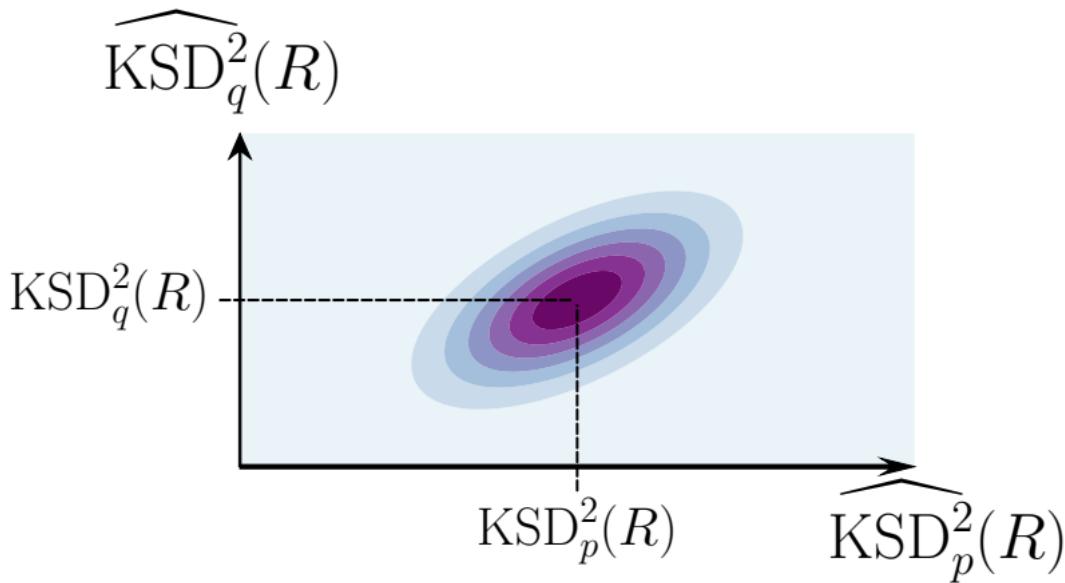
- Two generative models  $P$  and  $Q$ , data  $\{x_i\}_{i=1}^n \sim R$ .
- Neither model gives a perfect fit ( $P \neq R$  and  $Q \neq R$ ).



## Joint asymptotic normality

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

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



## 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:

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

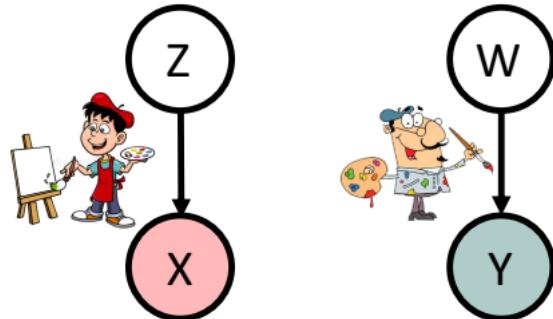
$\implies$  a statistical test with null hypothesis  $\text{KSD}_p^2(R) - \text{KSD}_q^2(R) \leq 0$  is straightforward.

## Latent variable models

Can we compare latent variable models with KSD?

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

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



Recall multi-dimensional Stein operator:

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

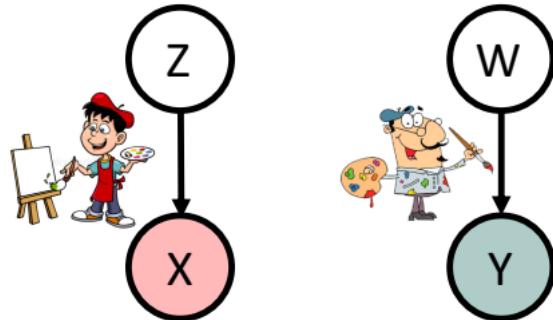
Expression  $(\alpha)$  requires marginal  $p(x)$ , often intractable...

## 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:

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

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

## Monte Carlo approximation

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

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

Estimate KSDs with approximate densities:

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

## Monte Carlo approximation

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 KSDs with approximate densities:

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

Recall

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

→ if  $m$  is large, can we simply substitute  $p_m$  and  $q_m$  ?

