

A Linear-Time Kernel Goodness-of-Fit Test

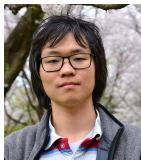
Wittawat Jitkrittum¹

Wenkai Xu¹

Zoltán Szabó²

Kenji Fukumizu³

Arthur Gretton¹



wittawat@gatsby.ucl.ac.uk

¹Gatsby Unit, University College London

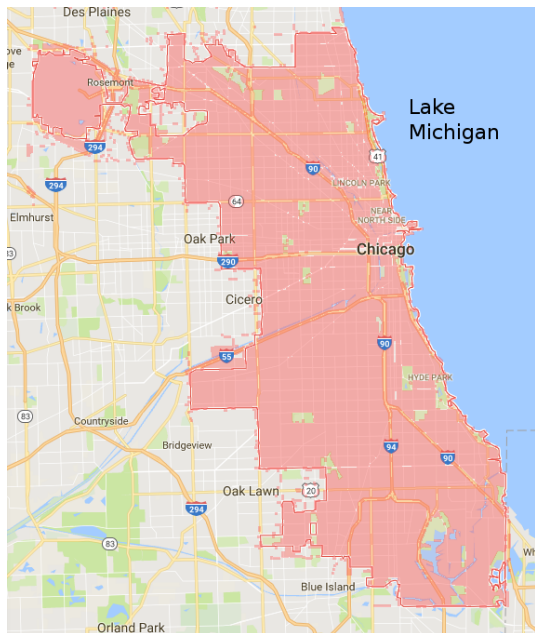
²CMAP, École Polytechnique

³The Institute of Statistical Mathematics, Tokyo

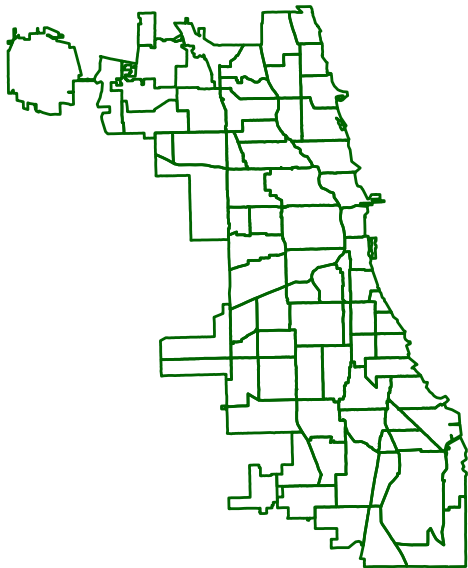
NIPS 2017, Long Beach

5 December 2017

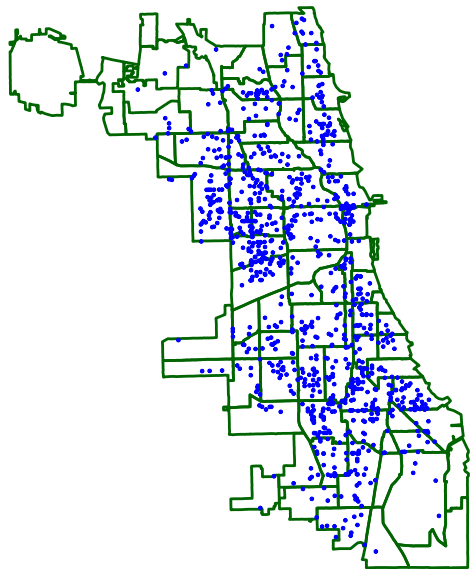
Model Criticism



Model Criticism

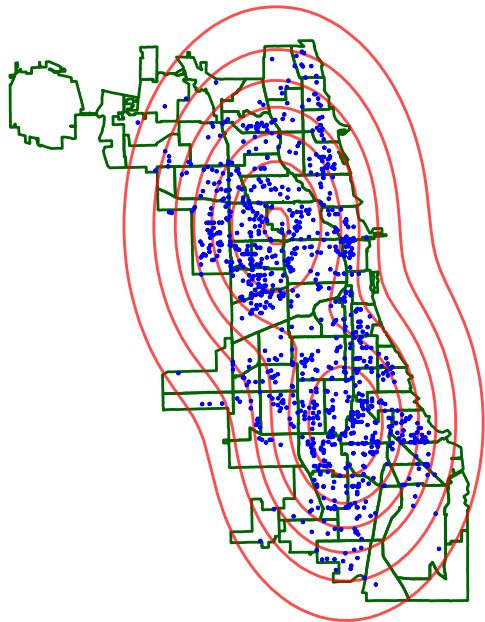


Model Criticism



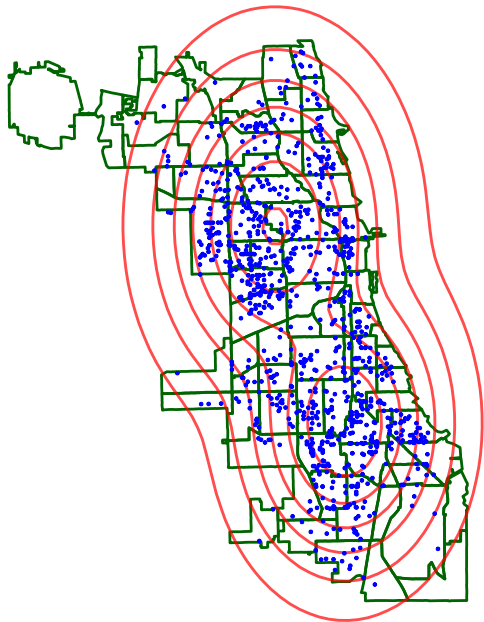
Data = robbery events in
Chicago in 2016.

Model Criticism



Is this a good **model**?

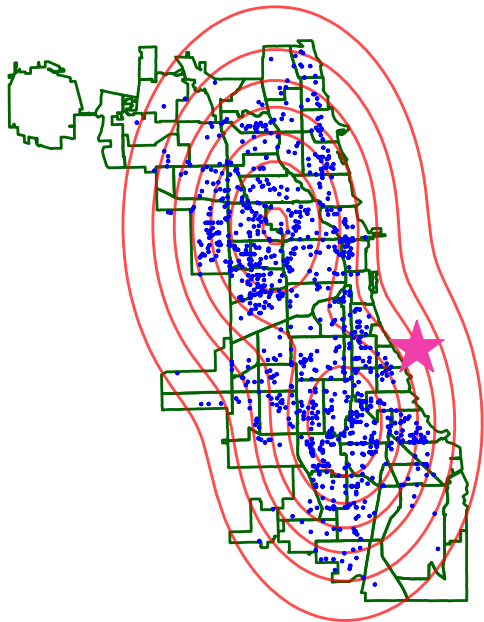
Model Criticism



Goals:

- 1 Test if a (complicated) **model** fits the **data**.
- 2 If it does not, show **a location** where it fails.

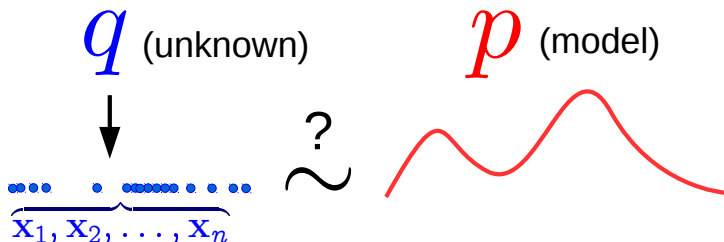
Model Criticism



Goals:

- 1 Test if a (complicated) **model** fits the **data**.
- 2 If it does not, show **a location** where it fails.

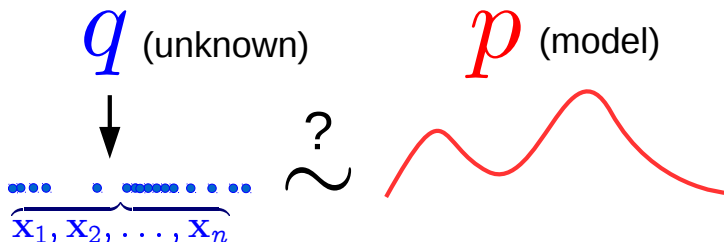
Problem Setting: Goodness-of-Fit Test



Test goal: Are data from the model p ?

- 1 Nonparametric.
- 2 Linear-time. Runtime is $\mathcal{O}(n)$. Fast.
- 3 Interpretable. Model criticism by finding .

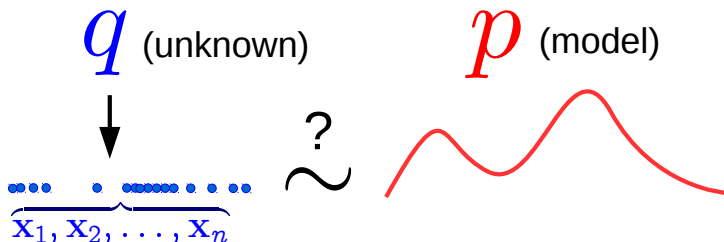
Problem Setting: Goodness-of-Fit Test




Test goal: Are **data** from the **model** p ?

- 1 Nonparametric.
- 2 Linear-time. Runtime is $\mathcal{O}(n)$. Fast.
- 3 Interpretable. Model criticism by finding  .

Problem Setting: Goodness-of-Fit Test



Test goal: Are **data** from the **model** p ?

- 1 **Nonparametric**.
- 2 **Linear-time**. Runtime is $\mathcal{O}(n)$. Fast.
- 3 **Interpretable**. Model criticism by finding  .

Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

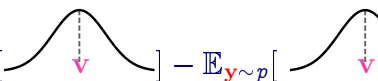
Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

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

Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

$$\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\text{bell curve} \right] - \mathbb{E}_{\mathbf{y} \sim p} \left[\text{bell curve} \right]$$


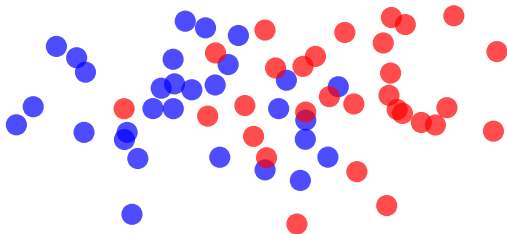
Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

$$\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\text{bell curve} \right] - \mathbb{E}_{\mathbf{y} \sim p} \left[\text{bell curve} \right]$$
$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

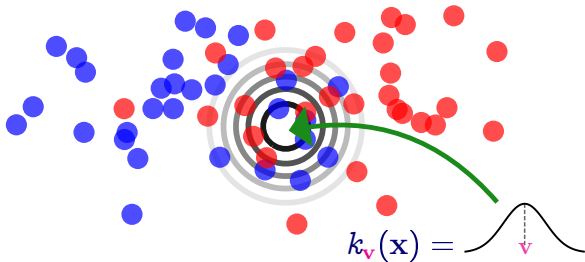


$$\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\text{bell curve}(\mathbf{v}) \right] - \mathbb{E}_{\mathbf{y} \sim p} \left[\text{bell curve}(\mathbf{v}) \right]$$
$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

score: 0.008



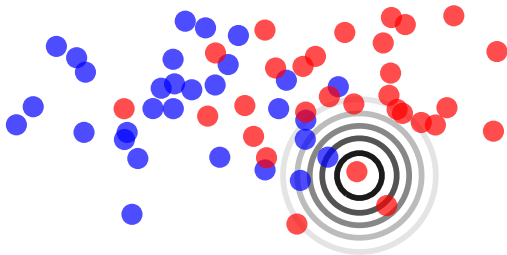
$$\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\text{kernel}(\mathbf{x}, \mathbf{v}) \right] - \mathbb{E}_{\mathbf{y} \sim p} \left[\text{kernel}(\mathbf{y}, \mathbf{v}) \right]$$

