

Scalable, Flexible and Active Learning on Distributions

Dougal J. Sutherland

Thesis Committee:

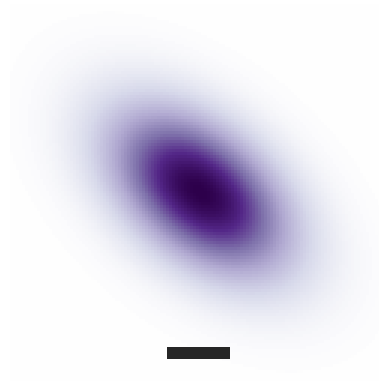
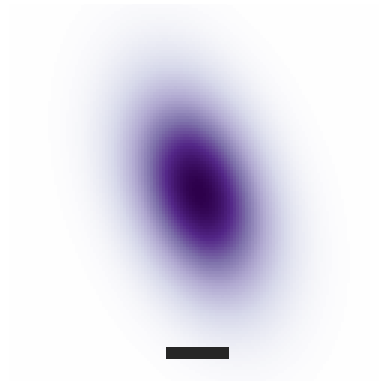
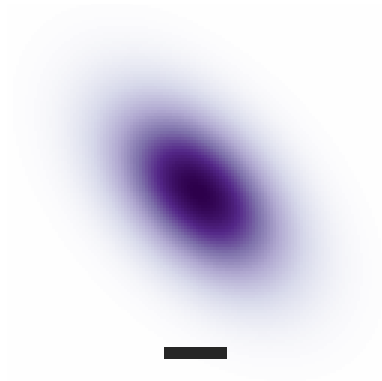
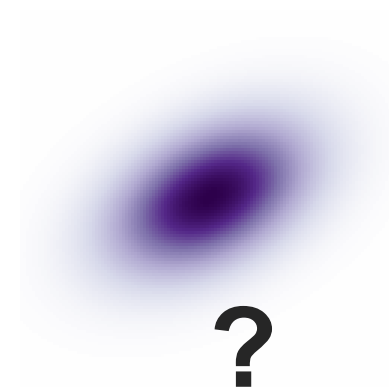
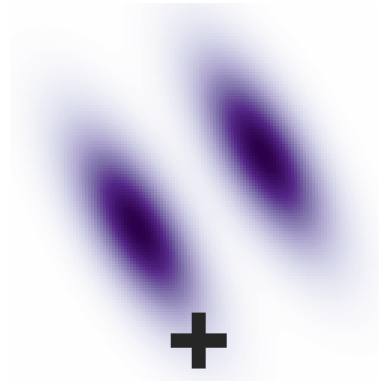
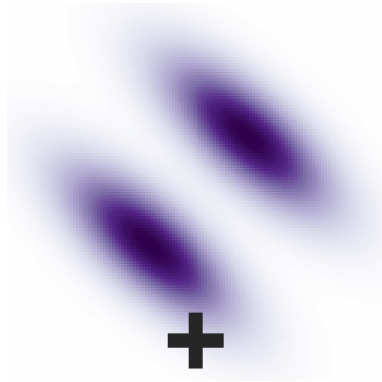
Jeff Schneider (Chair)

Nina Balcan

Barnabás Póczos

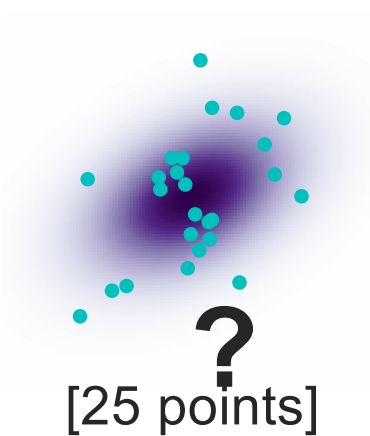
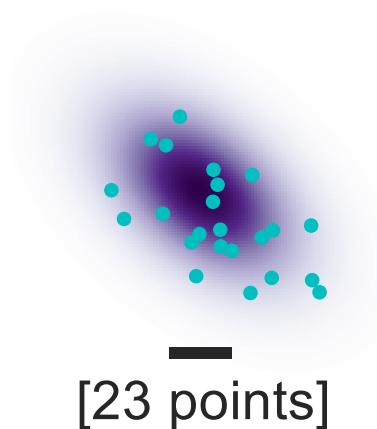
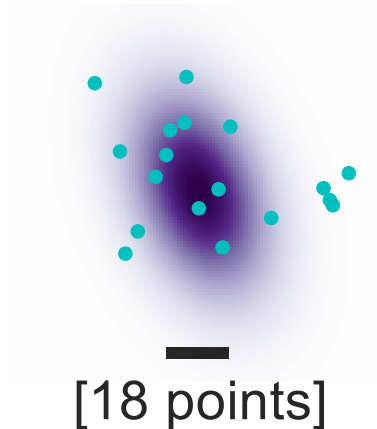
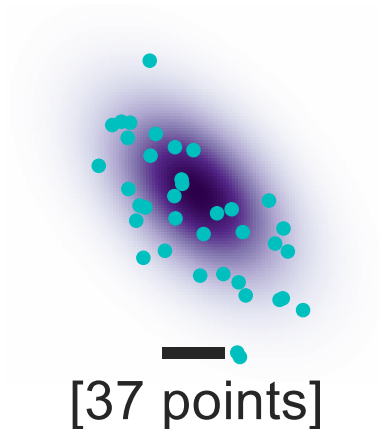
Arthur Gretton (UCL)

Learning on Distributions



We want to learn a distribution classification function....

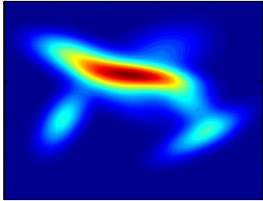
Learning on Distributions



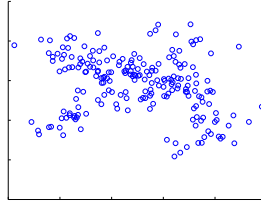
...based on sample sets.

Learning on Distributions

distribution



observed sample

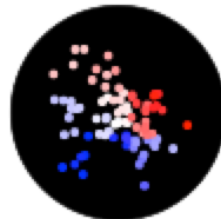
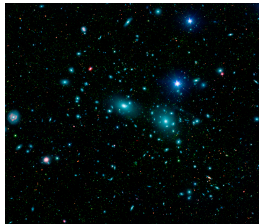


label

9 components



“seaside city”

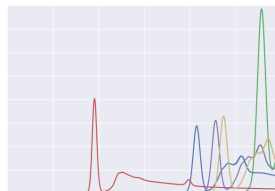


Mass $7 \times 10^{14} M_{\odot}$ and more...

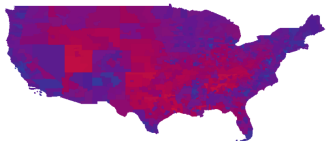
Ntampaka et al. (ApJ 2015, 2016)



Jin et al. (NSS 2016)



no Cs137 present



	AGEP	SEX	...	RACSOR	RACWHT
0	75	1	...	0	1
1	25	0	...	0	1
...

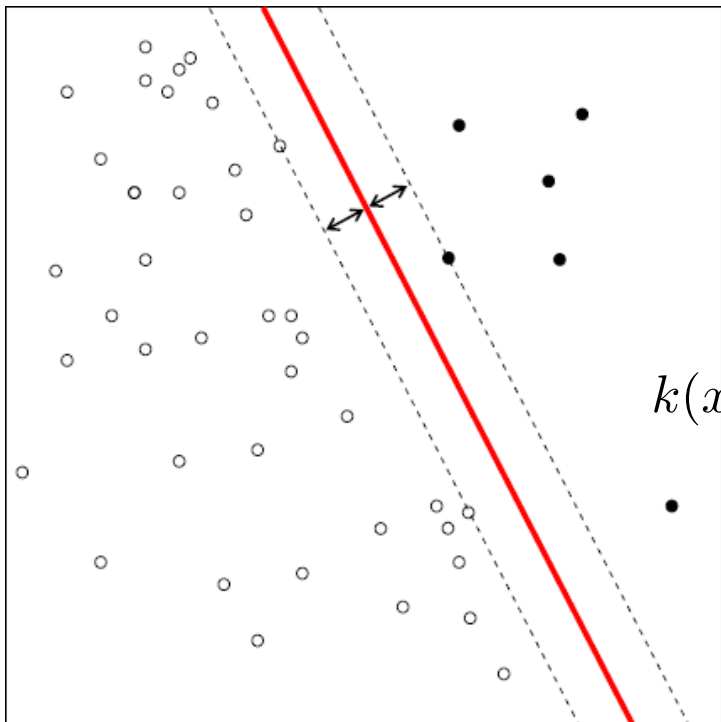
county voted
54% for Obama

Flaxman et al. (KDD 2015)

Contributions

- **Learning on distributions** with nonparametric kernels
- **Scalable** approximate kernel embeddings
 - Random Fourier features analysis
 - New embeddings for distribution kernels
- **Flexible** distribution kernels
 - Deep mean maps in computer vision
 - MMD kernel learning for testing
- **Active** pointillistic pattern search

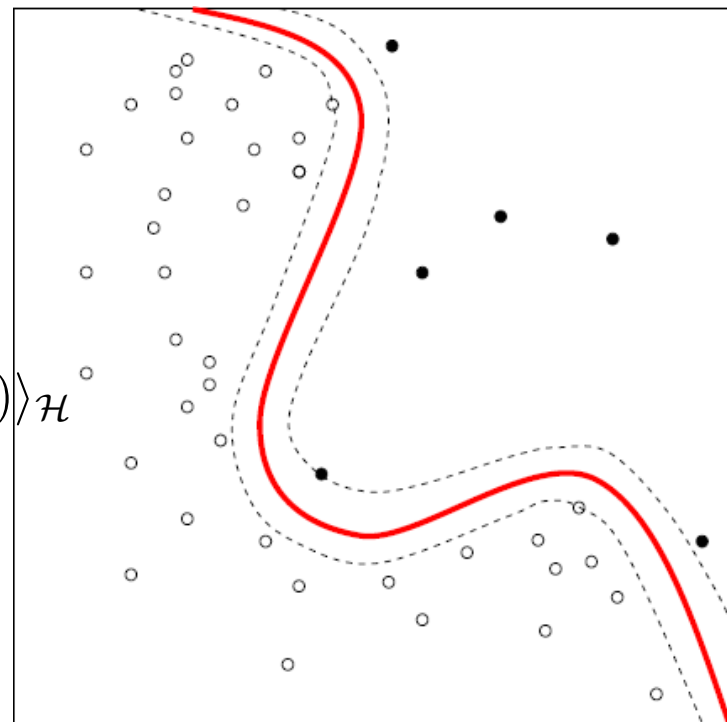
Kernel Methods



$$f(x) = w^T x + b$$

Linear models...

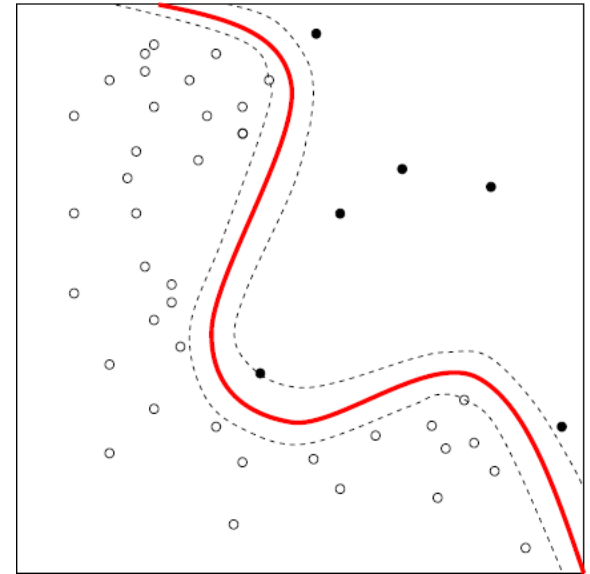
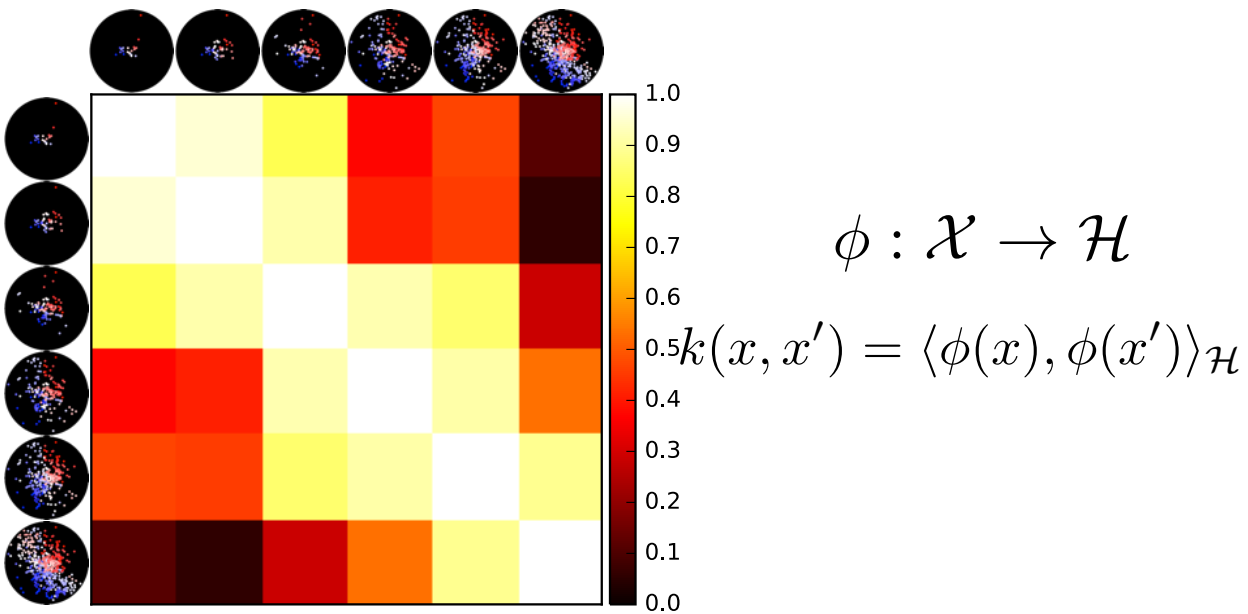
$$\phi : \mathbb{R}^d \rightarrow \mathcal{H}$$
$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$$



$$f(x) = \langle w, \phi(x) \rangle_{\mathcal{H}} + b$$
$$= \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b$$

...in Hilbert space.