## Simple proof of concept

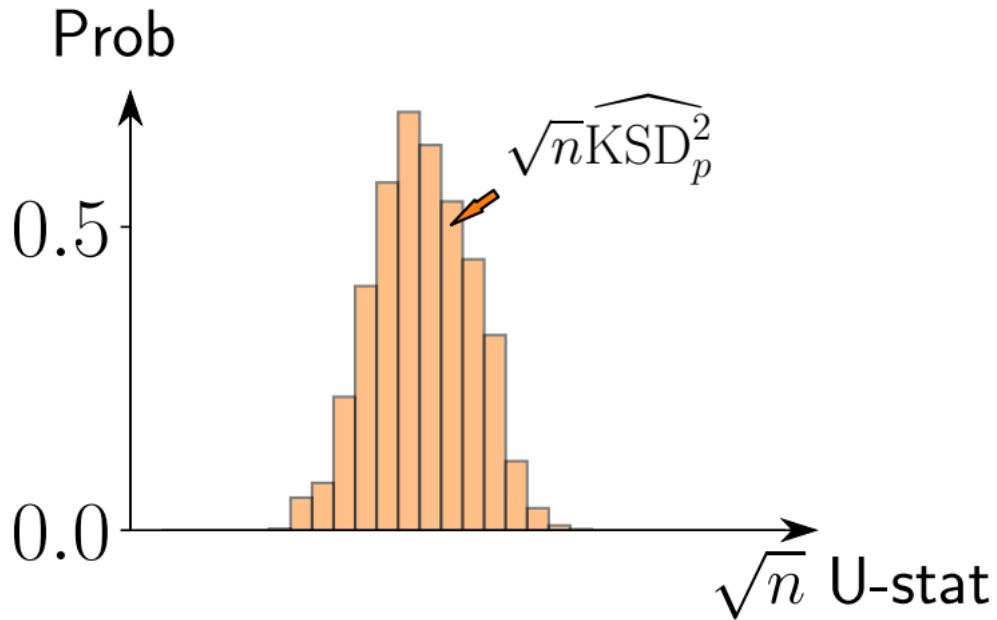
Check  $\widehat{\text{KSD}}_{\textcolor{red}{p}}^2(R) \approx \widehat{\text{KSD}}_{\textcolor{red}{p}_m}^2(R)$  with a toy model:

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

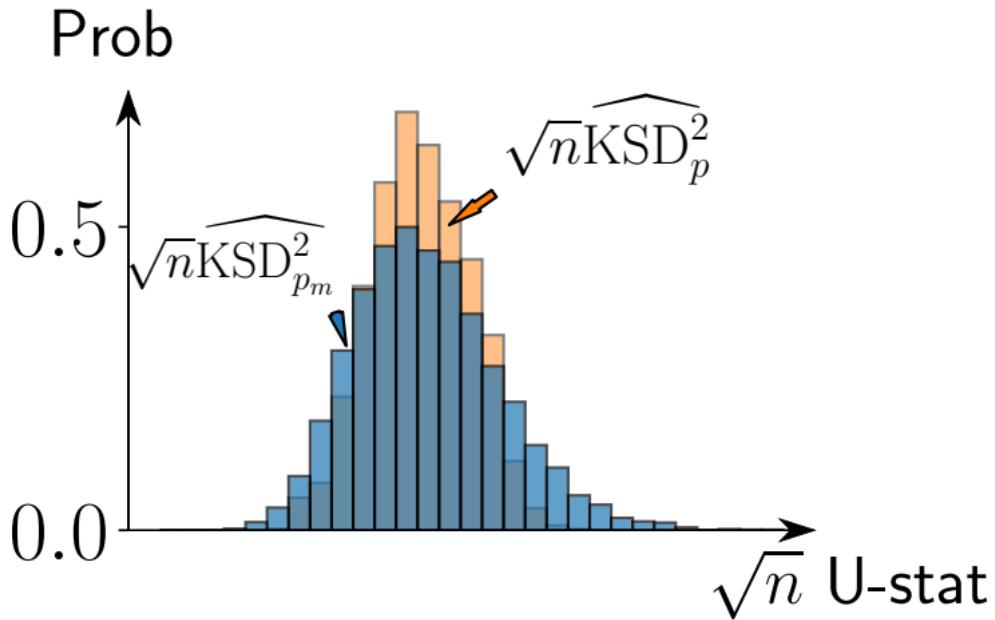
$$p(x|z) = \binom{N}{x} z^x (1-z)^{n-x}, \quad p(z) = \text{Beta}(a, b)$$