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

Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

score: 1.6

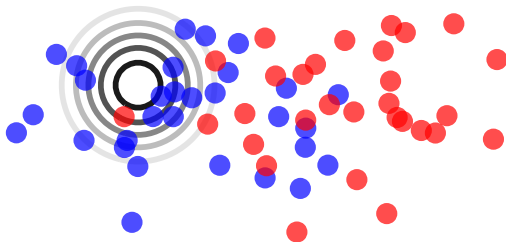


$$\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\text{curve}(\mathbf{v}) \right] - \mathbb{E}_{\mathbf{y} \sim p} \left[\text{curve}(\mathbf{v}) \right]$$
$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

score: 13

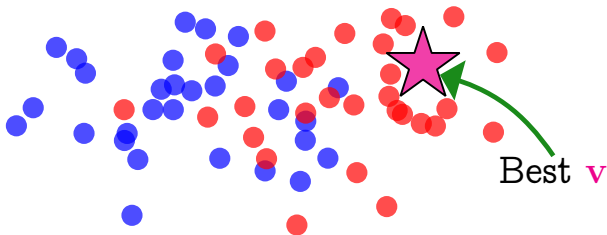


$$\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\text{curve}(\mathbf{v}) \right] - \mathbb{E}_{\mathbf{y} \sim p} \left[\text{curve}(\mathbf{v}) \right]$$
$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

score: 25

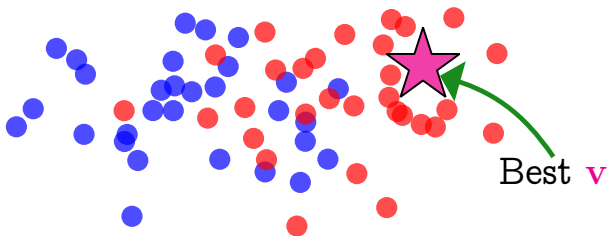


$$\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\text{bell curve}(\mathbf{x} - \mathbf{v}) \right] - \mathbb{E}_{\mathbf{y} \sim p} \left[\text{bell curve}(\mathbf{y} - \mathbf{v}) \right]$$
$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{standard deviation}(\mathbf{v})}.$$

Model Criticism by Maximum Mean Discrepancy [Gretton et al., 2012]

- Find a location \mathbf{v} at which q and p differ most [Jitkrittum et al., 2016].

score: 25



$$\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\text{bell curve}(\mathbf{v}) \right] - \mathbb{E}_{\mathbf{y} \sim p} \left[\text{bell curve}(\mathbf{v}) \right]$$

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

No sample from p .
Difficult to generate.

The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

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

The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

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


The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

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

The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]


Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

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


Idea: Define T_p such that $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$, for any \mathbf{v} .

The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

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


Idea: Define T_p such that $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$, for any \mathbf{v} .

The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

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

Idea: Define T_p such that $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$, for any \mathbf{v} .

The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

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

Idea: Define T_p such that $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$, for any \mathbf{v} .

Proposal: Good \mathbf{v} should have high

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

The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.


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

Idea: Define T_p such that $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$, for any \mathbf{v} .

Proposal: Good \mathbf{v} should have high

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

signal-to-noise
ratio



The Stein Witness Function [Liu et al., 2016, Chwialkowski et al., 2016]

Problem: No sample from p . Cannot estimate $\mathbb{E}_{\mathbf{y} \sim p}[k_{\mathbf{v}}(\mathbf{y})]$.

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

Idea: Define T_p such that $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$, for any \mathbf{v} .

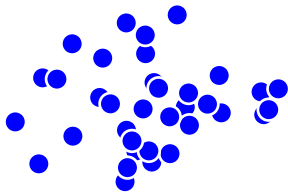
Proposal: Good \mathbf{v} should have high

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

signal-to-noise
ratio

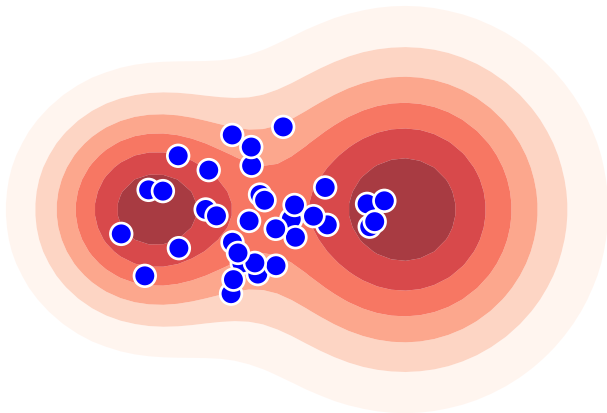
- $\text{witness}(\mathbf{v})$ and $\text{standard deviation}(\mathbf{v})$ can be estimated in linear-time.

Proposal: Model Criticism with the Stein Witness



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

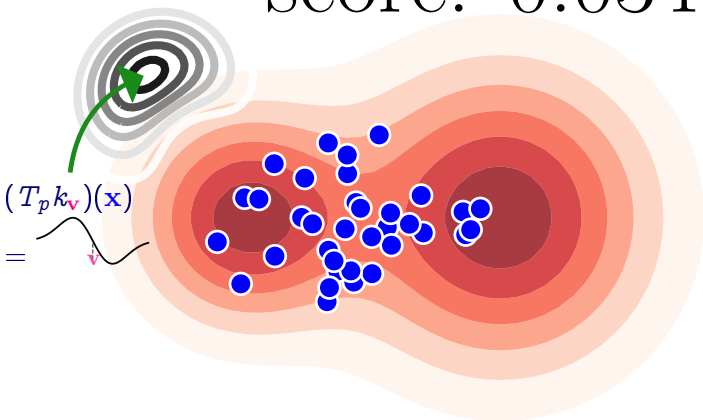
Proposal: Model Criticism with the Stein Witness