Kernel Methods



Can use a kernel
on any domain.

$$\begin{aligned} f(x) &= \langle w, \phi(x) \rangle_{\mathcal{H}} + b \\ &= \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \end{aligned}$$

Linear models...in Hilbert space.

Kernels on Distributions

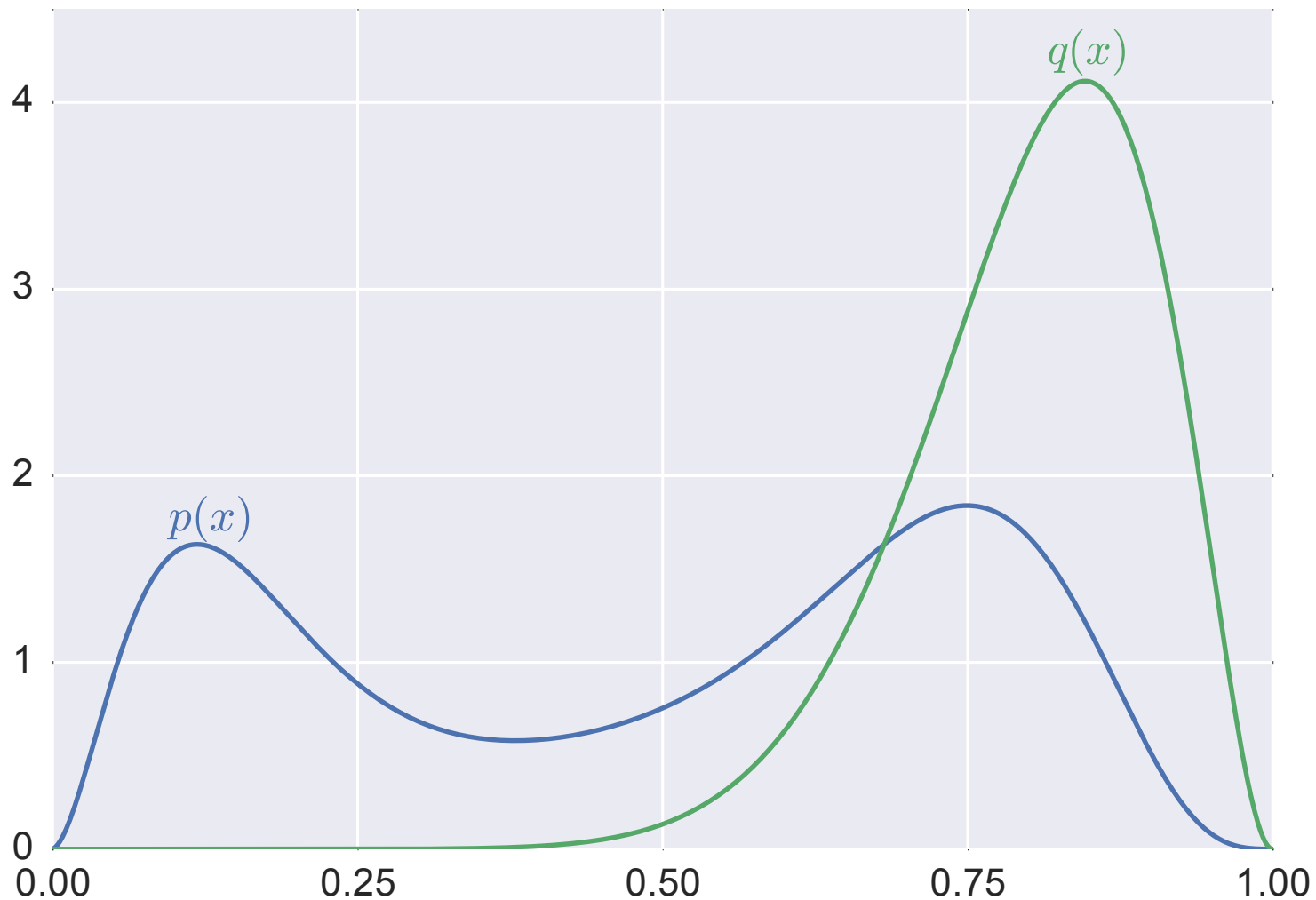
We'll use a kernel on distributions based on a distance ρ :

$$K(\mathbb{P}, \mathbb{Q}) = \exp \left(-\frac{1}{2\sigma^2} \rho^2(\mathbb{P}, \mathbb{Q}) \right)$$

The popular Gaussian RBF kernel has this form, with ρ the Euclidean distance between vectors.

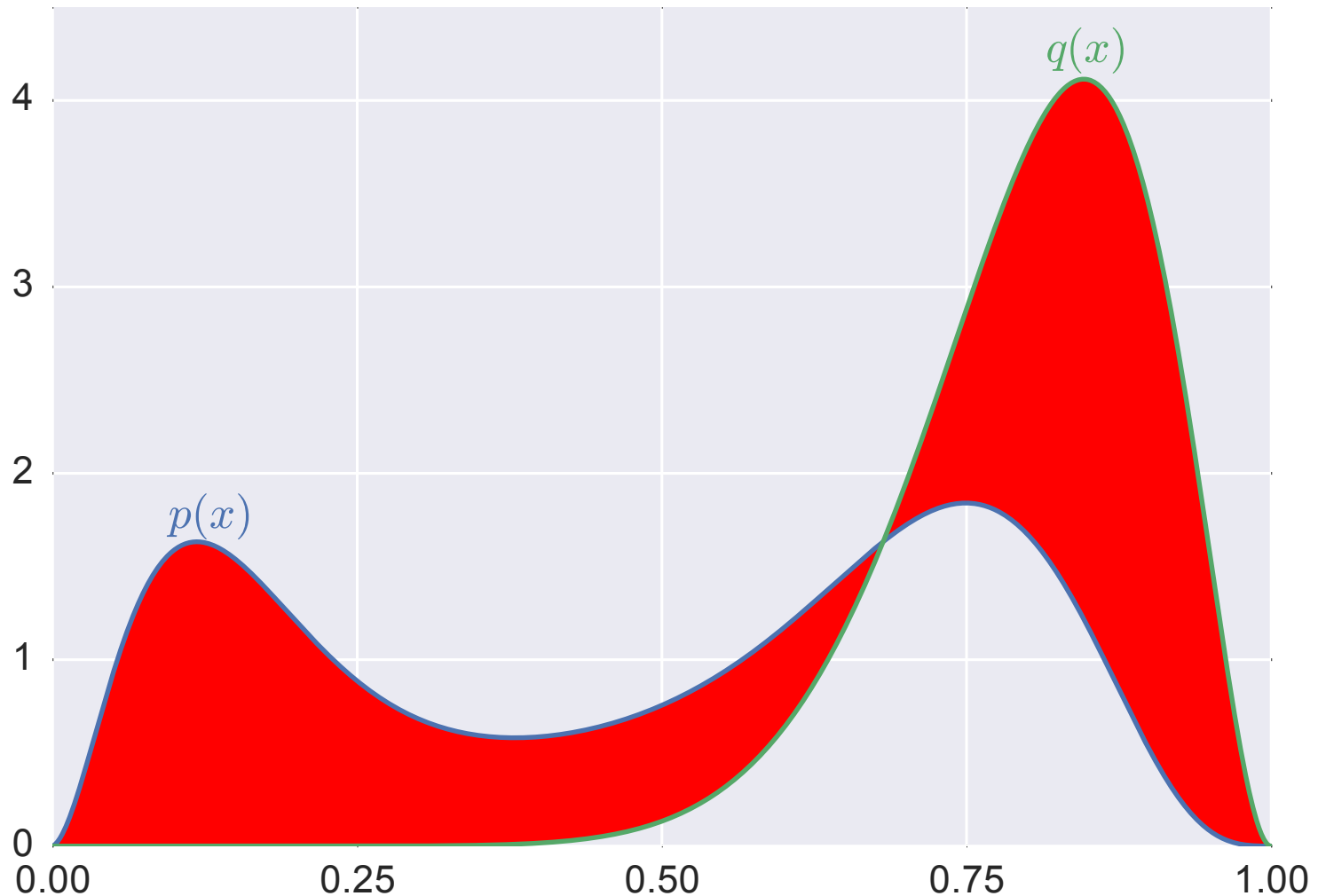
A valid kernel as long as ρ is *Hilbertian*.

Distances on Distributions



Distances on Distributions

$$\text{TV}(\mathbb{P}, \mathbb{Q}) = \int \frac{1}{2} |p(x) - q(x)| \, dx$$



Distances on Distributions

Total Variation

$$\text{TV}(p, q) = \int \frac{1}{2} |p(x) - q(x)| \, dx$$

L_2

$$L_2^2(p, q) = \int (p(x) - q(x))^2 \, dx$$

Hellinger

$$H^2(p, q) = \int \frac{1}{2} \left(\sqrt{p(x)} - \sqrt{q(x)} \right)^2 \, dx$$

Kullback-Leibler

$$\text{KL}(p \| q) = \int p(x) \log \frac{p(x)}{q(x)} \, dx$$

Rényi- α

$$R_\alpha(p \| q) = \frac{1}{\alpha - 1} \int p(x)^\alpha q(x)^{1-\alpha} \, dx$$

Jensen-Shannon

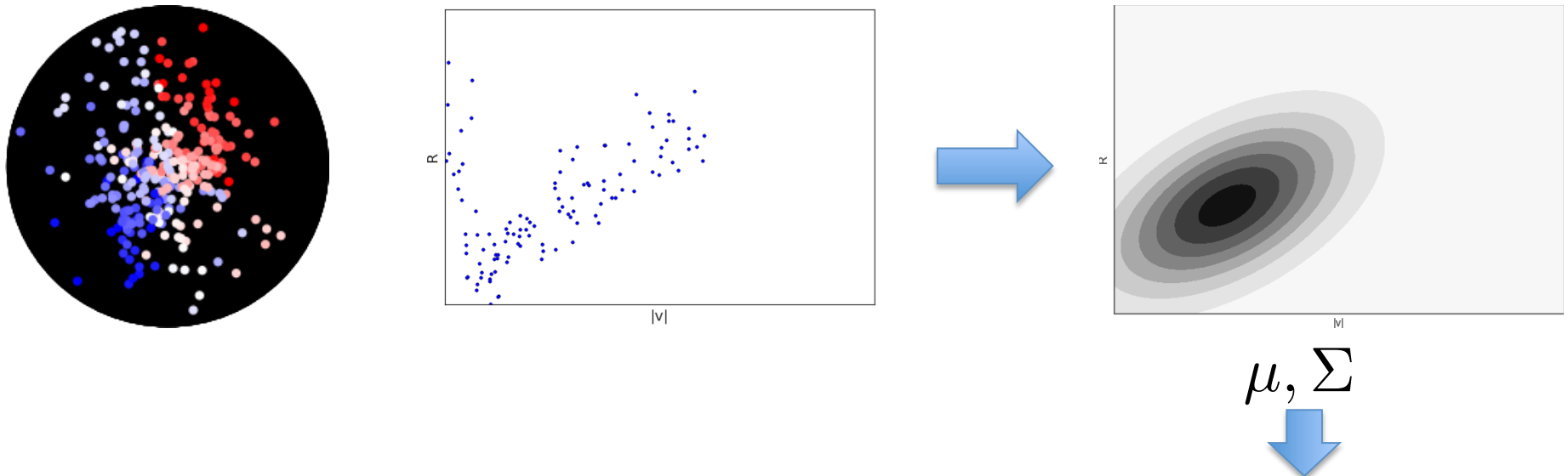
$$\text{JS}(p, q) = \frac{1}{2} \text{KL} \left(p \left\| \frac{p+q}{2} \right. \right) + \frac{1}{2} \text{KL} \left(q \left\| \frac{p+q}{2} \right. \right)$$