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

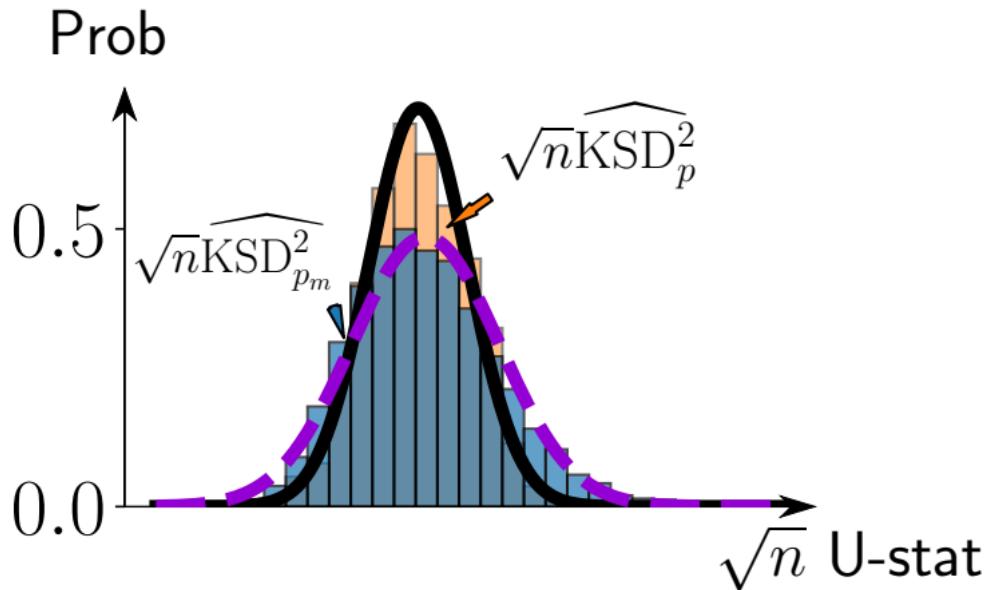
## Effect of sampling the latents (Beta-binomial)



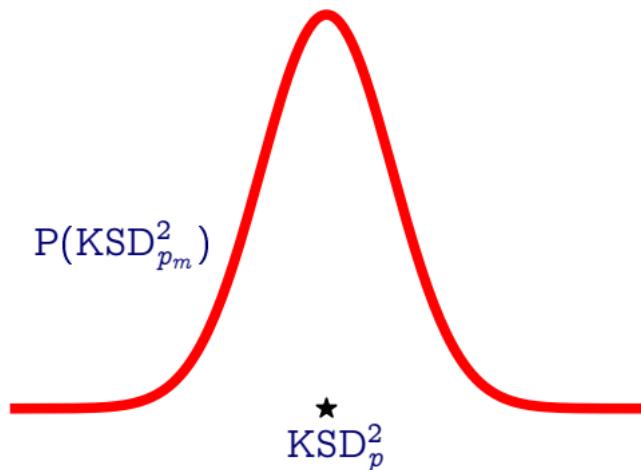
## Effect of sampling the latents (Beta-binomial)



## Effect of sampling the latents (Beta-binomial)

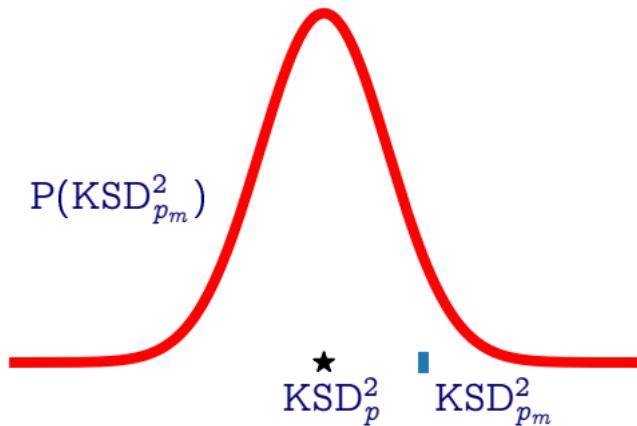


## Why this happens



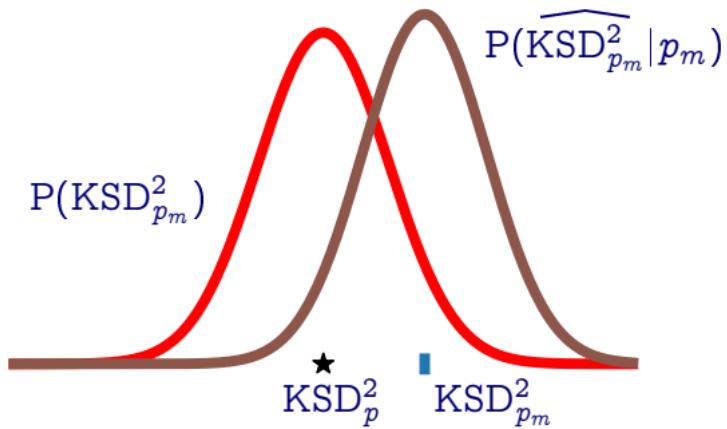
$KSD_{\textcolor{red}{p}_m}^2(\textcolor{blue}{R})$  is normally distributed around  $KSD_{\textcolor{red}{p}}^2(\textcolor{blue}{R})$   
(approximation error)