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

Proposal: Model Criticism with the Stein Witness

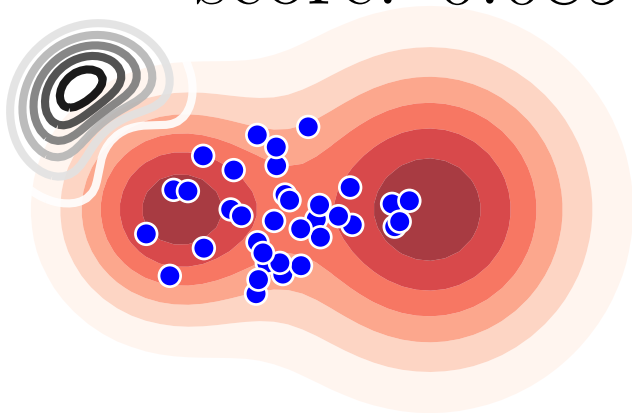
score: 0.034



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

Proposal: Model Criticism with the Stein Witness

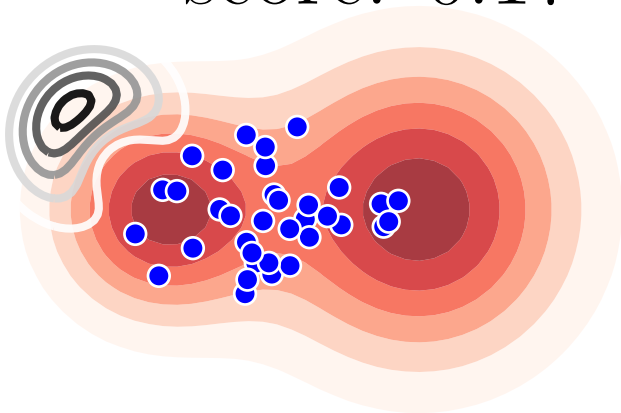
score: 0.089



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

Proposal: Model Criticism with the Stein Witness

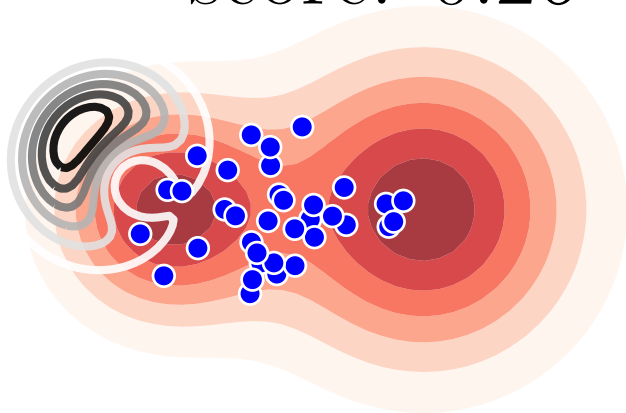
score: 0.17



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

Proposal: Model Criticism with the Stein Witness

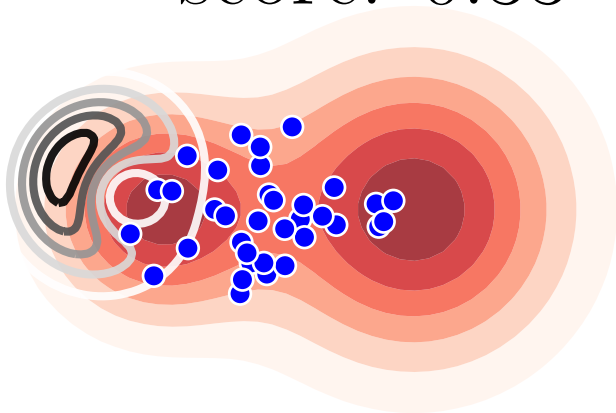
score: 0.26



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

Proposal: Model Criticism with the Stein Witness

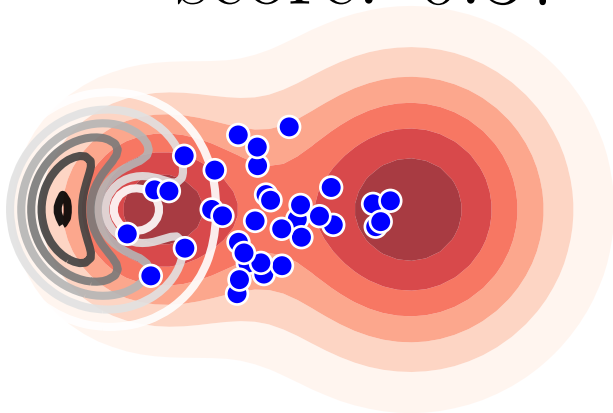
score: 0.33



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

Proposal: Model Criticism with the Stein Witness

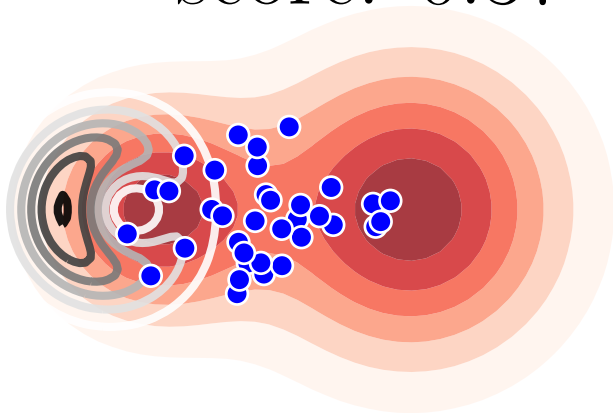
score: 0.37



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

Proposal: Model Criticism with the Stein Witness

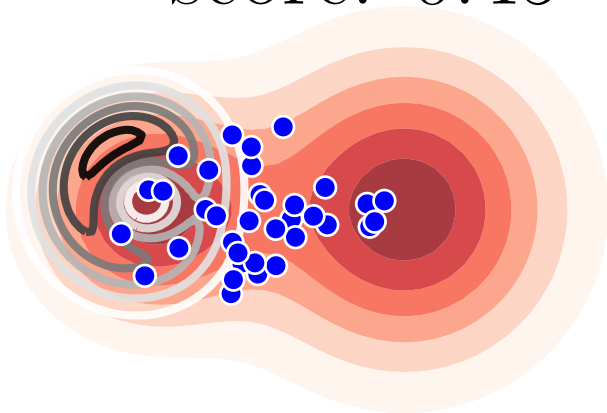
score: 0.37



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

Proposal: Model Criticism with the Stein Witness

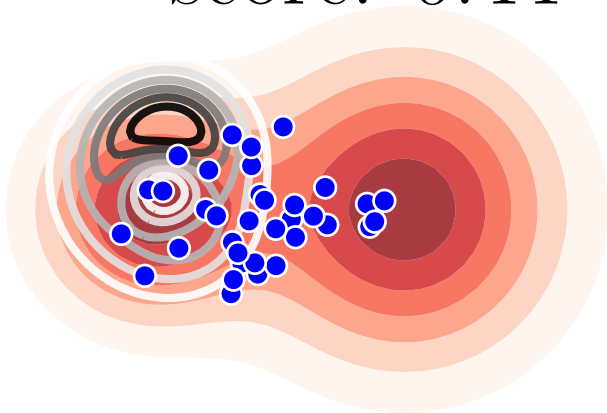
score: 0.45



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

Proposal: Model Criticism with the Stein Witness

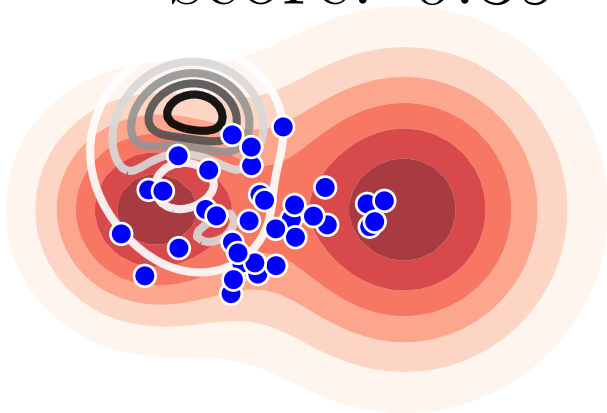
score: 0.44



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

Proposal: Model Criticism with the Stein Witness

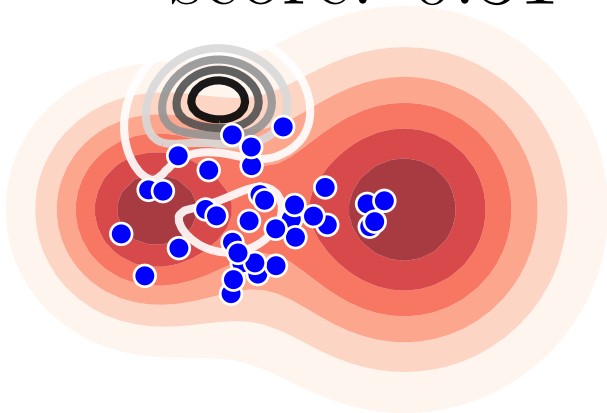
score: 0.39



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

Proposal: Model Criticism with the Stein Witness

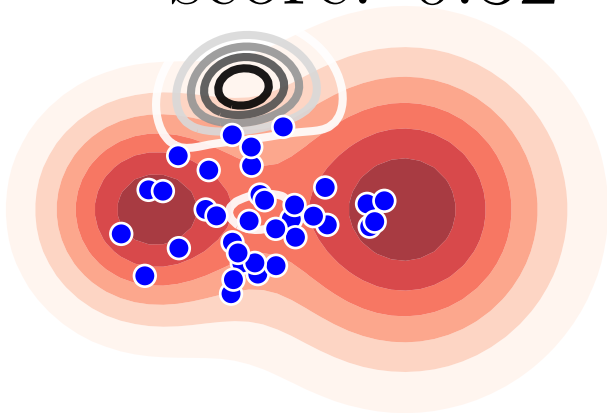
score: 0.31



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

Proposal: Model Criticism with the Stein Witness

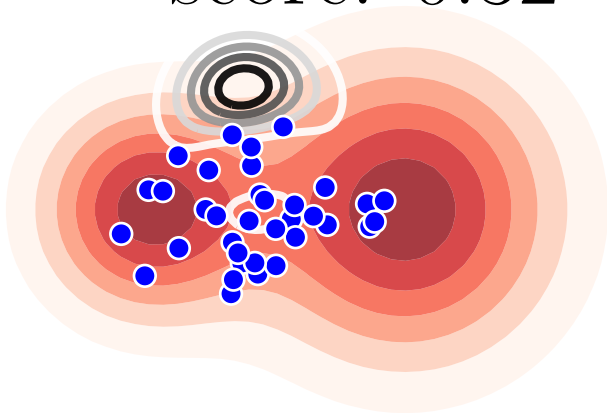
score: 0.32



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

Proposal: Model Criticism with the Stein Witness

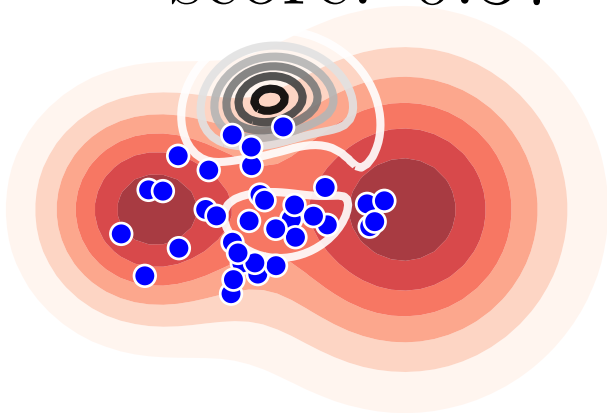
score: 0.32



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

Proposal: Model Criticism with the Stein Witness

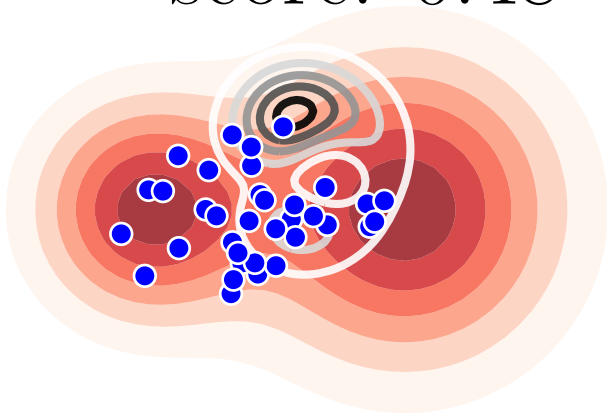
score: 0.37



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

Proposal: Model Criticism with the Stein Witness

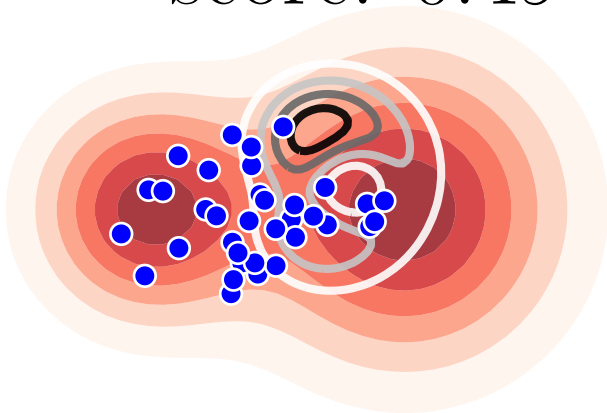
score: 0.48



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

Proposal: Model Criticism with the Stein Witness

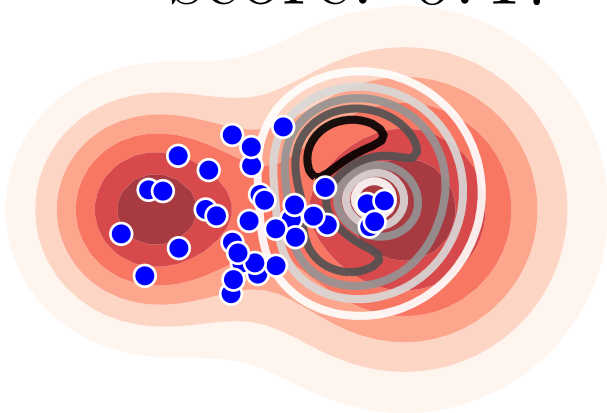
score: 0.49



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

Proposal: Model Criticism with the Stein Witness

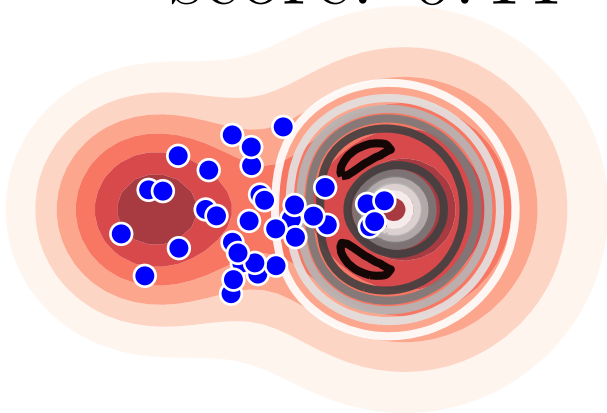
score: 0.47



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

Proposal: Model Criticism with the Stein Witness

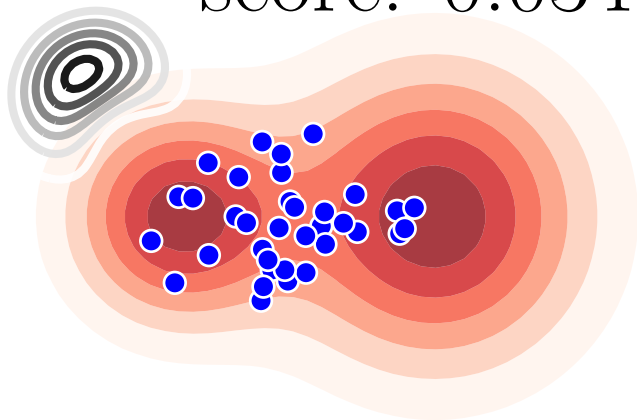
score: 0.44



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

Proposal: Model Criticism with the Stein Witness

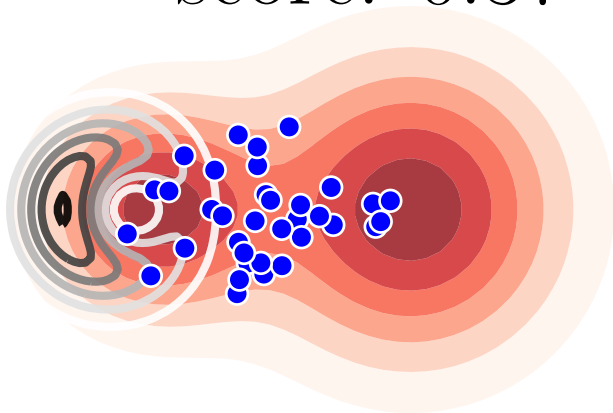