Maximum Mean Discrepancy

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) = \sup_{f \in \mathcal{H}_k} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y)$$

Estimators of Distributional Distances

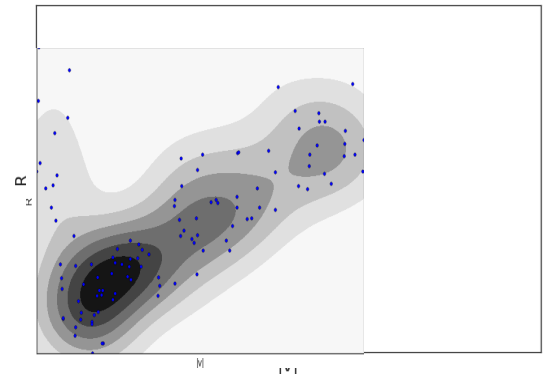
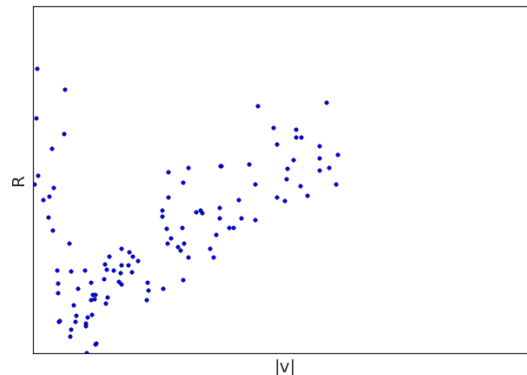
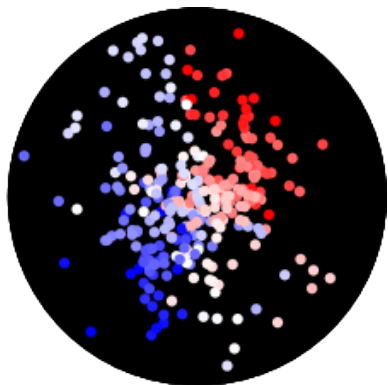
- Fit a parametric model and compute distances.
 - Some distances have closed form for some models.
 - Model introduces approximation error



$$L_2(\mathcal{N}(\mu, \Sigma), \mathcal{N}(\mu', \Sigma')) = \frac{1}{|4\pi\Sigma|^{\frac{1}{2}}} + \frac{1}{|4\pi\Sigma'|^{\frac{1}{2}}} - 2 \frac{\exp\left(-\frac{1}{2}(\mu - \mu')^T(\Sigma + \Sigma')^{-1}(\mu - \mu')\right)}{|2\pi(\Sigma + \Sigma')|^{\frac{1}{2}}}$$

Estimators of Distributional Distances

- Fit a nonparametric model and compute distances.
 - Histograms
 - Kernel density estimation
 - k -nearest neighbor density estimation
 - Basis function projections
 - Empirical distribution (for MMD)



Algorithm: Learning on Distributions

Given sample sets $X_i \sim \mathbb{P}_i$, a distance ρ , and a b.w. σ :

1. Estimate $\hat{\rho}(X_i, X_j)$ for all i, j nonparametrically.

2. Assemble into a kernel matrix

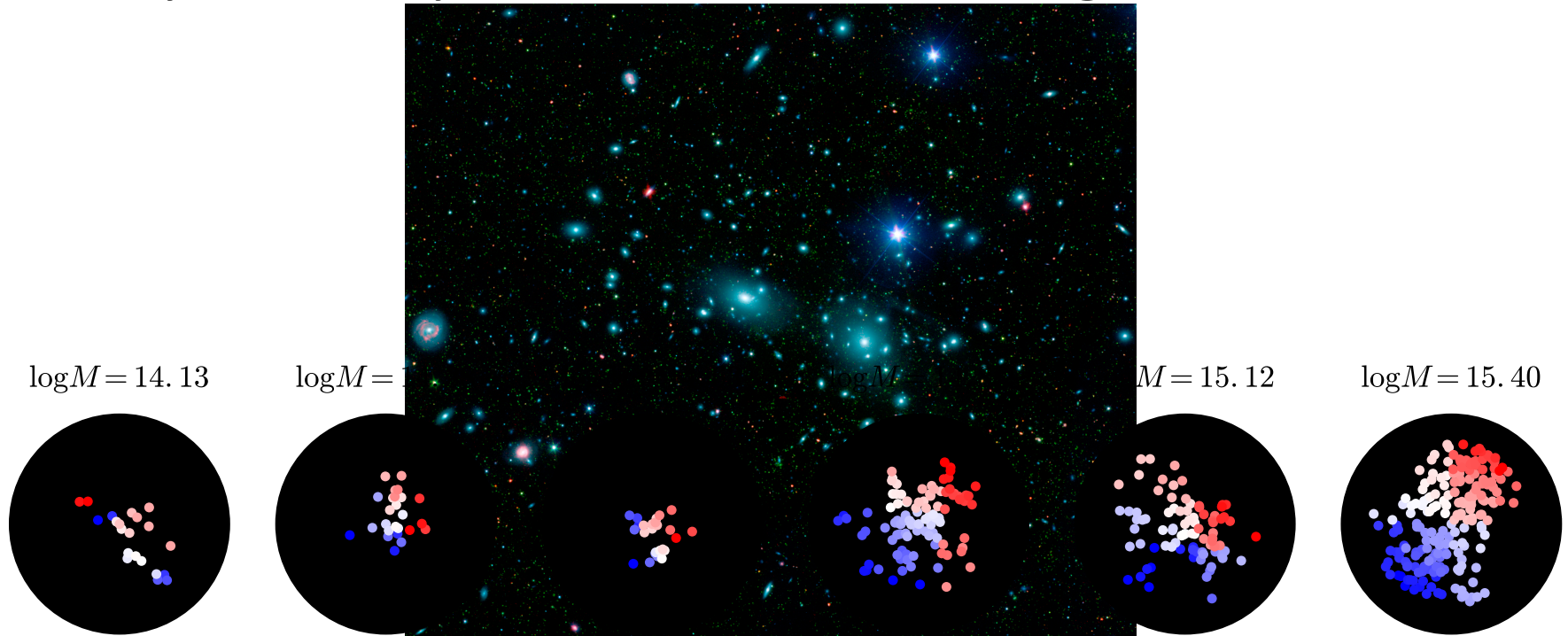
$$K_{ij} = \exp \left(-\frac{1}{2\sigma^2} \hat{\rho}(X_i, X_j) \right).$$

3. Run an SVM / GP / ridge regression / ... with K .

Application: Galaxy Cluster Mass

Galaxy clusters are fundamental in the study the universe. Their properties can tell us a lot about cosmology.

But they're mostly dark matter; measuring their mass is hard.



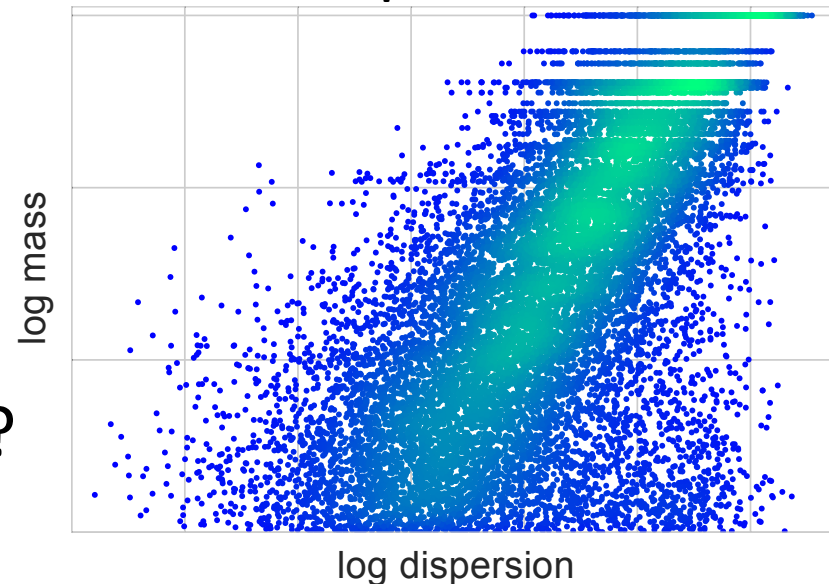
[Coma cluster](#), NASA

Application: Galaxy Cluster Mass

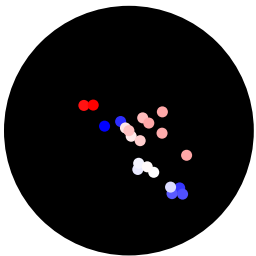
Fritz Zwicky (1933): Under reasonable assumptions, the velocity dispersion has power-law relationship to total mass.

Not so great...

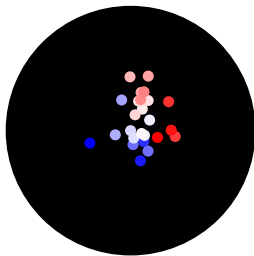
- Assumptions violated
- Which galaxies are in cluster?