## Why this happens



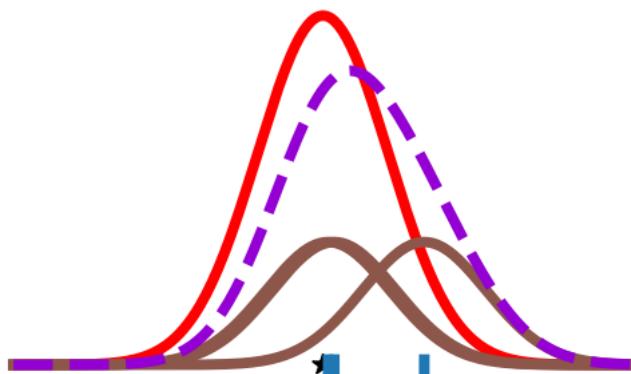
Approximation  $p_m$  gives a random draw  $KSD_{p_m}^2(R)$

## Why this happens



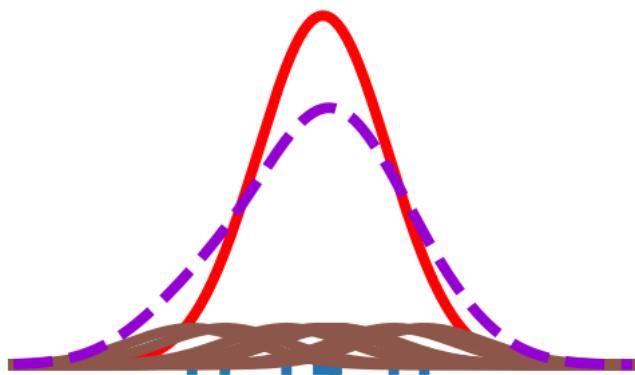
$\widehat{P(KSD^2_{p_m})}(R)$  is normally distributed around  $KSD^2_{p_m}(R)$

## Why this happens



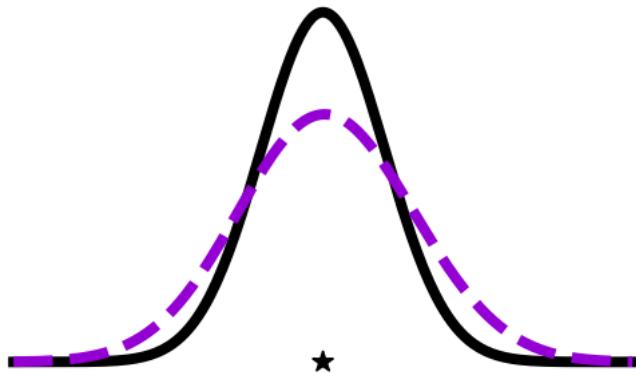
Distribution of  $\widehat{\text{KSD}}_{p_m}^2(R)$  is  
averaged over random draws of  $\text{KSD}_{p_m}^2(R)$