score: 0.034



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

Proposal: Model Criticism with the Stein Witness

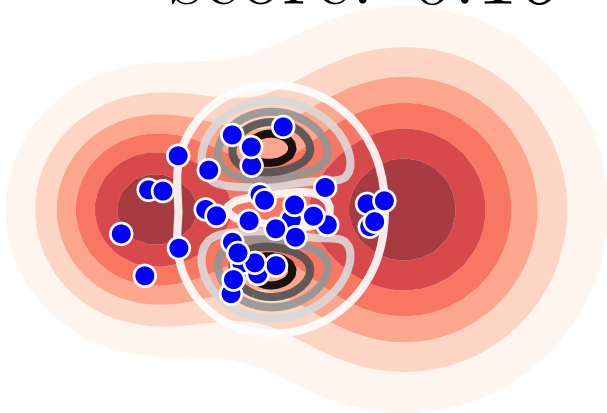
score: 0.37



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

Proposal: Model Criticism with the Stein Witness

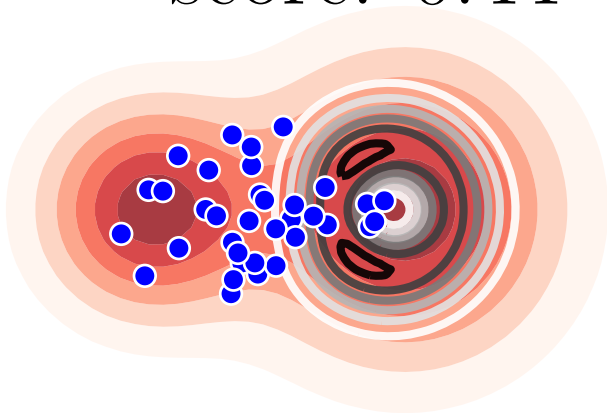
score: 0.16



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

Proposal: Model Criticism with the Stein Witness

score: 0.44



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

What is $T_p k_v$?

Recall $\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q}(T_p k_{\mathbf{v}})(\mathbf{x}) - \cancel{\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y})}$

What is $T_p k_v$?

Recall $\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q}(T_p k_{\mathbf{v}})(\mathbf{x}) - \cancel{\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y})}$

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k_{\mathbf{v}}(\mathbf{y}) p(\mathbf{y})].$$

Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

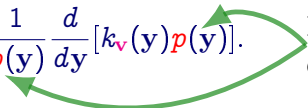
[Liu et al., 2016, Chwialkowski et al., 2016]

What is $T_p k_v$?

Recall $\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q}(T_p k_v)(\mathbf{x}) - \mathbb{E}_{\mathbf{y} \sim p}(T_p k_v)(\mathbf{y})$

$$(T_p k_v)(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k_v(\mathbf{y}) p(\mathbf{y})].$$

Normalizer
cancels



Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_v)(\mathbf{y}) = 0$.

[Liu et al., 2016, Chwialkowski et al., 2016]

Theorem: Maximizing

$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{uncertainty}(\mathbf{v})}$$

- increases true positive rate
= $\mathbb{P}(\text{detect difference when } p \neq q)$,
- does not affect false positive rate.

-
- General form of $\text{score}(\dots)$ can consider more than one location \mathbf{v} .

Theorem: Maximizing

$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{uncertainty}(\mathbf{v})}$$

- increases true positive rate
= $\mathbb{P}(\text{detect difference when } p \neq q)$,
- does not affect false positive rate.

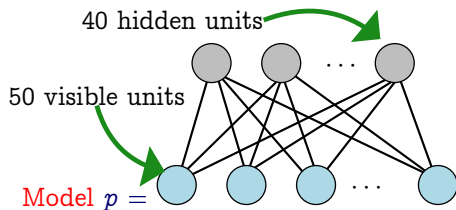
-
- General form of $\text{score}(\dots)$ can consider more than one location \mathbf{v} .

Theorem: Maximizing

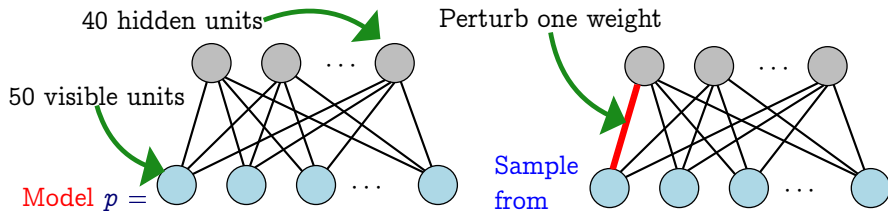
$$\text{score}(\mathbf{v}) = \frac{|\text{witness}(\mathbf{v})|}{\text{uncertainty}(\mathbf{v})}$$

- increases true positive rate
= $\mathbb{P}(\text{detect difference when } p \neq q)$,
 - does not affect false positive rate.
-
- General form of $\text{score}(\dots)$ can consider more than one location \mathbf{v} .

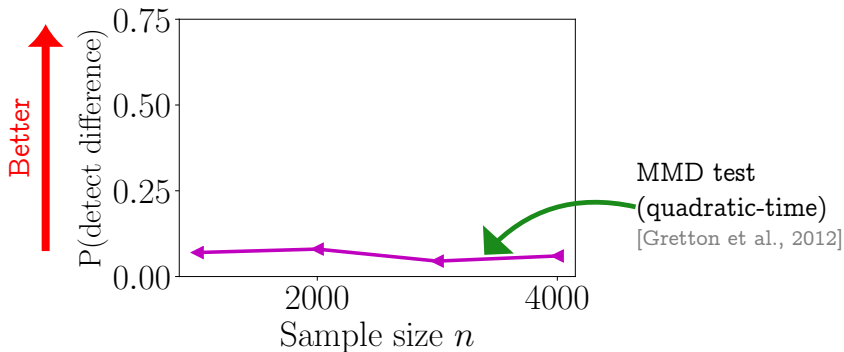
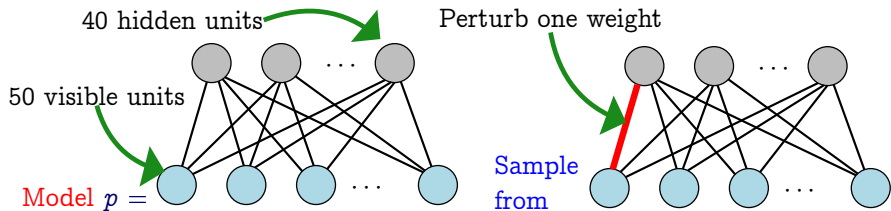
Experiment: Restricted Boltzmann Machine (RBM)



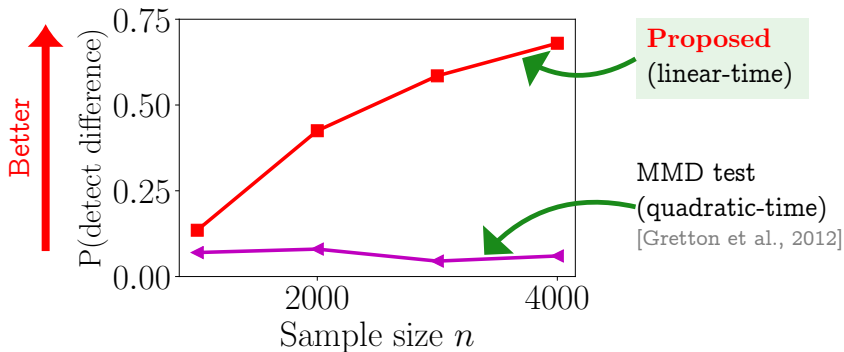
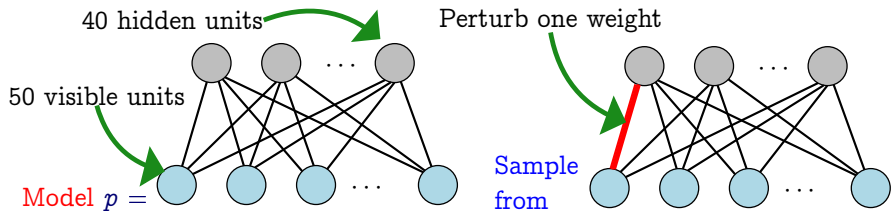
Experiment: Restricted Boltzmann Machine (RBM)



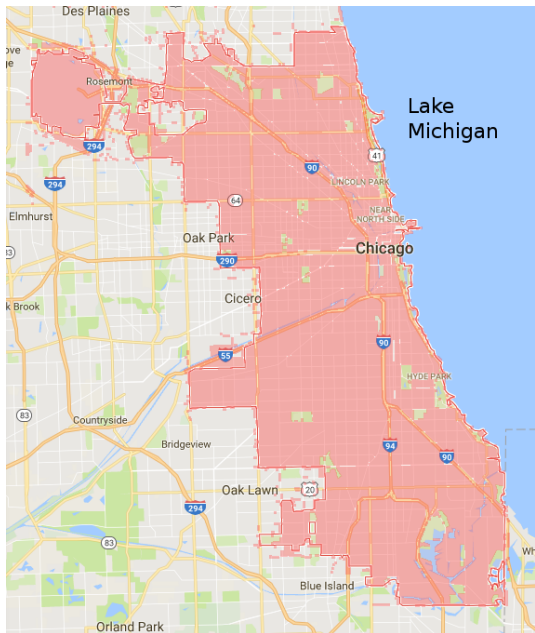
Experiment: Restricted Boltzmann Machine (RBM)



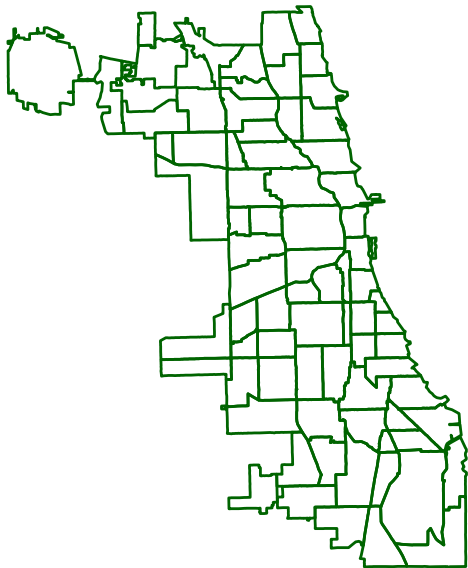
Experiment: Restricted Boltzmann Machine (RBM)