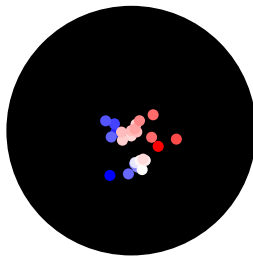
$\log M = 14.13$



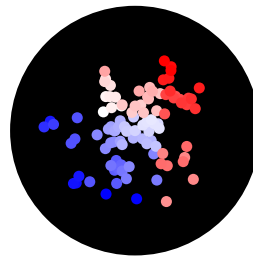
$\log M = 14.38$



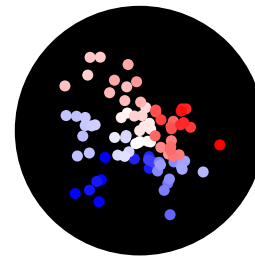
$\log M = 14.63$



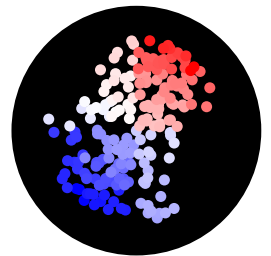
$\log M = 14.88$



$\log M = 15.12$

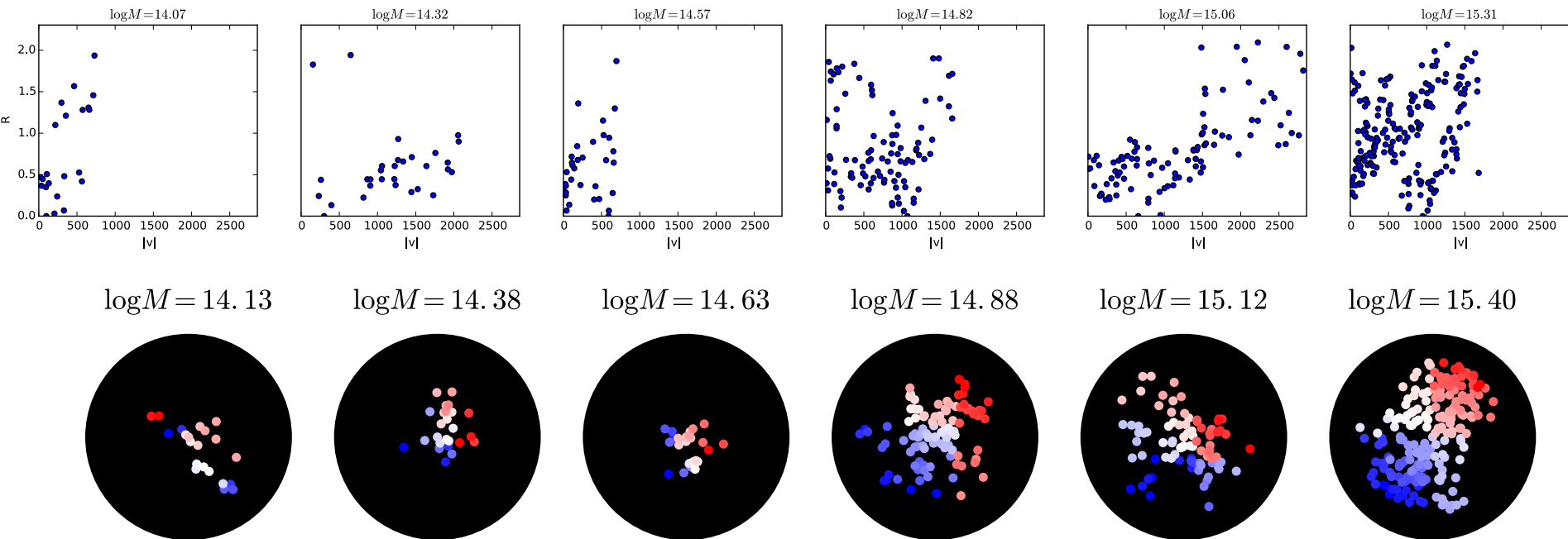


$\log M = 15.40$

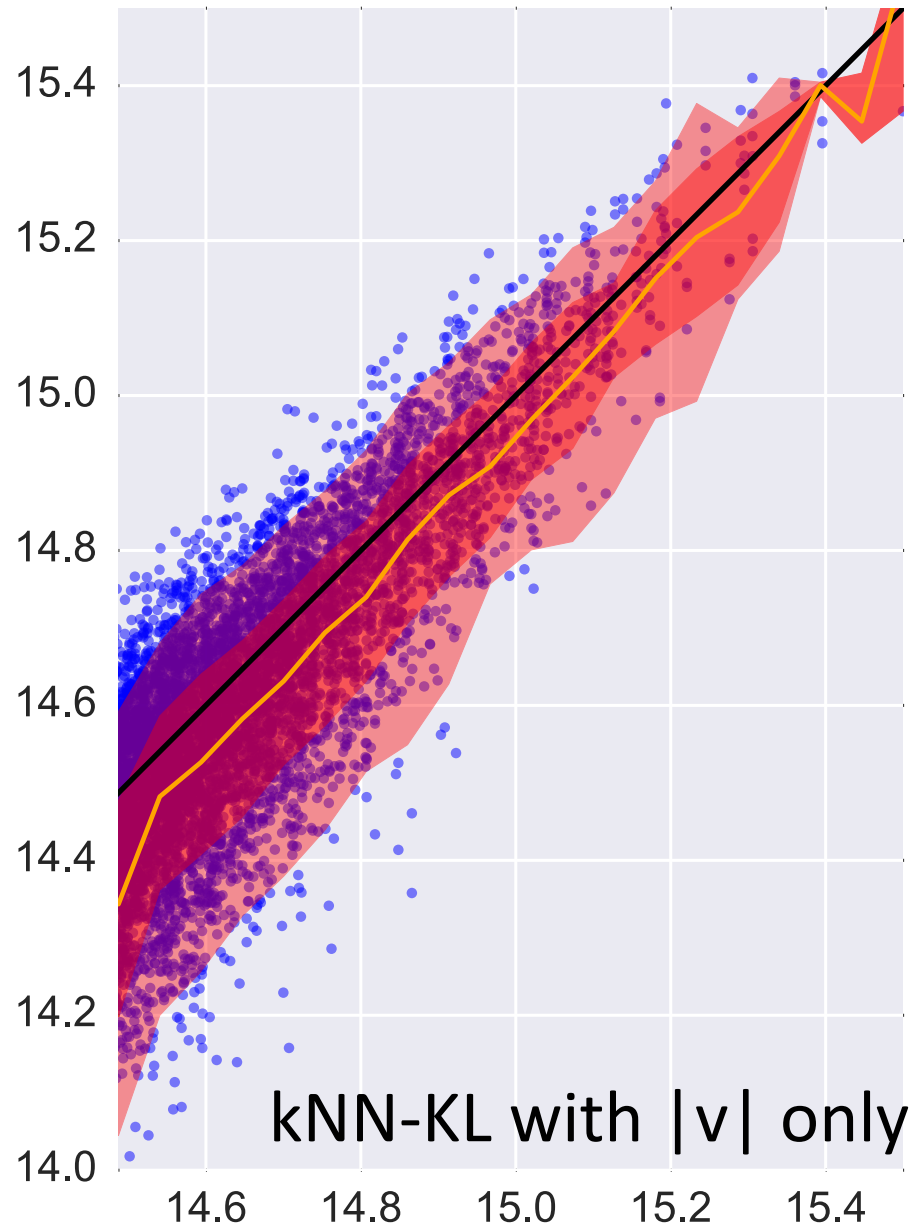
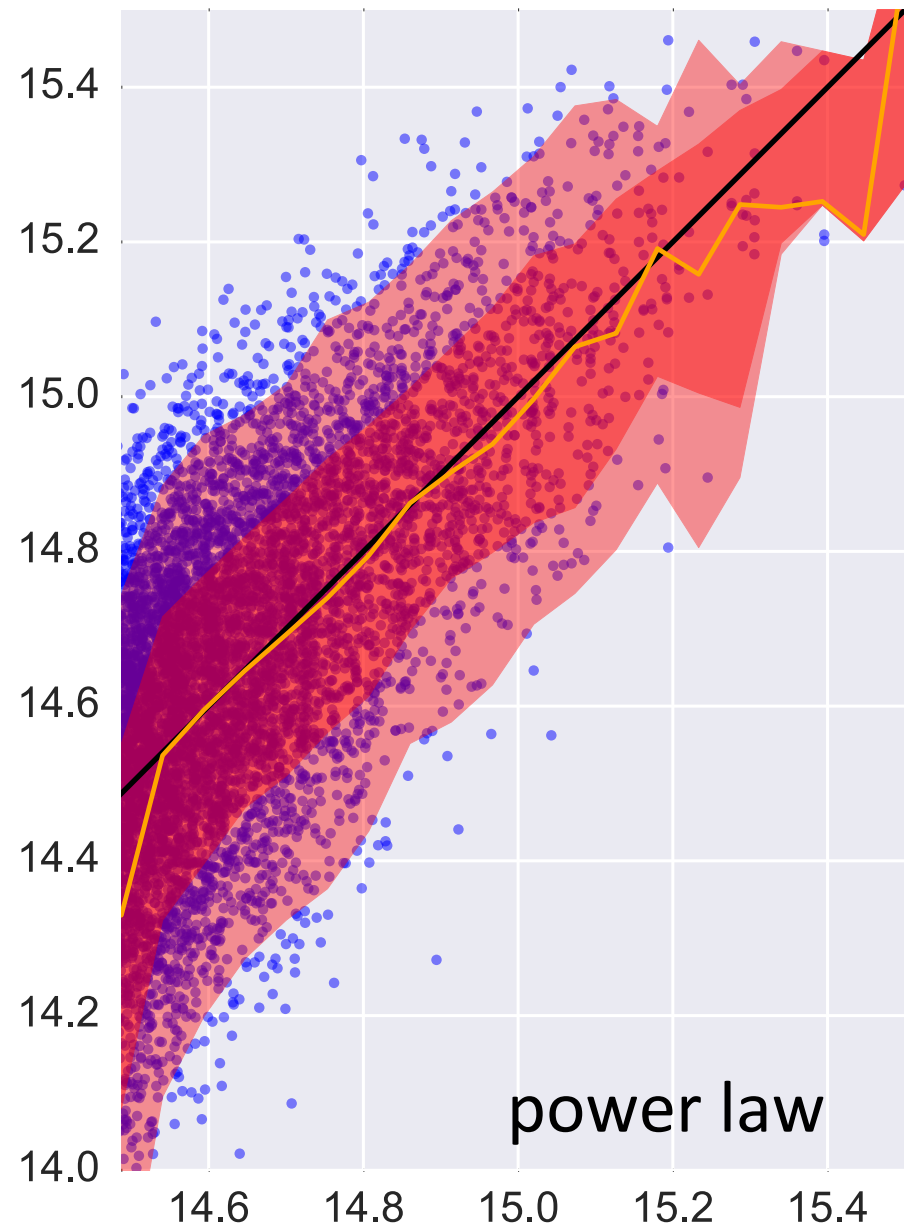


Application: Galaxy Cluster Mass

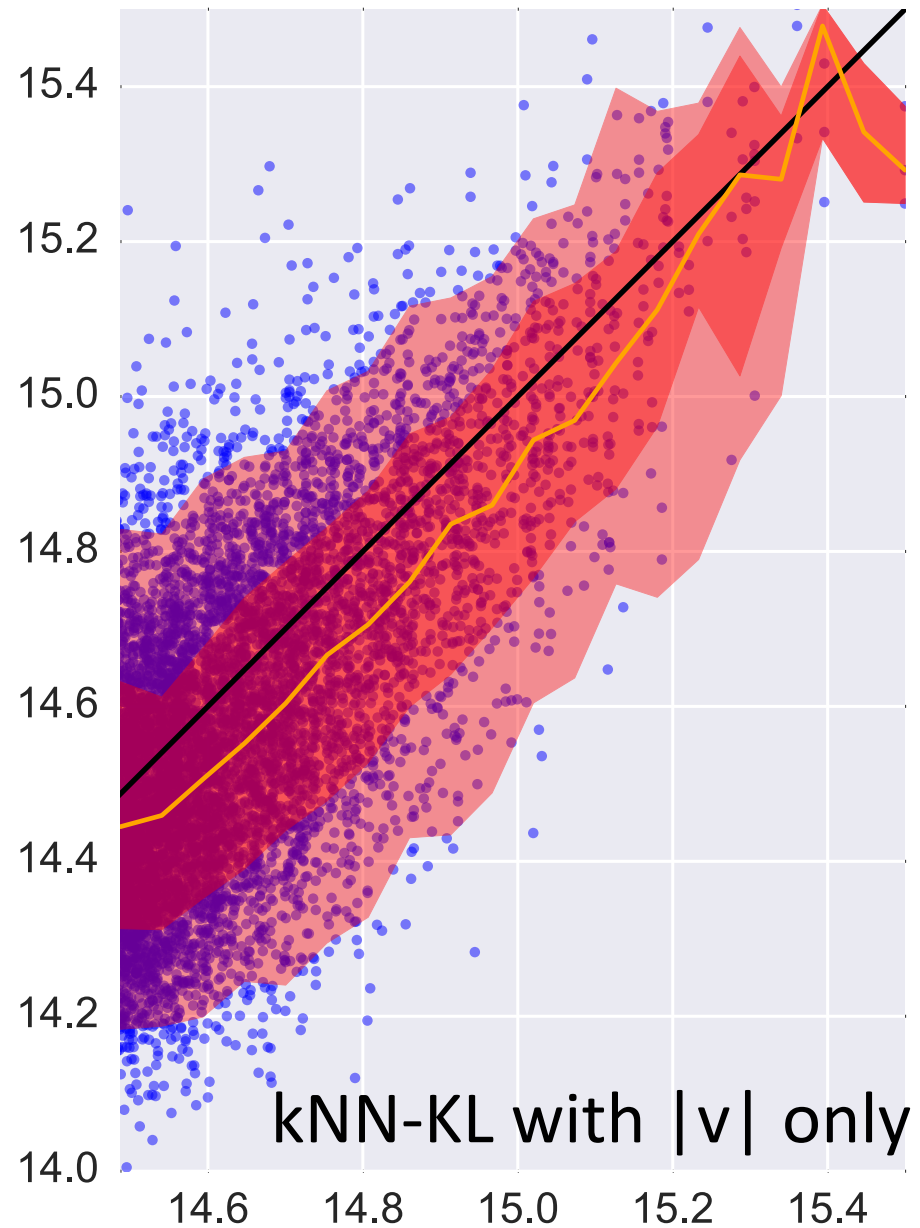
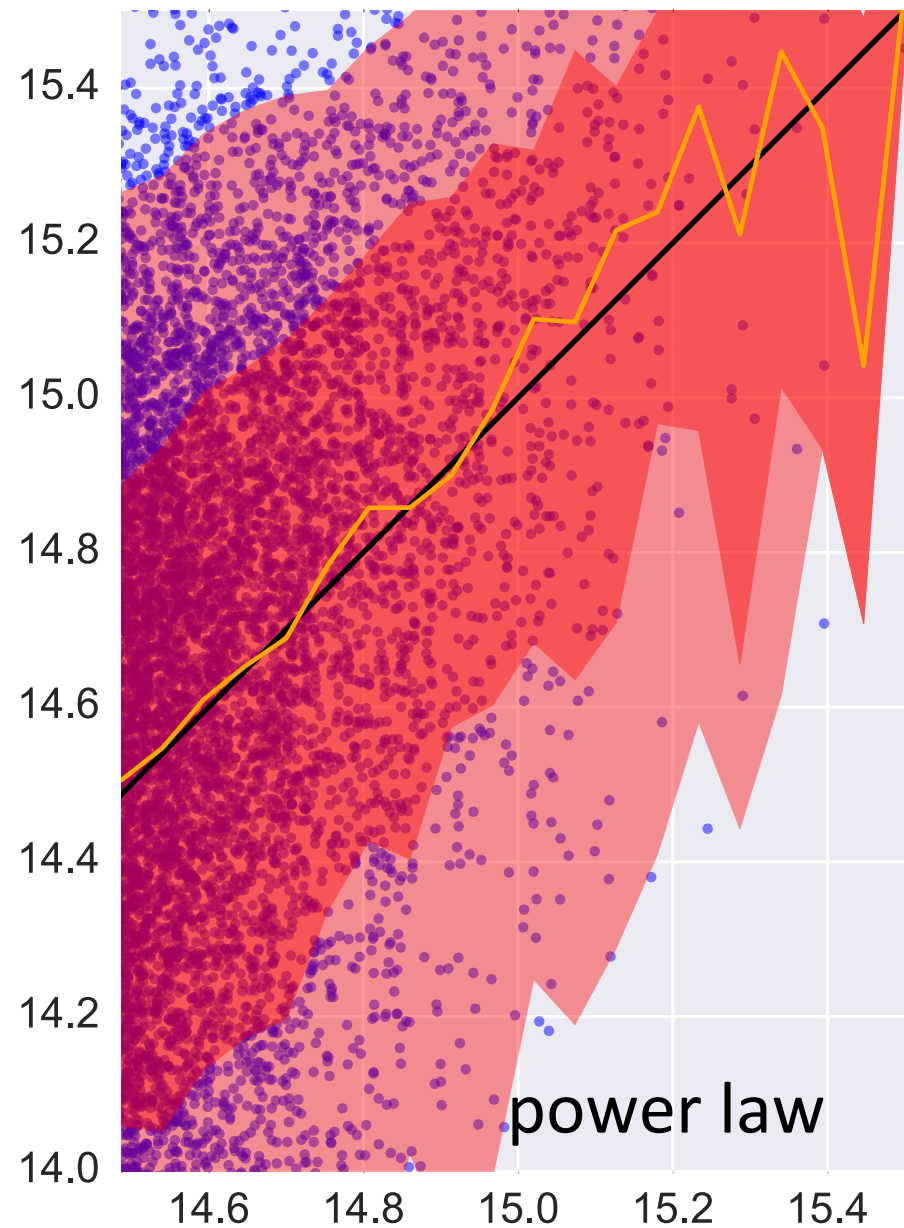
Alternative approach: consider each cluster as a *distribution* of galaxy features, and regress from these distributions to total mass.