## Why this happens



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

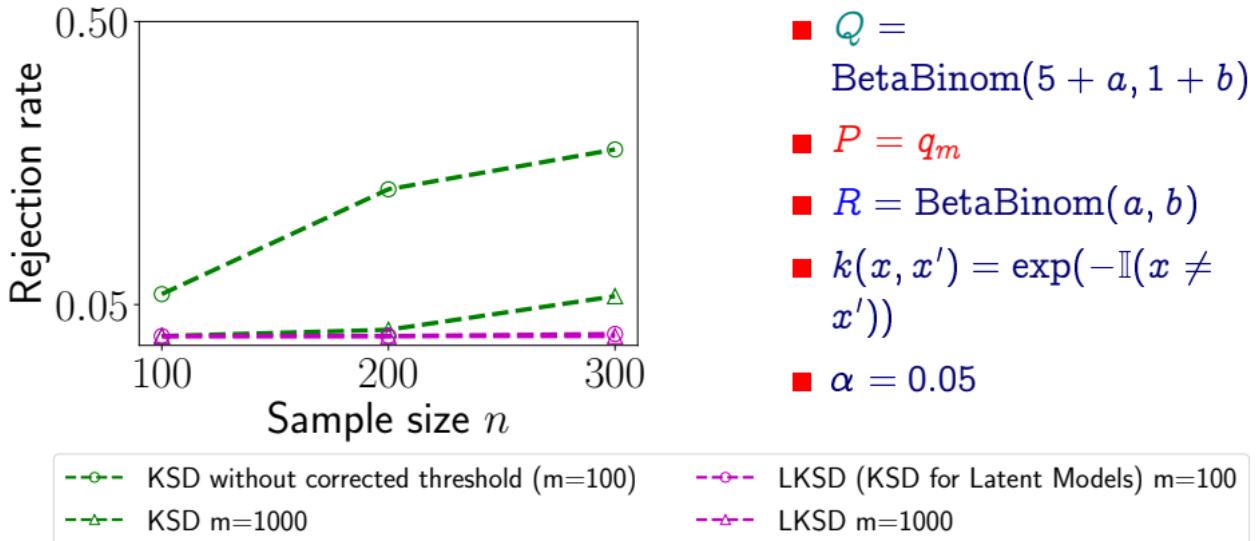
## Why this happens



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

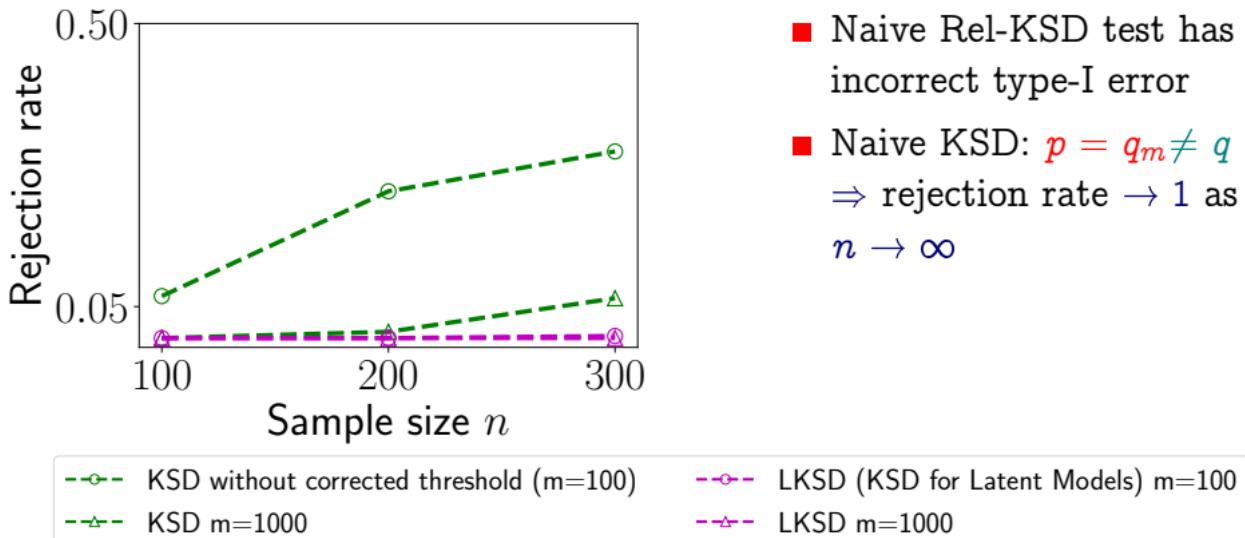
## Correction for this effect

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



## Correction for this effect

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



- Naive Rel-KSD test has incorrect type-I error
- Naive KSD:  $p = q_m \neq q$   
⇒ rejection rate → 1 as  $n \rightarrow \infty$

## 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

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

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

where

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

Fine print:

- $h_{\textcolor{red}{p}}(x, x') - h_{\textcolor{teal}{q}}(x, x')$  has a finite third moment
- Additional technical conditions

## Relative test, further detail

**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  $\nu_3(H_m) = \mathbf{E}_{x,x'} |H_m(x, x') - \mathbf{E}_{x,x'} H_m(x, x')|^3$
- $Y_m := \sqrt{m} \left( \mathbf{E}_n[U_{n,m}|H_m] - \mathbf{E}_n[U_n] \right) \xrightarrow{d} Y$