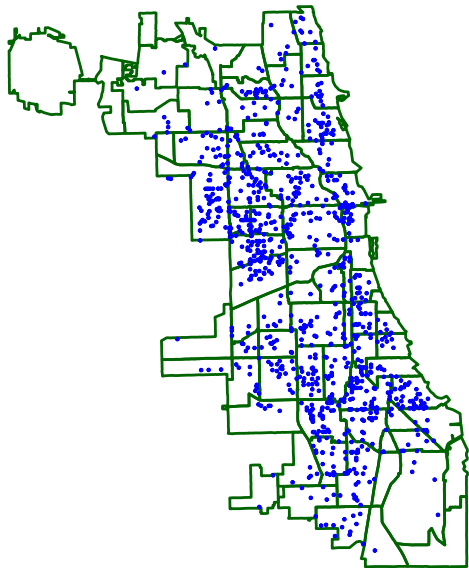
Interpretable Features: Chicago Crime



Interpretable Features: Chicago Crime

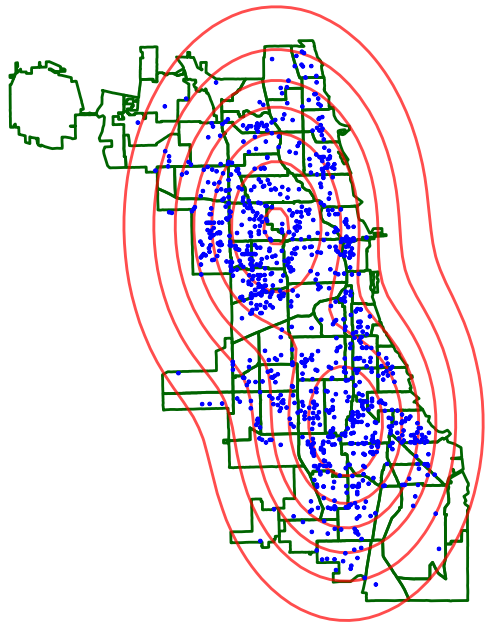


Interpretable Features: Chicago Crime



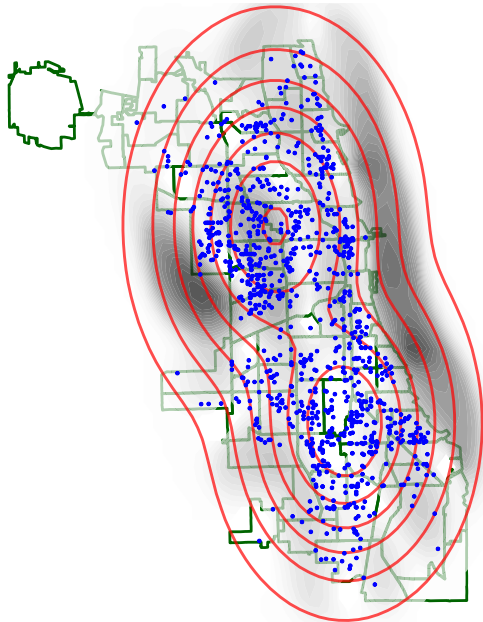
- $n = 11957$ robbery events in Chicago in 2016.
 - lat/long coordinates = sample from q .
- Model spatial density with Gaussian mixtures.

Interpretable Features: Chicago Crime



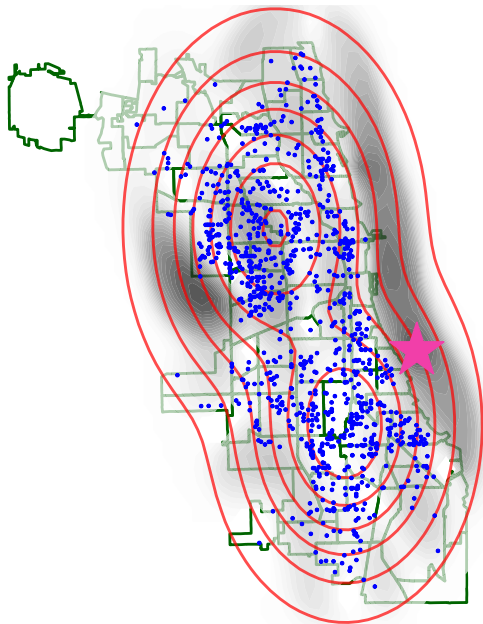
Model $p = 2$ -component Gaussian mixture.

Interpretable Features: Chicago Crime



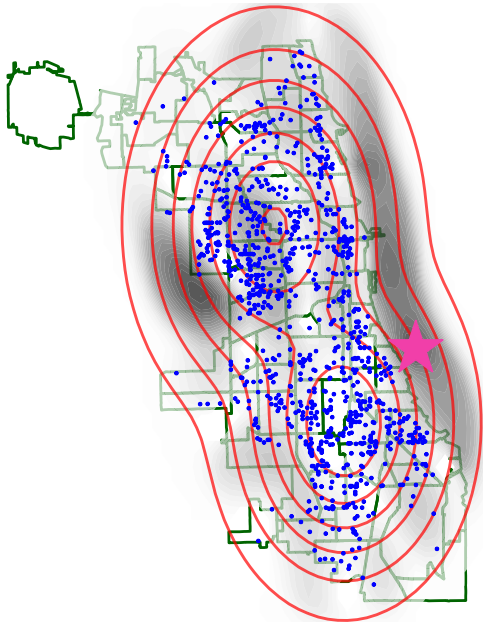
Score surface

Interpretable Features: Chicago Crime



★ = optimized \mathbf{v} .

Interpretable Features: Chicago Crime

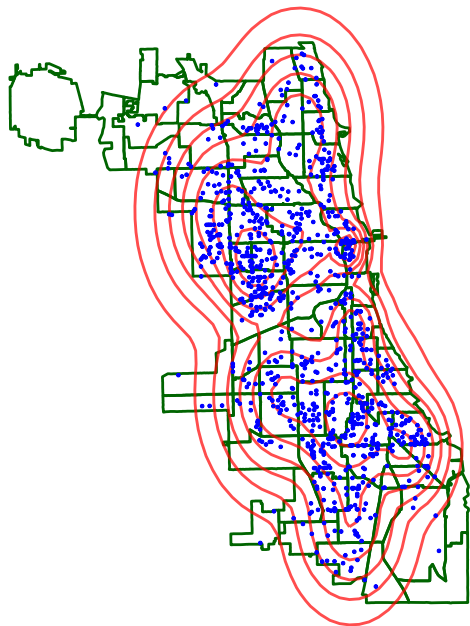


★ = optimized \mathbf{v} .

No robbery in Lake Michigan.

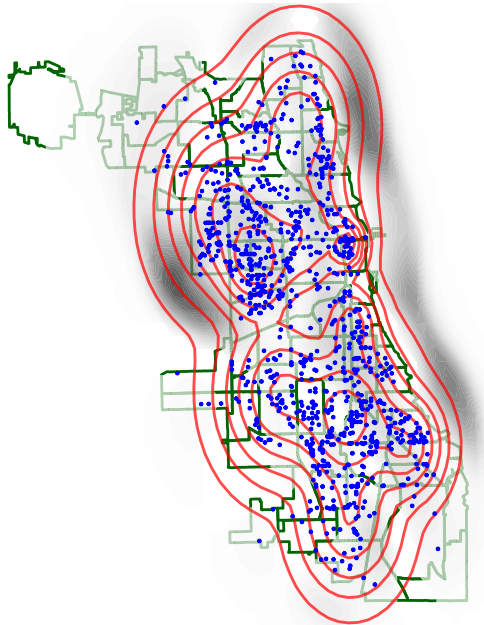


Interpretable Features: Chicago Crime



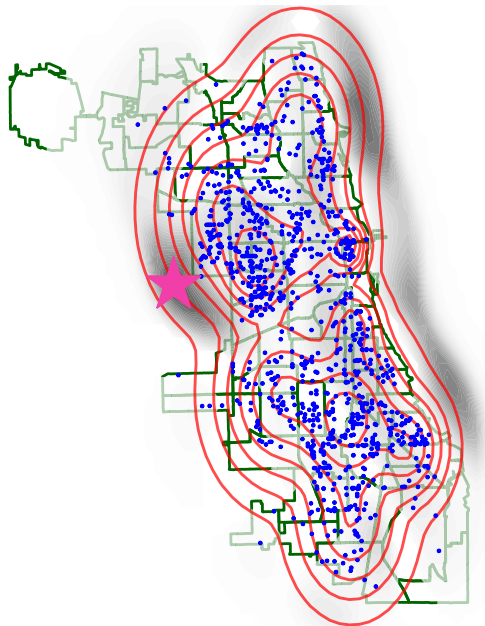
Model $p = 10$ -component Gaussian mixture.

Interpretable Features: Chicago Crime



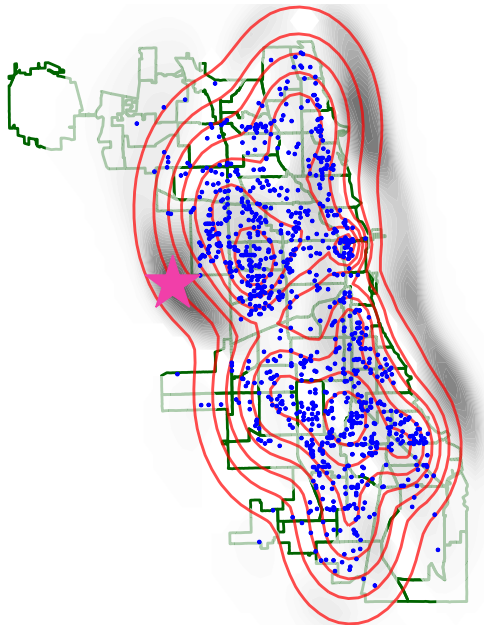
Capture the right tail better.

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.

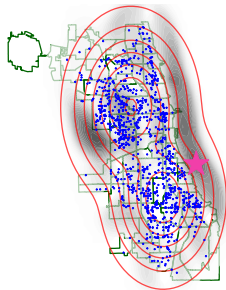
Conclusions

Proposed a new goodness-of-fit test.

- 1 Nonparametric. Normalizer not needed.
- 2 Linear-time
- 3 Interpretable

Poster #57 at Pacific Ballroom tonight.

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



Questions?

Thank you

FSSD and KSD in 1D Gaussian Case

Consider $p = \mathcal{N}(0, 1)$ and $q = \mathcal{N}(\mu_q, \sigma_q^2)$.

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

$$\text{FSSD}^2 = \frac{\sigma_k^2 e^{-\frac{(v - \mu_q)^2}{\sigma_k^2 + \sigma_q^2}} \left((\sigma_k^2 + 1) \mu_q + v (\sigma_q^2 - 1) \right)^2}{(\sigma_k^2 + \sigma_q^2)^3}.$$

- If $\mu_q \neq 0, \sigma_q^2 \neq 1$, and $v = -\frac{(\sigma_k^2 + 1)\mu_q}{(\sigma_q^2 - 1)}$, then $\text{FSSD}^2 = 0$!
 - This is why v should be drawn from a distribution with a density.
- For KSD, Gaussian kernel (bandwidth = κ^2).

$$S^2 = \frac{\mu_q^2 (\kappa^2 + 2\sigma_q^2) + (\sigma_q^2 - 1)^2}{(\kappa^2 + 2\sigma_q^2) \sqrt{\frac{2\sigma_q^2}{\kappa^2} + 1}}.$$

FSSD and KSD in 1D Gaussian Case

Consider $p = \mathcal{N}(0, 1)$ and $q = \mathcal{N}(\mu_q, \sigma_q^2)$.