Cluster Mass: Known Membership



Cluster Mass: With Interlopers



Contributions

- **Learning on distributions** with nonparametric kernels
- **Scalable** approximate kernel embeddings
 - Random Fourier features analysis
 - New embeddings for distribution kernels
- **Flexible** distribution kernels
 - Deep mean maps in computer vision
 - MMD kernel learning for testing
- **Active** pointillistic pattern search

Scalability

These methods need an $n \times n$ Gram matrix:

$$\mathbf{K} = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_n) \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(x_n, x_1) & K(x_n, x_2) & \dots & K(x_n, x_n) \end{bmatrix}$$

The matrix is n^2 ; operations can be as slow as $O(n^3)$.

Linear-kernel models usually scale like $O(n)$.

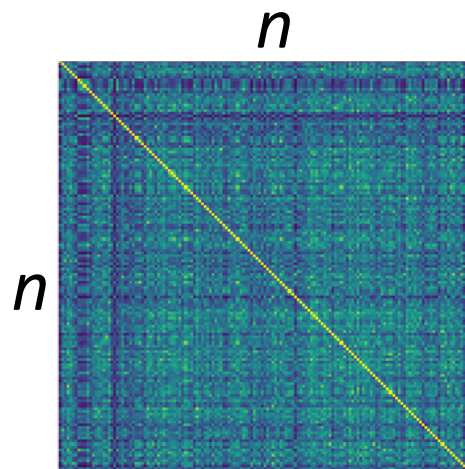
Approximate Embeddings

Traditional kernel methods:

$$\{x_i\} \quad \{\varphi(x_i)\} \subset \mathcal{H}$$

$$\langle \varphi(x_i), \varphi(x_j) \rangle_{\mathcal{H}} = k(x_i, x_j)$$

$$f(x) = \langle w, \varphi(x) \rangle + b = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b$$



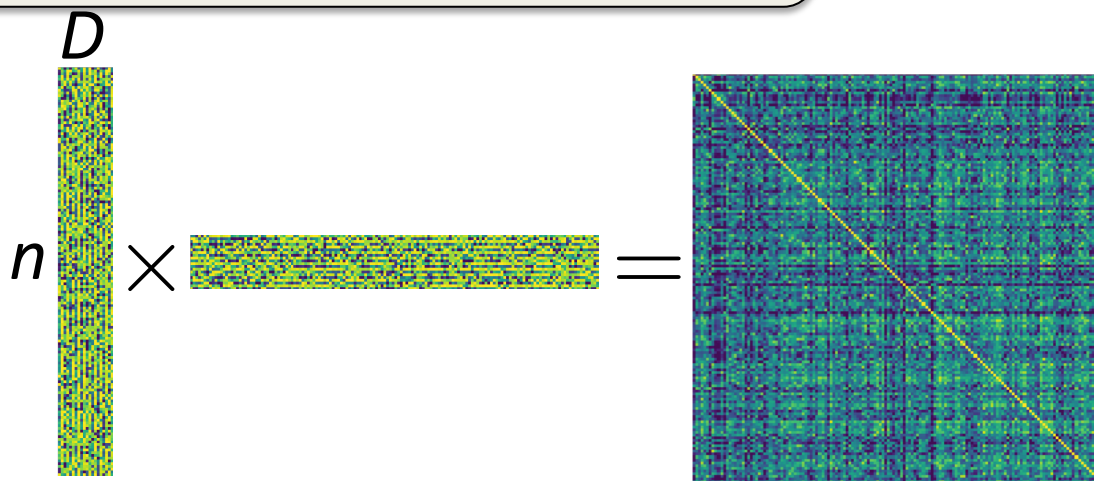
Approx

$$K(\text{img}_1, \text{img}_2) \approx z(\text{point_cloud}_1)^T z(\text{point_cloud}_2) \approx$$

$$\{x_i\} \quad \{z(x_i)\} \subset \mathbb{R}^D$$

$$z(x_i)^T z(x_j)$$

$$f(x) = w^T z(x) + b$$



Random Fourier Features

Rahimi and Recht (2007) developed random Fourier features, a.k.a. “random kitchen sinks”:

$$\begin{aligned} k(x, y) &= \underline{k}(\overbrace{\Delta}^{x-y}) = \exp\left(-\frac{1}{2\sigma^2}\|\Delta\|^2\right) \\ &= \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}} \quad (\dim \mathcal{H} = \infty) \\ &\approx z(x)^{\top} z(y) \quad (z : \mathbb{R}^d \rightarrow \mathbb{R}^D) \end{aligned}$$

Random Fourier Features

Rahimi and Recht (2007) developed random Fourier features, a.k.a. “random kitchen sinks”:

$$k(x, y) = \underline{k}(\overbrace{\Delta}^{x-y}) = \exp\left(-\frac{1}{2\sigma^2}\|\Delta\|^2\right) \approx z(x)^\top z(y)$$

$$\Omega(\omega) := \int \underline{k}(\Delta) \exp(-\mathrm{i}\omega^\top \Delta) \, \mathrm{d}\Delta \propto \exp\left(-\frac{\sigma^2}{2}\|\omega\|^2\right)$$

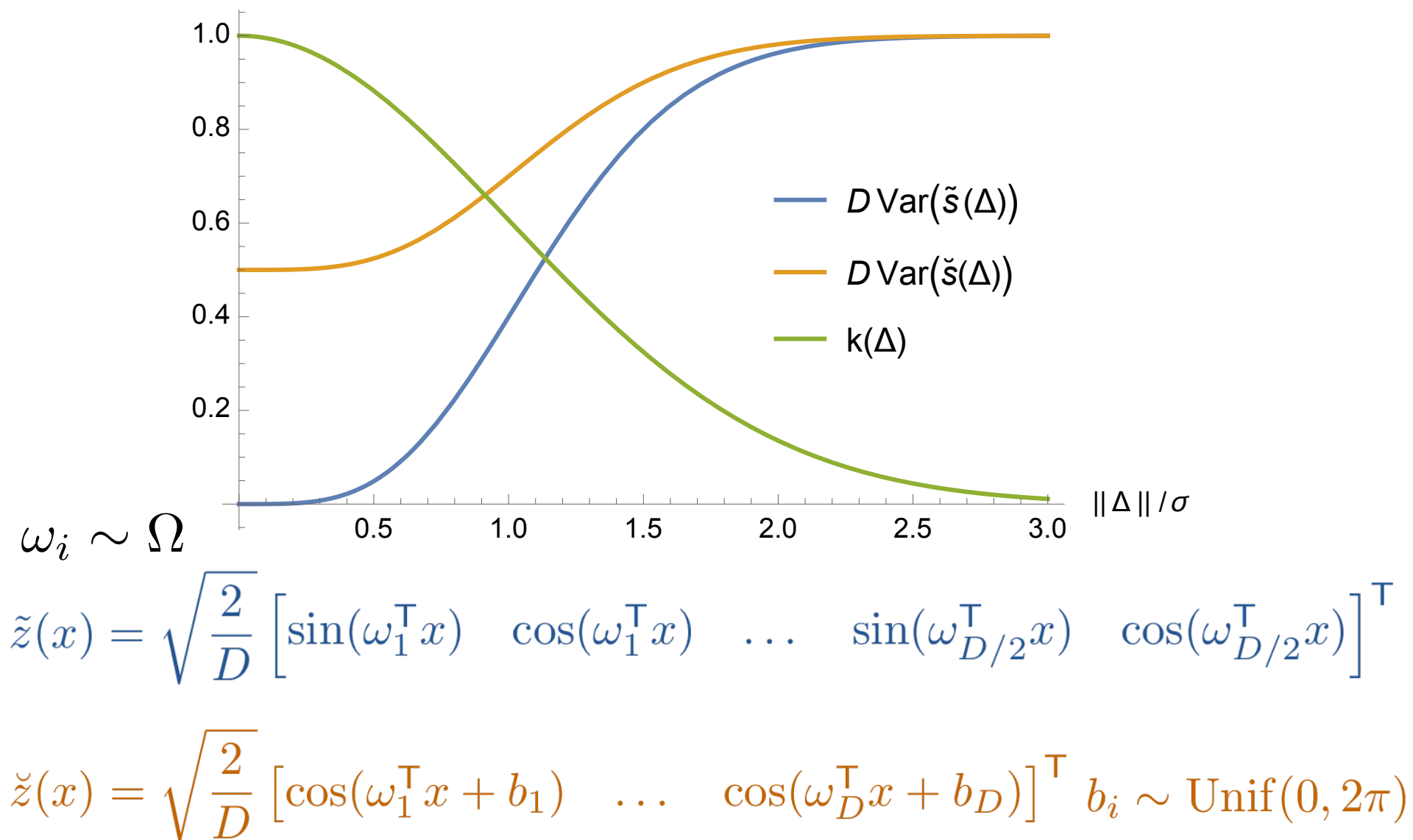
Bochner’s theorem: Ω is a (scaled) probability distribution.

$$\omega_i \sim \Omega$$

$$\tilde{z}(x) = \sqrt{\frac{2}{D}} \begin{bmatrix} \sin(\omega_1^\top x) & \cos(\omega_1^\top x) & \dots & \sin(\omega_{D/2}^\top x) & \cos(\omega_{D/2}^\top x) \end{bmatrix}^\top$$

$$\check{z}(x) = \sqrt{\frac{2}{D}} \begin{bmatrix} \cos(\omega_1^\top x + b_1) & \dots & \cos(\omega_D^\top x + b_D) \end{bmatrix}^\top \quad b_i \sim \text{Unif}(0, 2\pi)$$