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

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

## Experiment: sensitivity to model difference

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

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

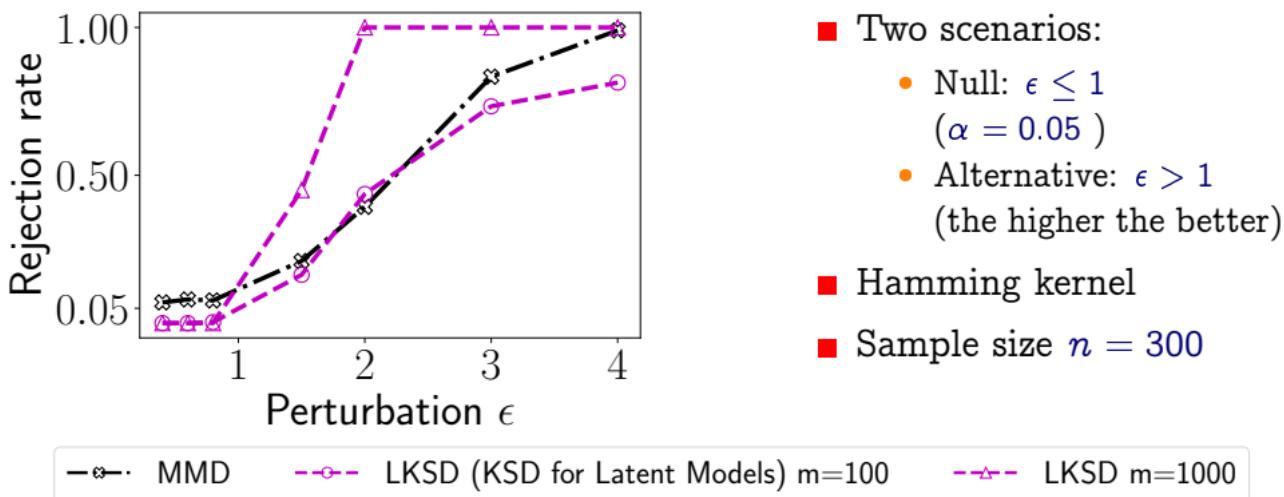
- Models:  $P = SBN(W + \epsilon[1, 0, \dots, 0])$ ,  $Q = SBN(W + [1, 0, \dots, 0])$
- Only the first column of weight  $W$  is perturbed by  $\epsilon$

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

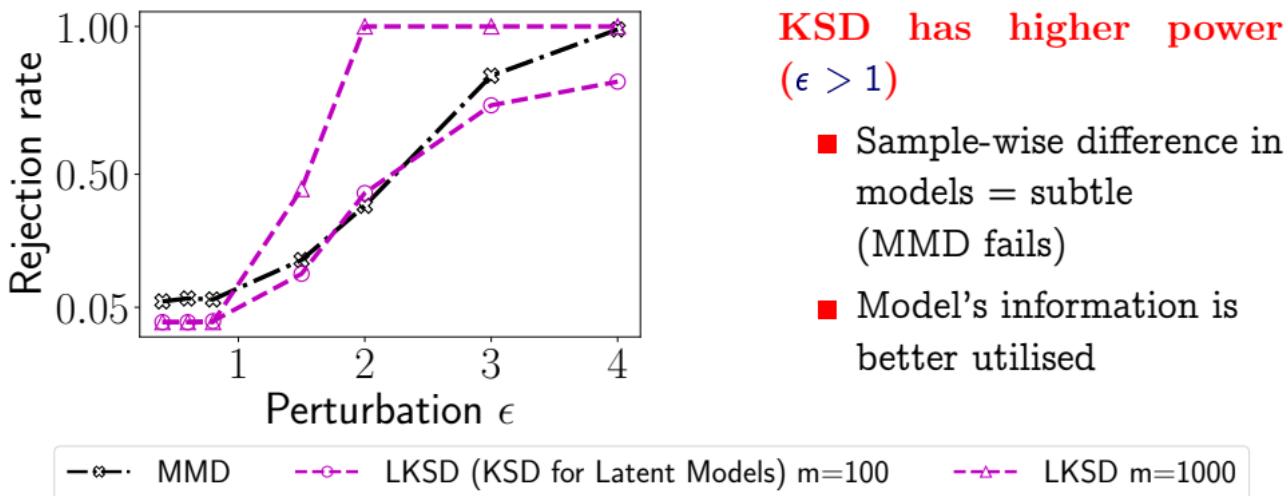


## Experiment: sensitivity to model difference

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



## Papers referenced

A Linear-Time Kernel Goodness-of-Fit Test.

Wittawat Jitkrittum, Wenkai Xu, Zoltán Szabó, Kenji Fukumizu,  
Arthur Gretton

<https://arxiv.org/abs/1705.07673>

■ Python code: <https://github.com/wittawatj/kernel-gof>

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?



# Efficiency comparison, linear-time tests

## Bahadur Slope and Bahadur Efficiency

- Bahadur slope  $\approx$  rate of p-value  $\rightarrow 0$  under  $H_1$  as  $n \rightarrow \infty$ .
- Measure a test's sensitivity to the departure from  $H_0$ .

$$H_0 : \theta = 0,$$

$$H_1 : \theta \neq 0.$$

- Typically  $\text{pval}_n \approx \exp\left(-\frac{1}{2}c(\theta)n\right)$  where  $c(\theta) > 0$  under  $H_1$ , and  $c(0) = 0$ . [?].
- $c(\theta)$  higher  $\implies$  more sensitive. Good.

### Bahadur slope

$$c(\theta) := -2 \operatorname{plim}_{n \rightarrow \infty} \frac{\log(1 - F(T_n))}{n},$$

where  $F(t) = \text{CDF}$  of  $T_n$  under  $H_0$ .

- Bahadur efficiency = ratio of slopes of two tests.