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

$$\text{FSSD}^2 = \frac{\sigma_k^2 e^{-\frac{(v - \mu_q)^2}{\sigma_k^2 + \sigma_q^2}} \left((\sigma_k^2 + 1) \mu_q + v (\sigma_q^2 - 1) \right)^2}{(\sigma_k^2 + \sigma_q^2)^3}.$$

- If $\mu_q \neq 0, \sigma_q^2 \neq 1$, and $v = -\frac{(\sigma_k^2 + 1)\mu_q}{(\sigma_q^2 - 1)}$, then $\text{FSSD}^2 = 0$!
 - This is why v should be drawn from a distribution with a density.
- For KSD, Gaussian kernel (bandwidth = κ^2).

$$S^2 = \frac{\mu_q^2 (\kappa^2 + 2\sigma_q^2) + (\sigma_q^2 - 1)^2}{(\kappa^2 + 2\sigma_q^2) \sqrt{\frac{2\sigma_q^2}{\kappa^2} + 1}}.$$

What is $T_p k_v$?

Recall $\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q}(T_p k_{\mathbf{v}})(\mathbf{x}) - \cancel{\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y})}$

What is $T_p k_v$?

Recall $\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q}(T_p k_{\mathbf{v}})(\mathbf{x}) - \cancel{\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y})}$

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

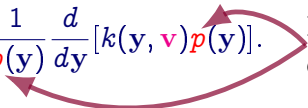
[Liu et al., 2016, Chwialkowski et al., 2016]

What is $T_p k_v$?

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

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Normalizer
cancels



Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

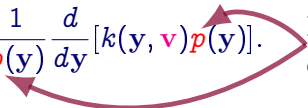
[Liu et al., 2016, Chwialkowski et al., 2016]

What is $T_p k_v$?

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

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Normalizer
cancels



Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

[Liu et al., 2016, Chwialkowski et al., 2016]

Proof:

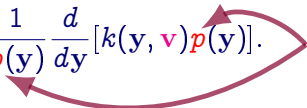
$$\mathbb{E}_{\mathbf{y} \sim p} [(T_p k_{\mathbf{v}})(\mathbf{y})]$$

What is $T_p k_v$?

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

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Normalizer
cancels



Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

[Liu et al., 2016, Chwialkowski et al., 2016]

Proof:

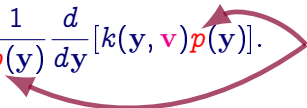
$$\mathbb{E}_{\mathbf{y} \sim p} [(T_p k_{\mathbf{v}})(\mathbf{y})] = \int_{-\infty}^{\infty} [(T_p k_{\mathbf{v}})(\mathbf{y})] p(\mathbf{y}) d\mathbf{y}$$

What is $T_p k_v$?

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

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Normalizer cancels



Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

[Liu et al., 2016, Chwialkowski et al., 2016]

Proof:

$$\mathbb{E}_{\mathbf{y} \sim p} [(T_p k_{\mathbf{v}})(\mathbf{y})] = \int_{-\infty}^{\infty} \left[\frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k_{\mathbf{v}}(\mathbf{y}) p(\mathbf{y})] \right] p(\mathbf{y}) d\mathbf{y}$$

What is $T_p k_v$?

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

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Normalizer cancels

Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

[Liu et al., 2016, Chwialkowski et al., 2016]

Proof:

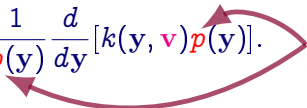
$$\mathbb{E}_{\mathbf{y} \sim p} [(T_p k_{\mathbf{v}})(\mathbf{y})] = \int_{-\infty}^{\infty} \left[\frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k_{\mathbf{v}}(\mathbf{y}) p(\mathbf{y})] \right] p(\mathbf{y}) d\mathbf{y}$$

What is $T_p k_v$?

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

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Normalizer
cancels



Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

[Liu et al., 2016, Chwialkowski et al., 2016]

Proof:

$$\begin{aligned} \mathbb{E}_{\mathbf{y} \sim p} [(T_p k_{\mathbf{v}})(\mathbf{y})] &= \int_{-\infty}^{\infty} \left[\frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k_{\mathbf{v}}(\mathbf{y}) p(\mathbf{y})] \right] p(\mathbf{y}) d\mathbf{y} \\ &= \int_{-\infty}^{\infty} \frac{d}{d\mathbf{y}} [k_{\mathbf{v}}(\mathbf{y}) p(\mathbf{y})] d\mathbf{y} \end{aligned}$$

What is $T_p k_v$?

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

$$(T_p k_{\mathbf{v}})(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Normalizer cancels

Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_{\mathbf{v}})(\mathbf{y}) = 0$.

[Liu et al., 2016, Chwialkowski et al., 2016]

Proof:

$$\begin{aligned} \mathbb{E}_{\mathbf{y} \sim p} [(T_p k_{\mathbf{v}})(\mathbf{y})] &= \int_{-\infty}^{\infty} \left[\frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k_{\mathbf{v}}(\mathbf{y}) p(\mathbf{y})] \right] p(\mathbf{y}) d\mathbf{y} \\ &= \int_{-\infty}^{\infty} \frac{d}{d\mathbf{y}} [k_{\mathbf{v}}(\mathbf{y}) p(\mathbf{y})] d\mathbf{y} \\ &= [k_{\mathbf{v}}(\mathbf{y}) p(\mathbf{y})]_{\mathbf{y}=-\infty}^{\mathbf{y}=\infty} \end{aligned}$$

What is $T_p k_v$?

Recall $\text{witness}(\mathbf{v}) = \mathbb{E}_{\mathbf{x} \sim q}(T_p k_v)(\mathbf{x}) - \mathbb{E}_{\mathbf{y} \sim p}(T_p k_v)(\mathbf{y})$

$$(T_p k_v)(\mathbf{y}) = \frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k(\mathbf{y}, \mathbf{v}) p(\mathbf{y})].$$

Normalizer cancels

Then, $\mathbb{E}_{\mathbf{y} \sim p}(T_p k_v)(\mathbf{y}) = 0$.

[Liu et al., 2016, Chwialkowski et al., 2016]

Proof:

$$\begin{aligned}\mathbb{E}_{\mathbf{y} \sim p} [(T_p k_v)(\mathbf{y})] &= \int_{-\infty}^{\infty} \left[\frac{1}{p(\mathbf{y})} \frac{d}{d\mathbf{y}} [k_v(\mathbf{y}) p(\mathbf{y})] \right] p(\mathbf{y}) d\mathbf{y} \\ &= \int_{-\infty}^{\infty} \frac{d}{d\mathbf{y}} [k_v(\mathbf{y}) p(\mathbf{y})] d\mathbf{y} \\ &= [k_v(\mathbf{y}) p(\mathbf{y})]_{\mathbf{y}=-\infty}^{\mathbf{y}=\infty} \\ &= 0\end{aligned}$$

(assume $\lim_{|\mathbf{y}| \rightarrow \infty} k(\mathbf{y}, \mathbf{v}) p(\mathbf{y}) = 0$)

FSSD is a Discrepancy Measure

Theorem 1.

Let $V = \{v_1, \dots, v_J\} \subset \mathbb{R}^d$ be drawn i.i.d. from a distribution η 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 (Stein witness not too rough) $\|g\|_k^2 < \infty$.
- 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) $\lim_{\|\mathbf{x}\| \rightarrow \infty} p(\mathbf{x})g(\mathbf{x}) = 0$.

Then, for any $J \geq 1$, η -almost surely

$$\text{FSSD}^2 = 0 \text{ if and only if } p = q.$$

- 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 $J = 5$.

What Are “Blind Spots”?

$$\begin{aligned} \mathbf{g}(\mathbf{v}) &:= \mathbb{E}_{\mathbf{x} \sim q} \left[\frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k_{\mathbf{v}}(\mathbf{x}) p(\mathbf{x})] \right] \\ &= \mathbb{E}_{\mathbf{x} \sim q} \left[\left(\frac{d}{d\mathbf{x}} \log p(\mathbf{x}) \right) k_{\mathbf{v}}(\mathbf{x}) + \partial_{\mathbf{x}} k_{\mathbf{v}}(\mathbf{x}) \right] \in \mathbb{R}^d. \end{aligned}$$

Consider $p = \mathcal{N}(0, 1)$ and $q = \mathcal{N}(0, \sigma_q^2)$. Use unit-width Gaussian kernel.

$$g(v) = \frac{v \exp\left(-\frac{v^2}{2+2\sigma_q^2}\right) (\sigma_q^2 - 1)}{(1 + \sigma_q^2)^{3/2}}$$

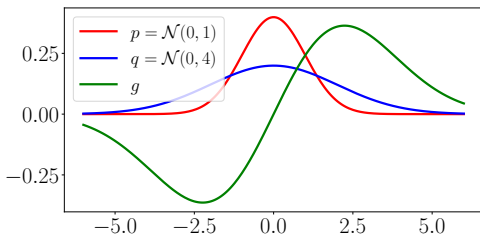
- If $v = 0$, then $\text{FSSD}^2 = g^2(v) = 0$ regardless of σ_q^2 .
- If $g \neq 0$, and k is real analytic, $R = \{v \mid g(v) = 0\}$ (blind spots) has 0 Lebesgue measure.
- So, if $v \sim$ a distribution with a density, then $v \notin R$.

What Are “Blind Spots”?

$$\begin{aligned} \mathbf{g}(\mathbf{v}) &:= \mathbb{E}_{\mathbf{x} \sim q} \left[\frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k_{\mathbf{v}}(\mathbf{x}) p(\mathbf{x})] \right] \\ &= \mathbb{E}_{\mathbf{x} \sim q} \left[\left(\frac{d}{d\mathbf{x}} \log p(\mathbf{x}) \right) k_{\mathbf{v}}(\mathbf{x}) + \partial_{\mathbf{x}} k_{\mathbf{v}}(\mathbf{x}) \right] \in \mathbb{R}^d. \end{aligned}$$

Consider $p = \mathcal{N}(0, 1)$ and $q = \mathcal{N}(0, \sigma_q^2)$. Use unit-width Gaussian kernel.

$$g(v) = \frac{v \exp\left(-\frac{v^2}{2+2\sigma_q^2}\right) (\sigma_q^2 - 1)}{(1 + \sigma_q^2)^{3/2}}$$



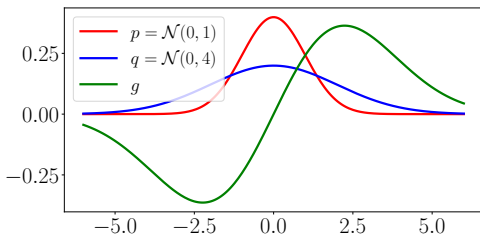
- If $v = 0$, then $\text{FSSD}^2 = g^2(v) = 0$ regardless of σ_q^2 .
- If $g \neq 0$, and k is real analytic, $R = \{v \mid g(v) = 0\}$ (blind spots) has 0 Lebesgue measure.
- So, if $v \sim$ a distribution with a density, then $v \notin R$.