Random Fourier Features



The L_∞ error is also tighter, in bounds and empirically. ²⁵

Contributions

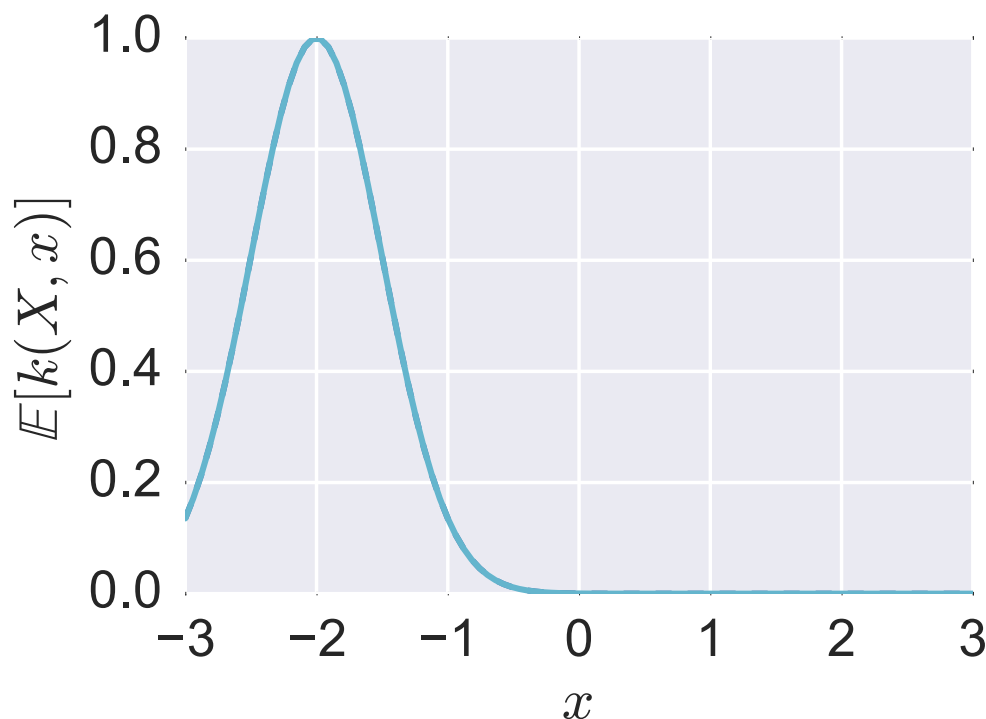
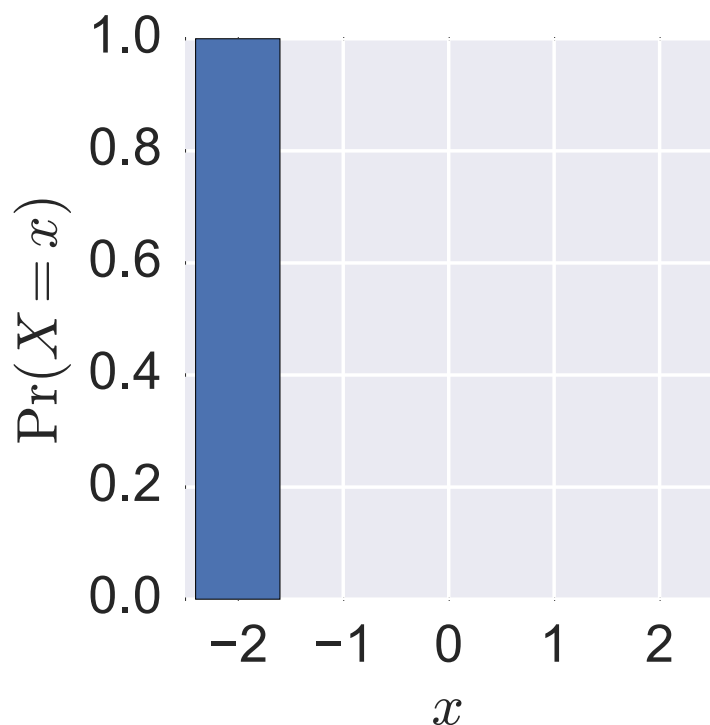
- **Learning on distributions** with nonparametric kernels
- **Scalable** approximate kernel embeddings
 - Random Fourier features analysis
 - New embeddings for distribution kernels
- **Flexible** distribution kernels
 - Deep mean maps in computer vision
 - MMD kernel learning for testing
- **Active** pointillistic pattern search

Mean embeddings of distributions

The *mean embedding* of a distribution in an RKHS:

$$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)]$$

Remember $\langle \varphi(x), \varphi(y) \rangle = k(x, y)$,
so we can think of $\varphi(x)$ as $k(x, \cdot)$.

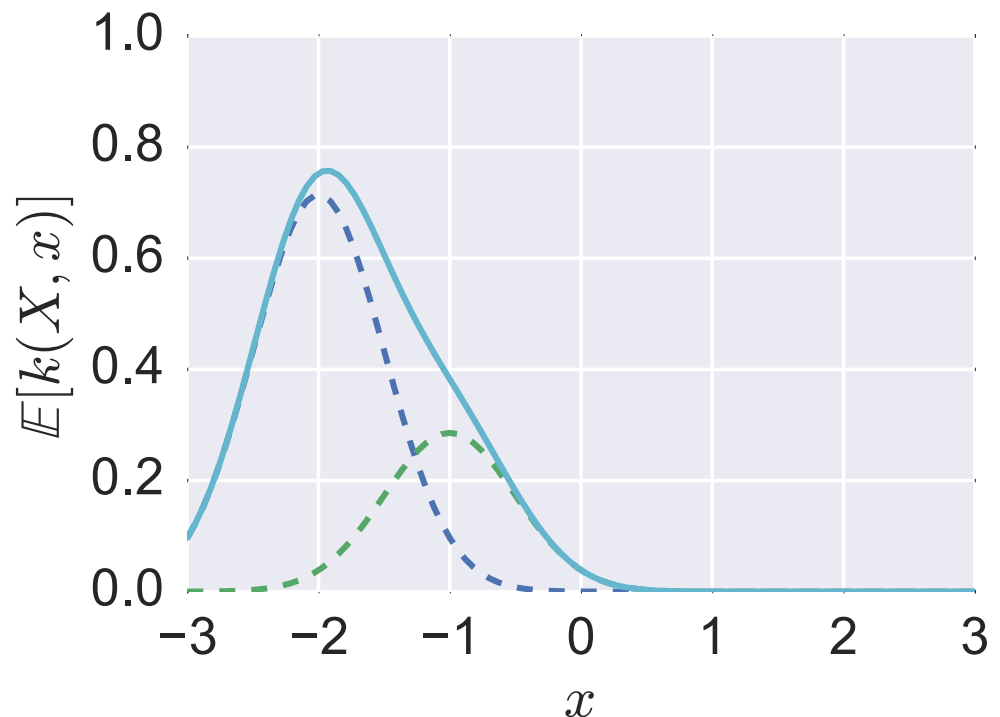
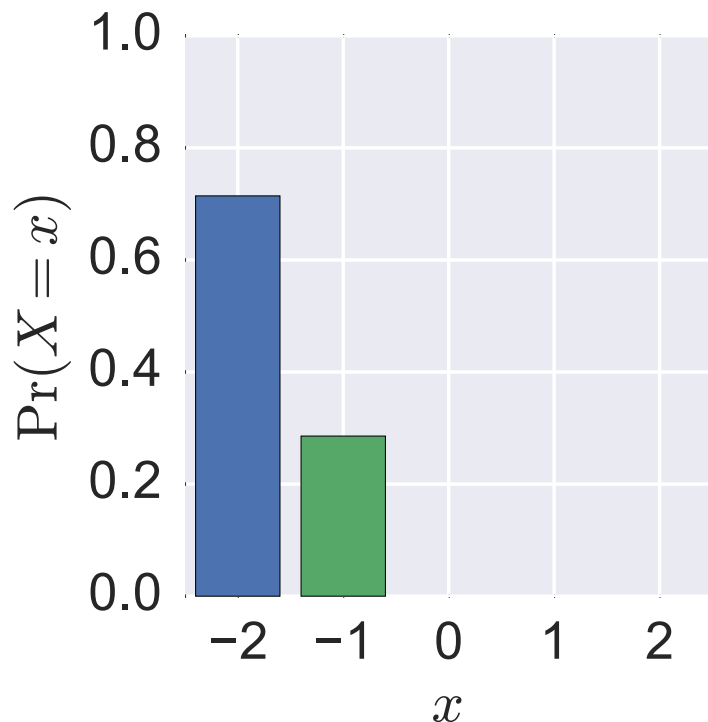


Mean embeddings of distributions

The *mean embedding* of a distribution in an RKHS:

$$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)]$$

Remember $\langle \varphi(x), \varphi(y) \rangle = k(x, y)$,
so we can think of $\varphi(x)$ as $k(x, \cdot)$.

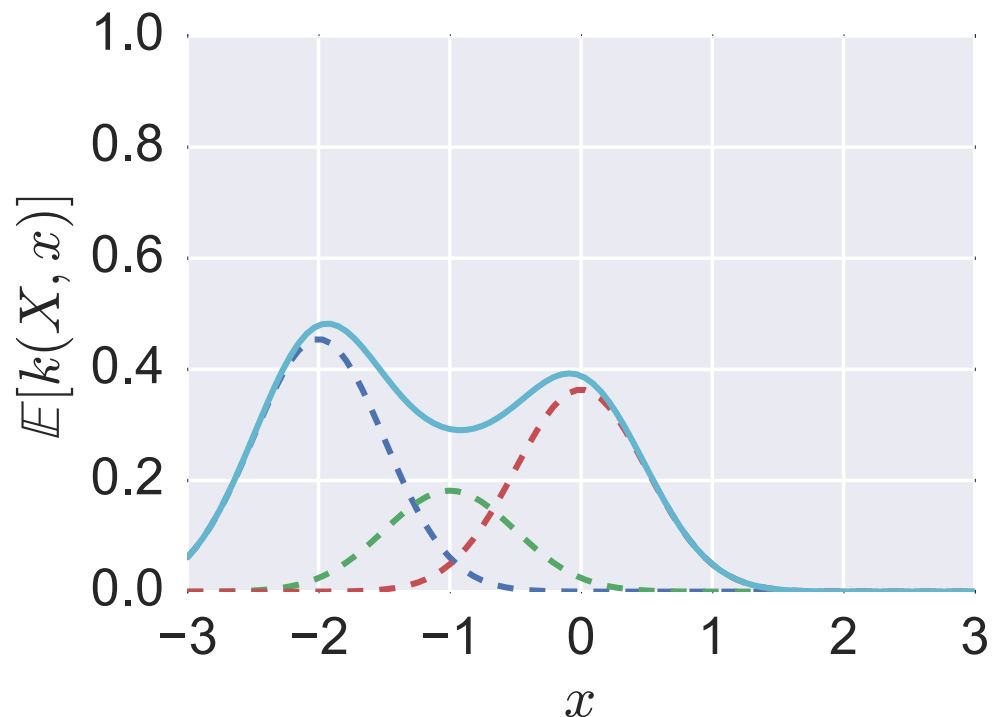
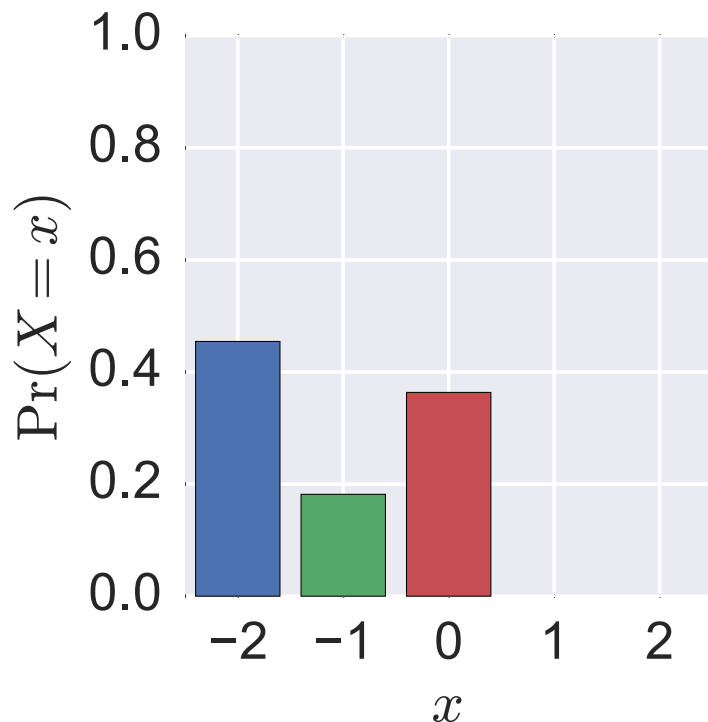


Mean embeddings of distributions

The *mean embedding* of a distribution in an RKHS:

$$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)]$$

Remember $\langle \varphi(x), \varphi(y) \rangle = k(x, y)$,
so we can think of $\varphi(x)$ as $k(x, \cdot)$.