## Bahadur Slope and Bahadur Efficiency

- Bahadur slope  $\approx$  rate of p-value  $\rightarrow 0$  under  $H_1$  as  $n \rightarrow \infty$ .
- Measure a test's sensitivity to the departure from  $H_0$ .

$$H_0 : \theta = 0,$$

$$H_1 : \theta \neq 0.$$

- Typically  $pval_n \approx \exp\left(-\frac{1}{2}c(\theta)n\right)$  where  $c(\theta) > 0$  under  $H_1$ , and  $c(0) = 0$ . [?].
- $c(\theta)$  higher  $\implies$  more sensitive. Good.

Bahadur slope

$$c(\theta) := -2 \operatorname{plim}_{n \rightarrow \infty} \frac{\log(1 - F(T_n))}{n},$$

where  $F(t) = \text{CDF of } T_n \text{ under } H_0$ .

- Bahadur efficiency = ratio of slopes of two tests.

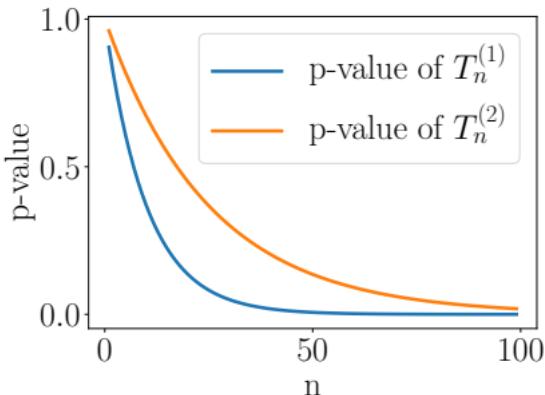
## Bahadur Slope and Bahadur Efficiency

- Bahadur slope  $\approx$  rate of p-value  $\rightarrow 0$  under  $H_1$  as  $n \rightarrow \infty$ .
- Measure a test's sensitivity to the departure from  $H_0$ .

$$H_0 : \theta = 0,$$

$$H_1 : \theta \neq 0.$$

- Typically  $p_{\text{val},n} \approx \exp\left(-\frac{1}{2}c(\theta)n\right)$  where  $c(\theta) > 0$  under  $H_1$ , and  $c(0) = 0$ . [?].
- $c(\theta)$  higher  $\implies$  more sensitive. Good.



Bahadur slope

$$c(\theta) := -2 \operatorname{plim}_{n \rightarrow \infty} \frac{\log(1 - F(T_n))}{n},$$

where  $F(t) = \text{CDF of } T_n \text{ under } H_0$ .

- Bahadur efficiency = ratio of slopes of two tests.

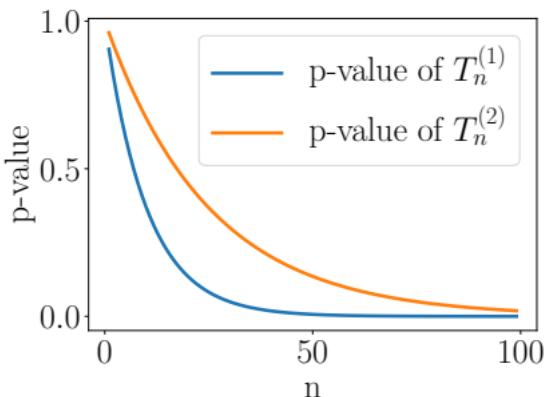
## Bahadur Slope and Bahadur Efficiency

- Bahadur slope  $\approx$  rate of p-value  $\rightarrow 0$  under  $H_1$  as  $n \rightarrow \infty$ .
- Measure a test's sensitivity to the departure from  $H_0$ .

$$H_0 : \theta = 0,$$

$$H_1 : \theta \neq 0.$$

- Typically  $p_{\text{val},n} \approx \exp\left(-\frac{1}{2}c(\theta)n\right)$  where  $c(\theta) > 0$  under  $H_1$ , and  $c(0) = 0$ . [?].
- $c(\theta)$  higher  $\implies$  more sensitive. Good.



Bahadur slope

$$c(\theta) := -2 \operatorname{plim}_{n \rightarrow \infty} \frac{\log(1 - F(T_n))}{n},$$

where  $F(t) = \text{CDF of } T_n \text{ under } H_0$ .

- Bahadur efficiency = ratio of slopes of two tests.

## Bahadur Slopes of FSSD and LKS

### Theorem 2.

The Bahadur slope of  $n\widehat{\text{FSSD}}^2$  is

$$c^{(\text{FSSD})} := \text{FSSD}^2 / \omega_1,$$

where  $\omega_1$  is the maximum eigenvalue of  $\Sigma_p := \text{cov}_{\mathbf{x} \sim p}[\tau(\mathbf{x})]$ .