What Are “Blind Spots”?

$$\begin{aligned} \mathbf{g}(\mathbf{v}) &:= \mathbb{E}_{\mathbf{x} \sim q} \left[\frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k_{\mathbf{v}}(\mathbf{x}) p(\mathbf{x})] \right] \\ &= \mathbb{E}_{\mathbf{x} \sim q} \left[\left(\frac{d}{d\mathbf{x}} \log p(\mathbf{x}) \right) k_{\mathbf{v}}(\mathbf{x}) + \partial_{\mathbf{x}} k_{\mathbf{v}}(\mathbf{x}) \right] \in \mathbb{R}^d. \end{aligned}$$

Consider $p = \mathcal{N}(0, 1)$ and $q = \mathcal{N}(0, \sigma_q^2)$. Use unit-width Gaussian kernel.

$$g(v) = \frac{v \exp\left(-\frac{v^2}{2+2\sigma_q^2}\right) (\sigma_q^2 - 1)}{(1 + \sigma_q^2)^{3/2}}$$



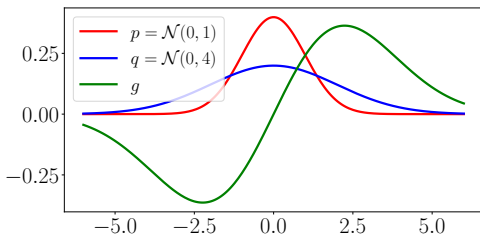
- If $v = 0$, then $\text{FSSD}^2 = g^2(v) = 0$ regardless of σ_q^2 .
- If $g \neq 0$, and k is real analytic, $R = \{v \mid g(v) = 0\}$ (blind spots) has 0 Lebesgue measure.
- So, if $v \sim$ a distribution with a density, then $v \notin R$.

What Are “Blind Spots”?

$$\begin{aligned} \mathbf{g}(\mathbf{v}) &:= \mathbb{E}_{\mathbf{x} \sim q} \left[\frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k_{\mathbf{v}}(\mathbf{x}) p(\mathbf{x})] \right] \\ &= \mathbb{E}_{\mathbf{x} \sim q} \left[\left(\frac{d}{d\mathbf{x}} \log p(\mathbf{x}) \right) k_{\mathbf{v}}(\mathbf{x}) + \partial_{\mathbf{x}} k_{\mathbf{v}}(\mathbf{x}) \right] \in \mathbb{R}^d. \end{aligned}$$

Consider $p = \mathcal{N}(0, 1)$ and $q = \mathcal{N}(0, \sigma_q^2)$. Use unit-width Gaussian kernel.

$$g(v) = \frac{v \exp\left(-\frac{v^2}{2+2\sigma_q^2}\right) (\sigma_q^2 - 1)}{(1 + \sigma_q^2)^{3/2}}$$



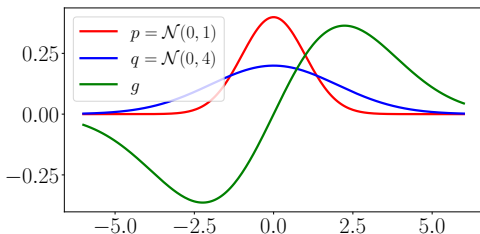
- If $v = 0$, then $\text{FSSD}^2 = g^2(v) = 0$ regardless of σ_q^2 .
- If $g \neq 0$, and k is real analytic, $R = \{v \mid g(v) = 0\}$ (blind spots) has 0 Lebesgue measure.
- So, if $v \sim$ a distribution with a density, then $v \notin R$.

What Are “Blind Spots”?

$$\begin{aligned} \mathbf{g}(\mathbf{v}) &:= \mathbb{E}_{\mathbf{x} \sim q} \left[\frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k_{\mathbf{v}}(\mathbf{x}) p(\mathbf{x})] \right] \\ &= \mathbb{E}_{\mathbf{x} \sim q} \left[\left(\frac{d}{d\mathbf{x}} \log p(\mathbf{x}) \right) k_{\mathbf{v}}(\mathbf{x}) + \partial_{\mathbf{x}} k_{\mathbf{v}}(\mathbf{x}) \right] \in \mathbb{R}^d. \end{aligned}$$

Consider $p = \mathcal{N}(0, 1)$ and $q = \mathcal{N}(0, \sigma_q^2)$. Use unit-width Gaussian kernel.

$$g(v) = \frac{v \exp\left(-\frac{v^2}{2+2\sigma_q^2}\right) (\sigma_q^2 - 1)}{(1 + \sigma_q^2)^{3/2}}$$



- If $v = 0$, then $\text{FSSD}^2 = g^2(v) = 0$ regardless of σ_q^2 .
- If $g \neq 0$, and k is real analytic, $R = \{v \mid g(v) = 0\}$ (blind spots) has 0 Lebesgue measure.
- So, if $v \sim$ a distribution with a density, then $v \notin R$.

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

- Recall $\xi(\mathbf{x}, \mathbf{v}) := \frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k(\mathbf{x}, \mathbf{v}) p(\mathbf{x})] \in \mathbb{R}^d$.
- $\tau(\mathbf{x}) :=$ vertically stack $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$. Features of \mathbf{x} .
- Mean feature: $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$.
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$ 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 Σ_p .

- 1 Under $H_0 : p = q$, asymptotically $n \widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1) \omega_i$.
 - 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 Σ_p ? No sample from p !

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

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

- Recall $\xi(\mathbf{x}, \mathbf{v}) := \frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k(\mathbf{x}, \mathbf{v}) p(\mathbf{x})] \in \mathbb{R}^d$.
- $\tau(\mathbf{x}) :=$ vertically stack $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$. Features of \mathbf{x} .
- Mean feature: $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$.
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$ 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 Σ_p .

- 1 Under $H_0 : p = q$, asymptotically $n \widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1) \omega_i$.

- 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 Σ_p ? No sample from p !

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

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

- Recall $\xi(\mathbf{x}, \mathbf{v}) := \frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k(\mathbf{x}, \mathbf{v}) p(\mathbf{x})] \in \mathbb{R}^d$.
- $\tau(\mathbf{x}) :=$ vertically stack $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$. Features of \mathbf{x} .
- Mean feature: $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$.
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$ 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 Σ_p .

- 1 Under $H_0 : p = q$, asymptotically $n \widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1) \omega_i$.
 - 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 Σ_p ? No sample from p !

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

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

- Recall $\xi(\mathbf{x}, \mathbf{v}) := \frac{1}{p(\mathbf{x})} \frac{d}{d\mathbf{x}} [k(\mathbf{x}, \mathbf{v}) p(\mathbf{x})] \in \mathbb{R}^d$.
- $\tau(\mathbf{x}) :=$ vertically stack $\xi(\mathbf{x}, \mathbf{v}_1), \dots, \xi(\mathbf{x}, \mathbf{v}_J) \in \mathbb{R}^{dJ}$. Features of \mathbf{x} .
- Mean feature: $\mu := \mathbb{E}_{\mathbf{x} \sim q}[\tau(\mathbf{x})]$.
- $\Sigma_r := \text{cov}_{\mathbf{x} \sim r}[\tau(\mathbf{x})] \in \mathbb{R}^{dJ \times dJ}$ 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 Σ_p .

- 1 Under $H_0 : p = q$, asymptotically $n \widehat{\text{FSSD}}^2 \xrightarrow{d} \sum_{i=1}^{dJ} (Z_i^2 - 1) \omega_i$.
 - 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 Σ_p ? No sample from p !

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

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 $c(0) = 0$ [Bahadur, 1960].
- $c(\theta)$ higher \implies more sensitive. Good.

Bahadur slope

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

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 $\text{pval}_n \approx \exp\left(-\frac{1}{2}c(\theta)n\right)$ where $c(\theta) > 0$ under H_1 , and $c(0) = 0$ [Bahadur, 1960].
- $c(\theta)$ higher \implies more sensitive. Good.

Bahadur slope

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

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

- Bahadur efficiency = ratio of slopes of two tests.

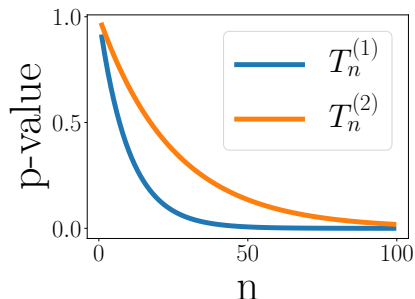
Bahadur Slope and Bahadur Efficiency

- 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 $\text{pval}_n \approx \exp\left(-\frac{1}{2}c(\theta)n\right)$ where $c(\theta) > 0$ under H_1 , and $c(0) = 0$ [Bahadur, 1960].
- $c(\theta)$ higher \implies more sensitive. Good.



Bahadur slope

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

where $F(t)$ = CDF of T_n under H_0 .

- Bahadur efficiency = ratio of slopes of two tests.

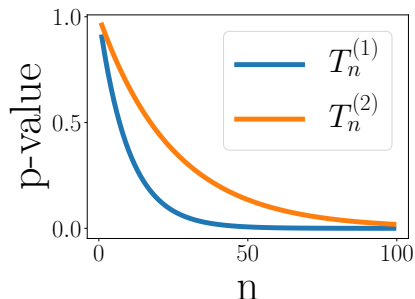
Bahadur Slope and Bahadur Efficiency

- 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 $\text{pval}_n \approx \exp\left(-\frac{1}{2}c(\theta)n\right)$ where $c(\theta) > 0$ under H_1 , and $c(0) = 0$ [Bahadur, 1960].
- $c(\theta)$ higher \Rightarrow more sensitive. Good.



Bahadur slope

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

where $F(t)$ = CDF of T_n under H_0 .

- Bahadur efficiency = ratio of slopes of two tests.

Gaussian Mean Shift Problem

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

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

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

Theorem 2 (FSSD is at least two times more efficient).

Fix $\sigma_k^2 = 1$ for $\widehat{n\text{FSSD}}^2$. Then, $\forall \mu_q \neq 0, \exists v \in \mathbb{R}, \forall \kappa^2 > 0$, we have Bahadur efficiency

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

Gaussian Mean Shift Problem

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

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

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

Theorem 2 (FSSD is at least two times more efficient).

Fix $\sigma_k^2 = 1$ for $\widehat{n\text{FSSD}}^2$. Then, $\forall \mu_q \neq 0, \exists v \in \mathbb{R}, \forall \kappa^2 > 0$, we have Bahadur efficiency

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

Gaussian Mean Shift Problem

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

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

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

Theorem 2 (FSSD is at least two times more efficient).