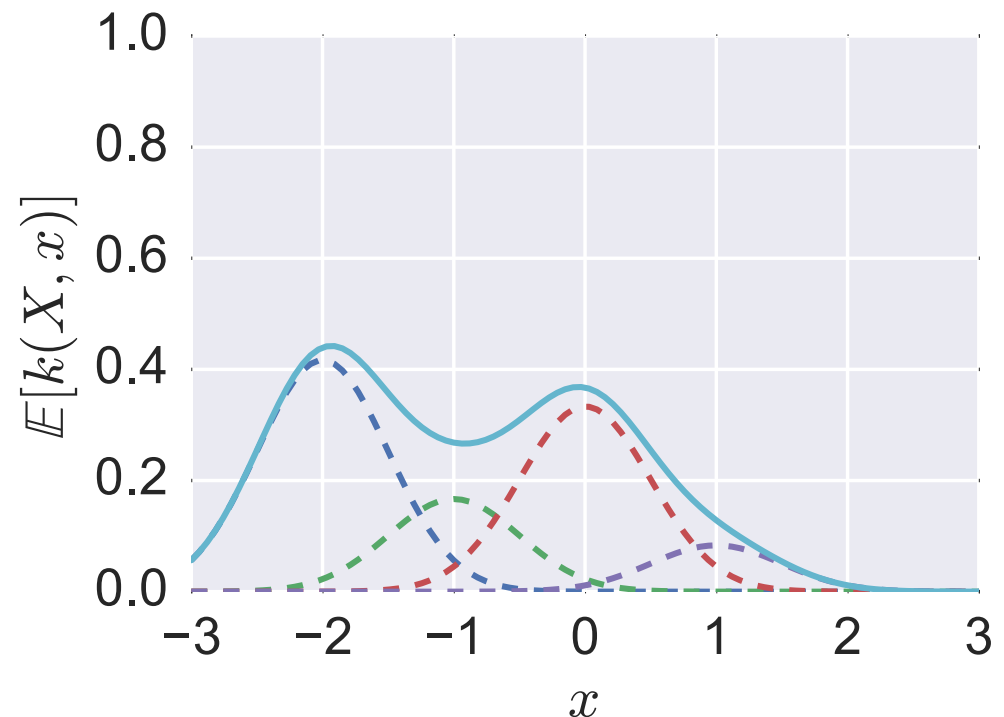
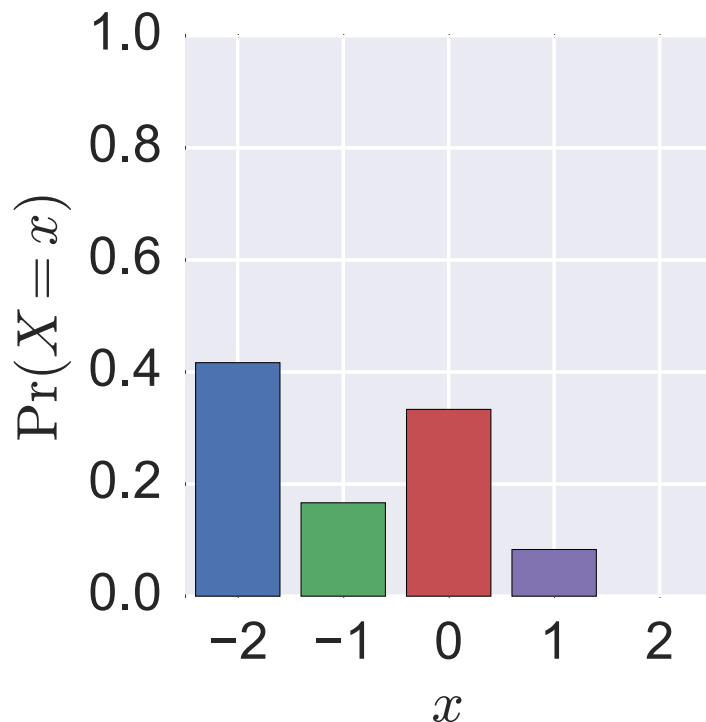


Mean embeddings of distributions

The *mean embedding* of a distribution in an RKHS:

$$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)]$$

Remember $\langle \varphi(x), \varphi(y) \rangle = k(x, y)$,
so we can think of $\varphi(x)$ as $k(x, \cdot)$.

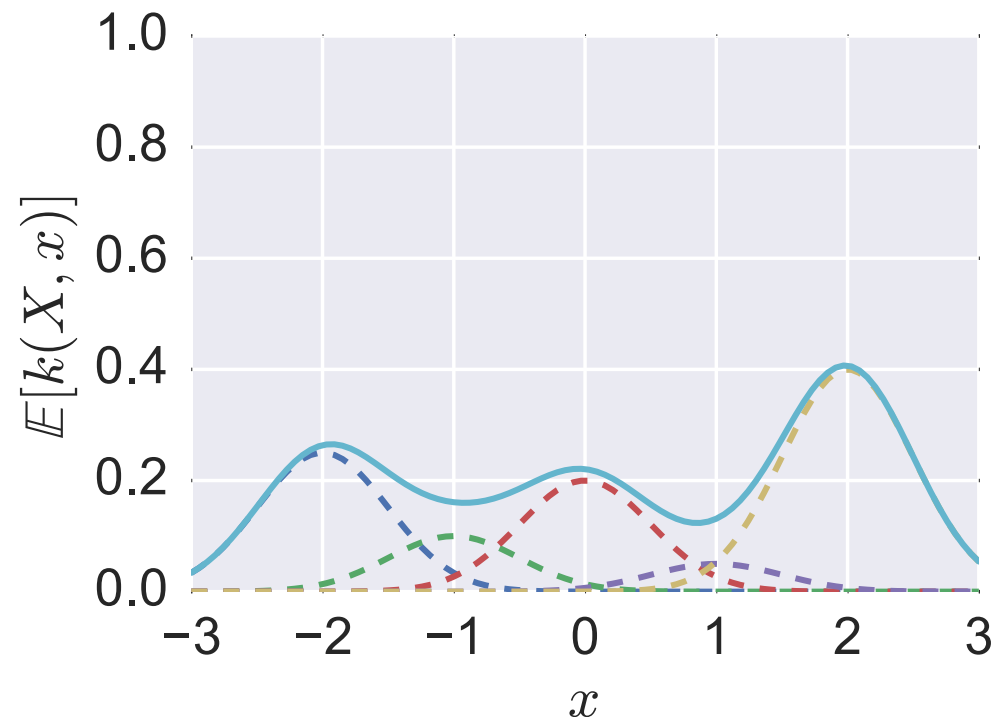
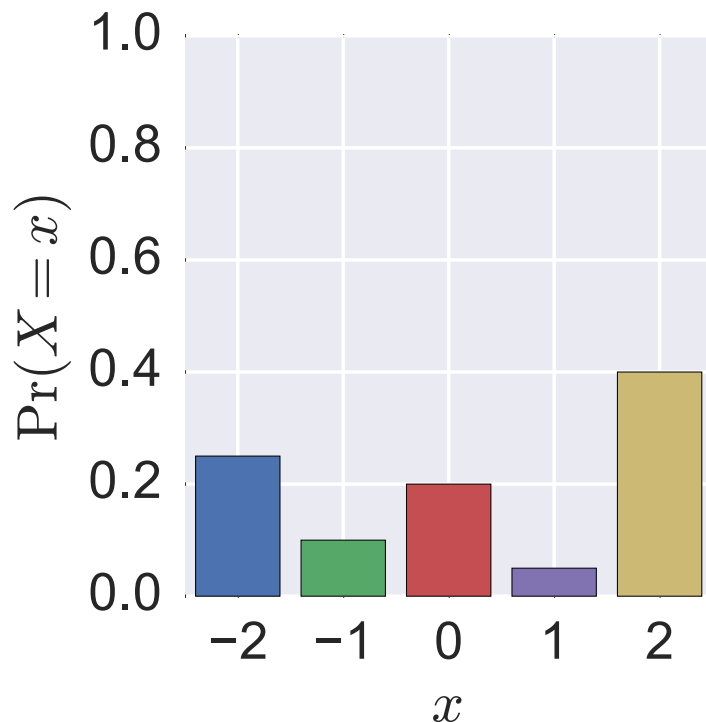


Mean embeddings of distributions

The *mean embedding* of a distribution in an RKHS:

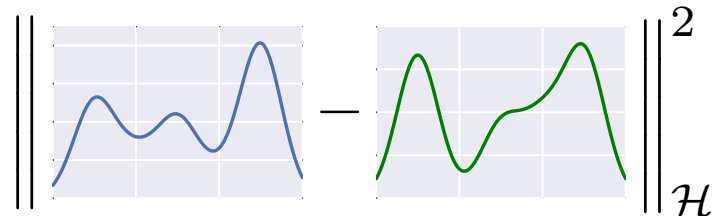
$$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)]$$

Remember $\langle \varphi(x), \varphi(y) \rangle = k(x, y)$,
so we can think of $\varphi(x)$ as $k(\cdot, x)$.



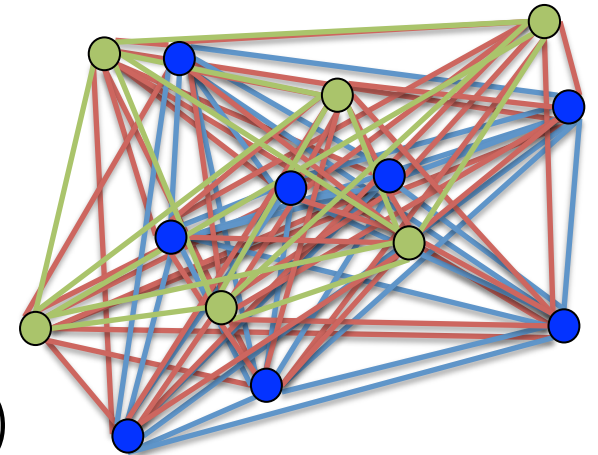
Maximum Mean Discrepancy (MMD)

The MMD is the distance between mean embeddings:

$$\begin{aligned}\mu_{\mathbb{P}} &= \mathbb{E}_{X \sim \mathbb{P}}[\varphi(X)] \\ \text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}}^2 \\ &= \langle \mu_{\mathbb{P}}, \mu_{\mathbb{P}} \rangle + \langle \mu_{\mathbb{Q}}, \mu_{\mathbb{Q}} \rangle - 2\langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle\end{aligned}$$


$$\begin{aligned}\langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle &= \langle \mathbb{E}_{X \sim P}[\varphi(X)], \mathbb{E}_{Y \sim Q}[\varphi(Y)] \rangle \\ &= \mathbb{E}_{\substack{X \sim P \\ Y \sim Q}} [\langle \varphi(X), \varphi(Y) \rangle] \\ &= \mathbb{E}_{\substack{X \sim P \\ Y \sim Q}} [k(X, Y)]\end{aligned}$$

Estimator: $\langle \hat{\mu}_{\mathbb{P}}, \hat{\mu}_{\mathbb{Q}} \rangle = \frac{1}{mn} \sum_{ij} k(X_i, Y_j)$



Embedding MMD

$$\widehat{\text{MMD}}(X, Y) = \|\hat{\mu}_{\mathbb{P}} - \hat{\mu}_{\mathbb{Q}}\|_{\mathcal{H}}$$

$$\langle \hat{\mu}_{\mathbb{P}}, \hat{\mu}_{\mathbb{Q}} \rangle = \frac{1}{mn} \sum_{ij} k(X_i, Y_j) = \frac{1}{m} \mathbb{1}^T \begin{matrix} \text{matrix} \end{matrix} \frac{1}{n} \mathbb{1}$$

$$\approx \frac{1}{mn} \sum_{ij} z(X_i)^T z(Y_j) = \frac{1}{m} \mathbb{1}^T \left(\begin{matrix} \text{column} \end{matrix} \times \begin{matrix} \text{row} \end{matrix} \right) \frac{1}{n} \mathbb{1}$$

$$= \left[\frac{1}{m} \sum_i z(X_i) \right]^T \left[\frac{1}{n} \sum_j z(Y_j) \right] = \left(\begin{matrix} \text{row} \end{matrix} \frac{1}{m} \mathbb{1} \right)^T \left(\begin{matrix} \text{row} \end{matrix} \frac{1}{n} \mathbb{1} \right)$$

$$= \bar{z}(X)^T \bar{z}(Y) = \begin{matrix} \text{row} \end{matrix}^T \begin{matrix} \text{row} \end{matrix}$$

$$\widehat{\text{MMD}}_z(X, Y) = \|\bar{z}(X) - \bar{z}(Y)\| = \left\| \begin{matrix} \text{row} \end{matrix} - \begin{matrix} \text{row} \end{matrix} \right\|$$

$$\exp \left(-\frac{1}{2\sigma^2} \widehat{\text{MMD}}_z(X, Y) \right) \approx z(\bar{z}(X))^T z(\bar{z}(Y))$$

Embedding L_2

Oliva, Neiswanger, Póczos, Schneider, Xing (AISTATS 2014)
gave an embedding for

$$K(\hat{p}, \hat{q}) = \exp \left(-\frac{1}{2\sigma^2} \|\hat{p} - \hat{q}\|_{L_2}^2 \right)$$

based on projection coefficients onto an orthonormal basis for L_2 .