### Theorem 3.

The Bahadur slope of the linear-time kernel Stein (LKS) statistic  $\sqrt{n}\widehat{S}_l^2$  is

$$c^{(\text{LKS})} = \frac{1}{2} \frac{\mathbb{E}_q h_p(\mathbf{x}, \mathbf{x}')^2}{\mathbb{E}_p [h_p^2(\mathbf{x}, \mathbf{x}')]},$$

where  $h_p$  is the U-statistic kernel of the KSD statistic.

- Let's consider a specific case ...

## Bahadur Slopes of FSSD and LKS

### Theorem 2.

The Bahadur slope of  $n\widehat{\text{FSSD}}^2$  is

$$c^{(\text{FSSD})} := \text{FSSD}^2 / \omega_1,$$

where  $\omega_1$  is the maximum eigenvalue of  $\Sigma_p := \text{cov}_{\mathbf{x} \sim p}[\tau(\mathbf{x})]$ .

### Theorem 3.

The Bahadur slope of the linear-time kernel Stein (LKS) statistic  $\sqrt{n}\widehat{S}_l^2$  is

$$c^{(\text{LKS})} = \frac{1}{2} \frac{\mathbb{E}_q h_p(\mathbf{x}, \mathbf{x}')^2}{\mathbb{E}_p [h_p^2(\mathbf{x}, \mathbf{x}')]},$$

where  $h_p$  is the U-statistic kernel of the KSD statistic.

- Let's consider a specific case ...

## Gaussian Mean Shift Problem

Consider  $p = \mathcal{N}(0, 1)$  and  $q = \mathcal{N}(\mu_q, 1)$ .

- Assume  $J = 1$  feature for  $n\widehat{\text{FSSD}}^2$ . Gaussian kernel (bandwidth =  $\sigma_k^2$ )

$$c^{(\text{FSSD})}(\mu_q, v, \sigma_k^2) = \frac{\sigma_k^2 (\sigma_k^2 + 2)^3 \mu_q^2 e^{\frac{v^2}{\sigma_k^2+2} - \frac{(v-\mu_q)^2}{\sigma_k^2+1}}}{\sqrt{\frac{2}{\sigma_k^2} + 1} (\sigma_k^2 + 1) (\sigma_k^6 + 4\sigma_k^4 + (v^2 + 5)\sigma_k^2 + 2)}.$$

- For LKS, Gaussian kernel (bandwidth =  $\kappa^2$ ).

$$c^{(\text{LKS})}(\mu_q, \kappa^2) = \frac{(\kappa^2)^{5/2} (\kappa^2 + 4)^{5/2} \mu_q^4}{2(\kappa^2 + 2)(\kappa^8 + 8\kappa^6 + 21\kappa^4 + 20\kappa^2 + 12)}.$$

## Gaussian Mean Shift Problem

Consider  $p = \mathcal{N}(0, 1)$  and  $q = \mathcal{N}(\mu_q, 1)$ .

- Assume  $J = 1$  feature for  $n\widehat{\text{FSSD}}^2$ . Gaussian kernel (bandwidth =  $\sigma_k^2$ )

$$c^{(\text{FSSD})}(\mu_q, v, \sigma_k^2) = \frac{\sigma_k^2 (\sigma_k^2 + 2)^3 \mu_q^2 e^{\frac{v^2}{\sigma_k^2+2} - \frac{(v-\mu_q)^2}{\sigma_k^2+1}}}{\sqrt{\frac{2}{\sigma_k^2} + 1} (\sigma_k^2 + 1) (\sigma_k^6 + 4\sigma_k^4 + (v^2 + 5)\sigma_k^2 + 2)}.$$

- For LKS, Gaussian kernel (bandwidth =  $\kappa^2$ ).

$$c^{(\text{LKS})}(\mu_q, \kappa^2) = \frac{(\kappa^2)^{5/2} (\kappa^2 + 4)^{5/2} \mu_q^4}{2(\kappa^2 + 2)(\kappa^8 + 8\kappa^6 + 21\kappa^4 + 20\kappa^2 + 12)}.$$

## Gaussian Mean Shift Problem

Theorem 4 (FSSD is at least two times more efficient).

- Fix  $\sigma_k^2 = 1$  for  $n\widehat{\text{FSSD}}^2$ .

Then,  $\forall \mu_q \neq 0$ ,  $\exists v \in \mathbb{R}$ ,  $\forall \kappa^2 > 0$ , we have Bahadur efficiency

$$\frac{c^{(\text{FSSD})}(\mu_q, v, \sigma_k^2)}{c^{(\text{LKS})}(\mu_q, \kappa^2)} > 2.$$