Fix $\sigma_k^2 = 1$ for $\widehat{n\text{FSSD}}^2$. Then, $\forall \mu_q \neq 0, \exists v \in \mathbb{R}, \forall \kappa^2 > 0$, we have Bahadur efficiency

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

Linear-Time Kernel Stein Discrepancy (LKS)

- [Liu et al., 2016] also proposed a linear version of KSD.
- For $\{\mathbf{x}_i\}_{i=1}^n \sim q$, KSD test statistic is

$$\frac{2}{n(n-1)} \sum_{i < j} h_p(\mathbf{x}_i, \mathbf{x}_j).$$

	1	2	3	4	5	6	7	8
1								
2								
3								
4								
5								
6								
7								
8								

- LKS test statistic is a “running average”

$$\frac{2}{n} \sum_{i=1}^{n/2} h_p(\mathbf{x}_{2i-1}, \mathbf{x}_{2i}).$$

	1	2	3	4	5	6	7	8
1								
2								
3								
4								
5								
6								
7								
8								

- Both unbiased. LKS has $\mathcal{O}(d^2 n)$ runtime.
- ~~X~~ LKS has high variance. Poor test power.

Bahadur Slopes of FSSD and LKS

Theorem 3.

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

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

where ω_1 is the maximum eigenvalue of $\Sigma_p := \text{cov}_{\mathbf{x} \sim p}[\tau(\mathbf{x})]$.

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.

Illustration: Optimization Objective

- Consider $J = 1$ location.
- Training objective $\frac{\widehat{\text{FSSD}}^2(\mathbf{v})}{\widehat{\sigma}_{H_1}(\mathbf{v})}$ (gray), p in wireframe, $\{\mathbf{x}_i\}_{i=1}^n \sim q$ in purple, ★ = best \mathbf{v} .

$$p = \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right) \text{ vs. } q = \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}\right).$$

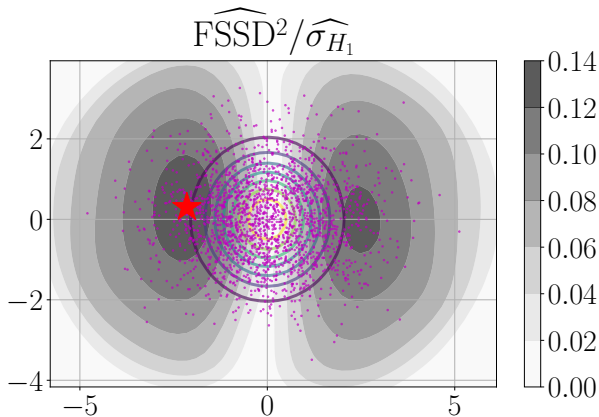
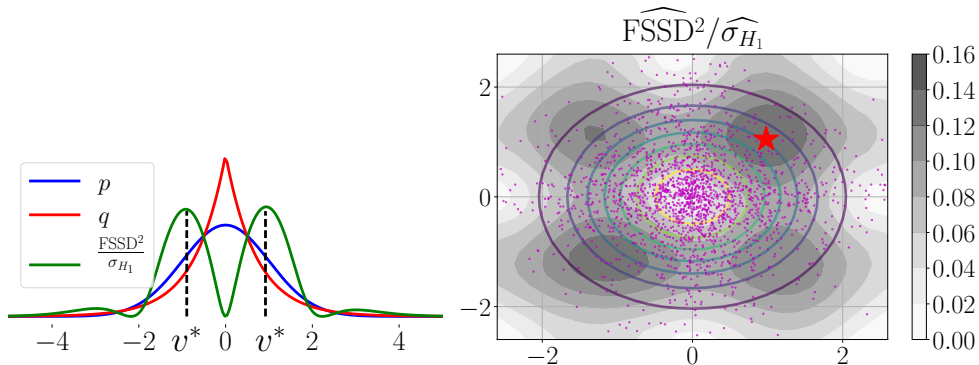


Illustration: Optimization Objective

- Consider $J = 1$ location.
- Training objective $\frac{\widehat{\text{FSSD}}^2(\mathbf{v})}{\widehat{\sigma}_{H_1}(\mathbf{v})}$ (gray), p in wireframe, $\{\mathbf{x}_i\}_{i=1}^n \sim q$ in purple, ★ = best \mathbf{v} .

$p = \mathcal{N}(\mathbf{0}, \mathbf{I})$ vs. $q = \text{Laplace}$ with same mean & variance.



Simulation Settings

- Gaussian kernel $k(\mathbf{x}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{v}\|_2^2}{2\sigma_k^2}\right)$

	Method	Description
1	FSSD-opt	Proposed. With optimization. $J = 5$.
2	FSSD-rand	Proposed. Random test locations.
3	KSD	Quadratic-time kernel Stein discrepancy [Liu et al., 2016, Chwialkowski et al., 2016]
4	LKS	Linear-time running average version of KSD.
5	MMD-opt	MMD two-sample test [Gretton et al., 2012]. With optimization.
6	ME-test	<u>Mean Embeddings</u> two-sample test [Jitkrittum et al., 2016]. With optimization.

- Two-sample tests need to draw sample from p .
- Tests with optimization use 20% of the data.
- $\alpha = 0.05$. 200 trials.

Simulation Settings

- Gaussian kernel $k(\mathbf{x}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{v}\|_2^2}{2\sigma_k^2}\right)$

	Method	Description
1	FSSD-opt	Proposed. With optimization. $J = 5$.
2	FSSD-rand	Proposed. Random test locations.
3	KSD	Quadratic-time kernel Stein discrepancy [Liu et al., 2016, Chwialkowski et al., 2016]
4	LKS	Linear-time running average version of KSD.
5	MMD-opt	MMD two-sample test [Gretton et al., 2012]. With optimization.
6	ME-test	<u>Mean Embeddings</u> two-sample test [Jitkrittum et al., 2016]. With optimization.

- Two-sample tests need to draw sample from p .
- Tests with optimization use 20% of the data.
- $\alpha = 0.05$. 200 trials.

Simulation Settings

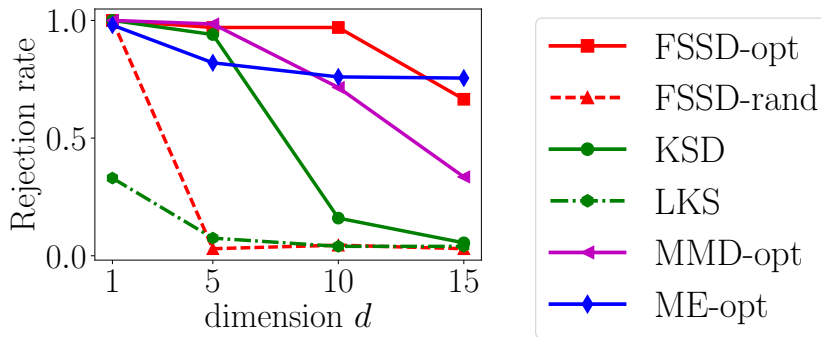
- Gaussian kernel $k(\mathbf{x}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{v}\|_2^2}{2\sigma_k^2}\right)$

	Method	Description
1	FSSD-opt	Proposed. With optimization. $J = 5$.
2	FSSD-rand	Proposed. Random test locations.
3	KSD	Quadratic-time kernel Stein discrepancy [Liu et al., 2016, Chwialkowski et al., 2016]
4	LKS	Linear-time running average version of KSD.
5	MMD-opt	MMD two-sample test [Gretton et al., 2012]. With optimization.
6	ME-test	<u>Mean Embeddings</u> two-sample test [Jitkrittum et al., 2016]. With optimization.

- Two-sample tests need to draw sample from p .
- Tests with optimization use 20% of the data.
- $\alpha = 0.05$. 200 trials.

Gaussian Vs. Laplace

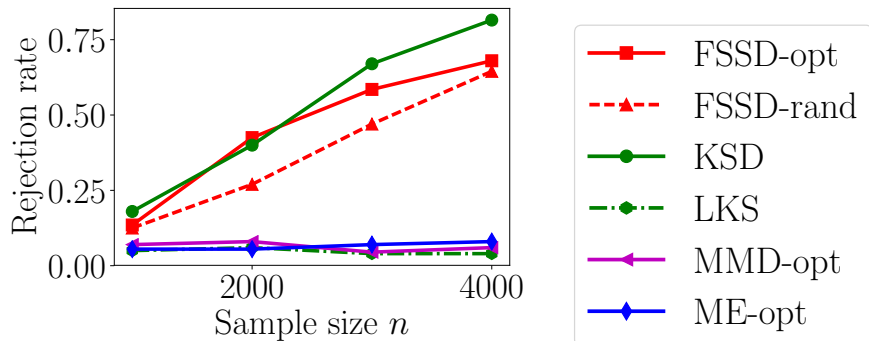
- p = Gaussian. q = Laplace. Same mean and variance. High-order moments differ.
- Sample size $n = 1000$.



- Optimization increases the power.
- Two-sample tests can perform well in this case (p, q clearly differ).

Harder RBM Problem

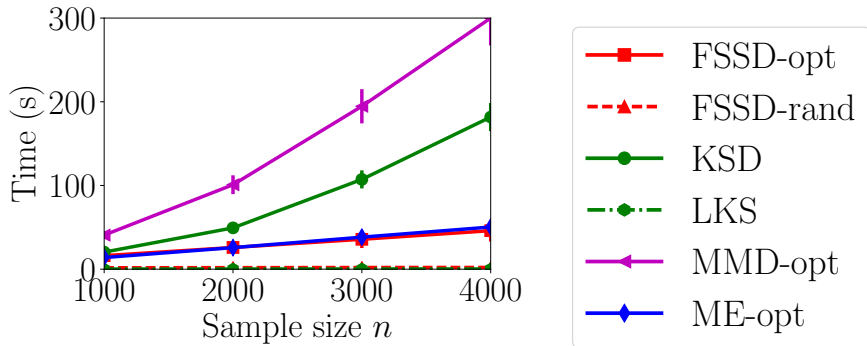
- Perturb only one entry of $\mathbf{B} \in \mathbb{R}^{50 \times 40}$ (in the RBM).
- $B_{1,1} \leftarrow B_{1,1} + \mathcal{N}(0, \sigma_{per}^2 = 0.1^2)$.



- Two-sample tests fail. Samples from p, q look roughly the same.
- FSSD-opt is comparable to KSD at low n . One order of magnitude faster.

Harder RBM Problem

- Perturb only one entry of $\mathbf{B} \in \mathbb{R}^{50 \times 40}$ (in the RBM).
- $B_{1,1} \leftarrow B_{1,1} + \mathcal{N}(0, \sigma_{per}^2 = 0.1^2)$.



- Two-sample tests fail. Samples from p, q look roughly the same.
- FSSD-opt is comparable to KSD at low n . One order of magnitude faster.

References I



Bahadur, R. R. (1960).

Stochastic comparison of tests.

The Annals of Mathematical Statistics, 31(2):276–295.



Chwialkowski, K., Strathmann, H., and Gretton, A. (2016).

A kernel test of goodness of fit.

In *ICML*, pages 2606–2615.





Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012).

A kernel two-sample test.

JMLR, 13:723–773.

References II

-  Jitkrittum, W., Szabó, Z., Chwialkowski, K. P., and Gretton, A. (2016).
Interpretable Distribution Features with Maximum Testing Power.
In *NIPS*, pages 181–189.
-  Liu, Q., Lee, J., and Jordan, M. (2016).
A kernelized Stein discrepancy for goodness-of-fit tests.
In *ICML*, pages 276–284.