Embedding HDDs

$$\rho^2(p, q) = \int_{\mathcal{X}} \kappa(p(x), q(x)) \, dx \quad \text{Homogeneous Density Distances}$$

ρ^2 can be:

- Jensen-Shannon
- Total Variation
- Squared Hellinger

Algorithm: HDD Embedding

$$\rho^2(p, q) = \int_{\mathcal{X}} \kappa(p(x), q(x)) \, dx \quad \text{Homogeneous Density Distances}$$

1. Approximately embed ρ into L_2
 - by embedding κ into L_2 (Fuglede 2005).

$$\rho(p, q) \approx \|\psi(p) - \psi(q)\|_{L_2^{2M}}$$

ρ on distributions $\approx L_2$ distance on random ψ functions

2. Approximately embed L_2 into \mathbf{R}^m .

$$\|\psi(p) - \psi(q)\|_{L_2^{2M}} \approx \|A(p) - A(q)\|_{\mathbb{R}^{2M|V|}}$$

L_2 distance on random ψ functions \approx Euclidean distance on A vectors

3. Use random Fourier features to embed K on \mathbf{R}^m into \mathbf{R}^D .

$$\exp\left(-\frac{1}{2\sigma^2} \rho^2(p, q)\right) \approx z(A(p))^T z(A(q))$$

RBF kernel with $\rho \approx$ random Fourier features of A vectors

HDD estimator convergence

The embedding is an estimator for the kernel.

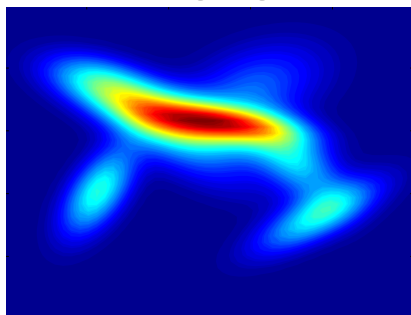
For fixed p and q , we have a finite-sample bound on

$$\Pr \left(\left| K(p, q) - z(\hat{A}(\hat{p}))^\top z(\hat{A}(\hat{q})) \right| \geq \varepsilon \right)$$

which behaves as expected:

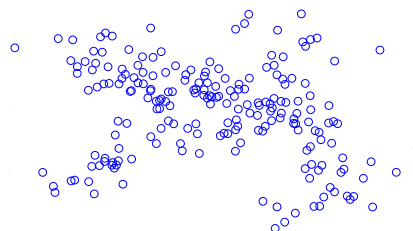
- Lower for smoother, lower-dimensional densities
- Lower for more samples
- More projection coefficients / samples from μ :
 - better approximation, harder integration

Application: Gaussian Mixtures



Underlying density

Sample

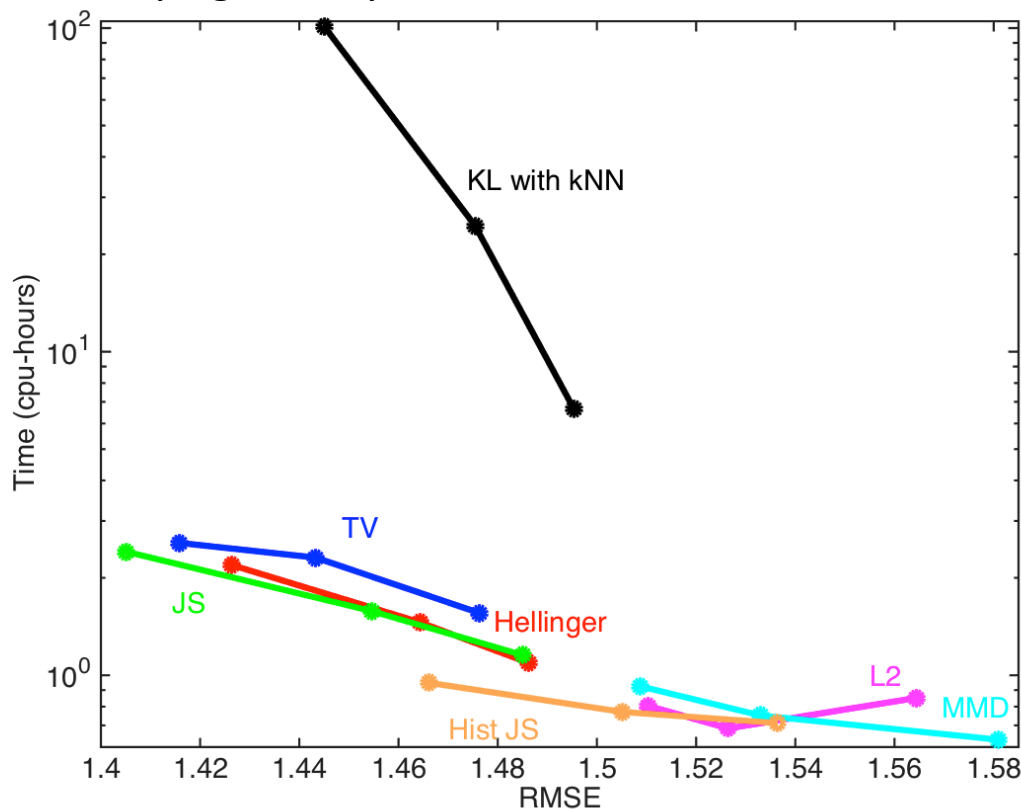


Observed sample

Regression

9

Prediction



Train with 4K, 8K, 16K distributions; test on 2K. 200 iid points each.

AIC: 2.7

BIC: 3.8

Mean: 2.8

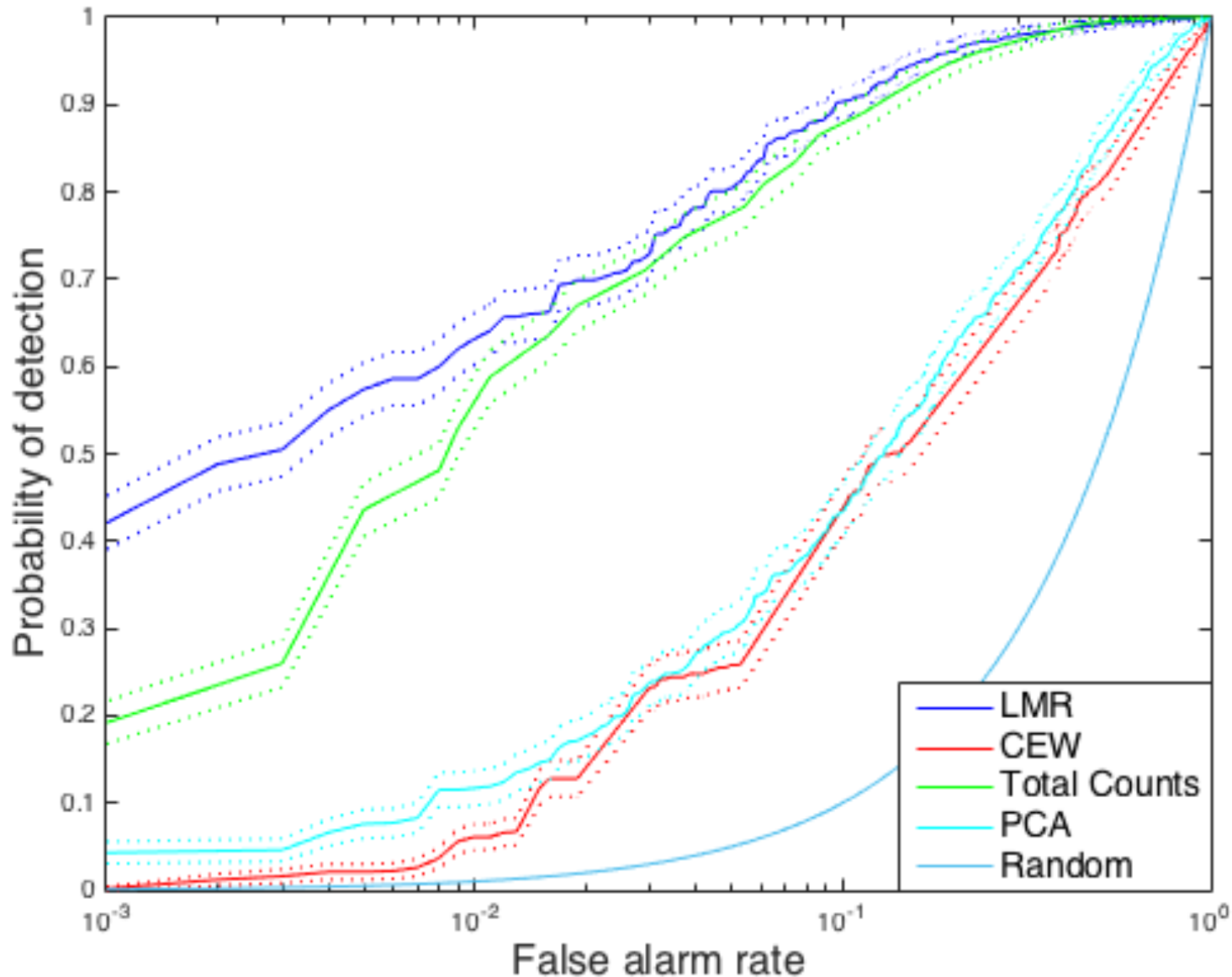
Application: Nuclear Threat Detection



[Rapiscan MP100](#)

Small sensors and cluttered environments make lots of challenges for traditional detection algorithms.

Application: Nuclear Threat Detection



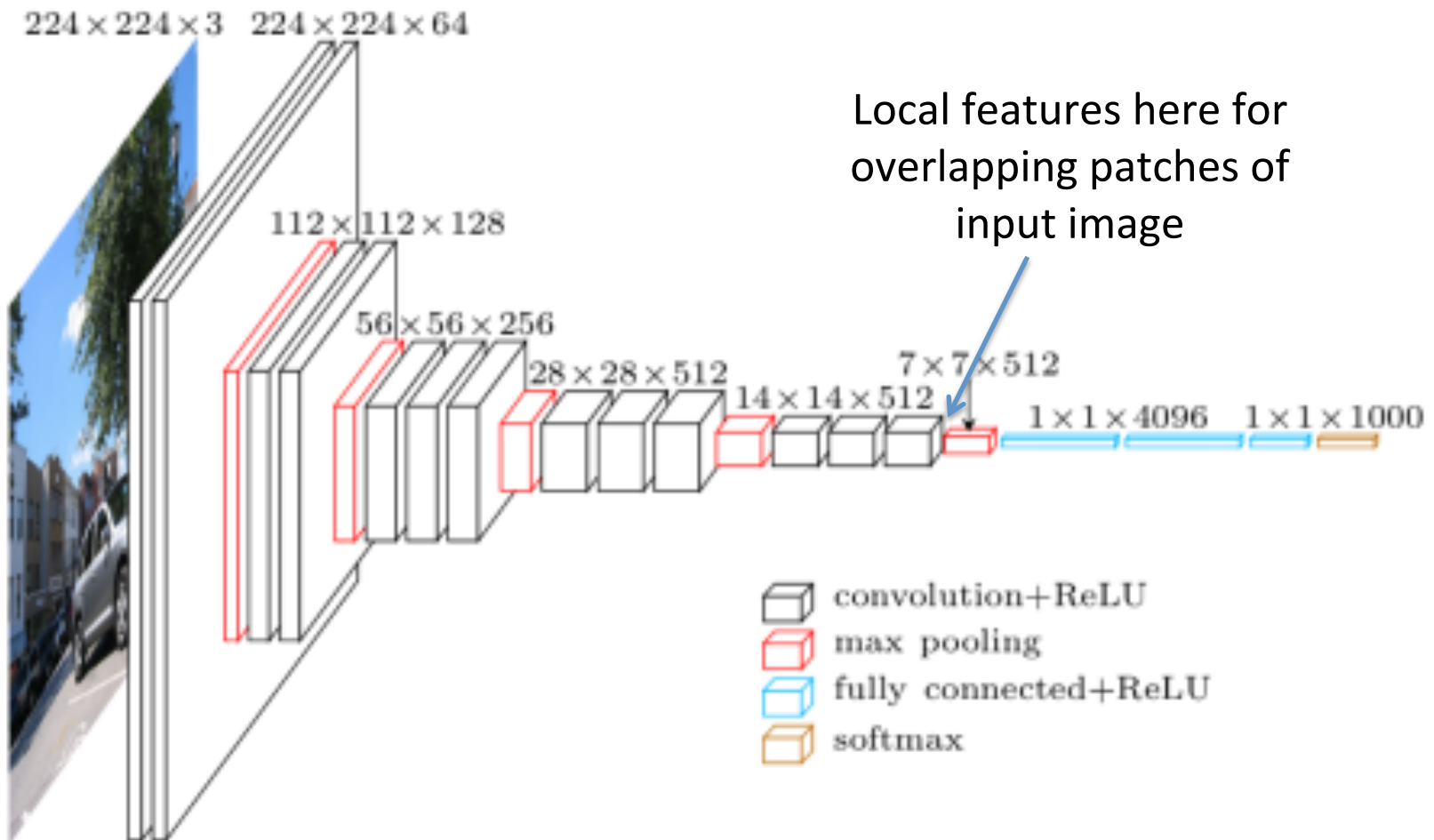
Application: Scene Classification



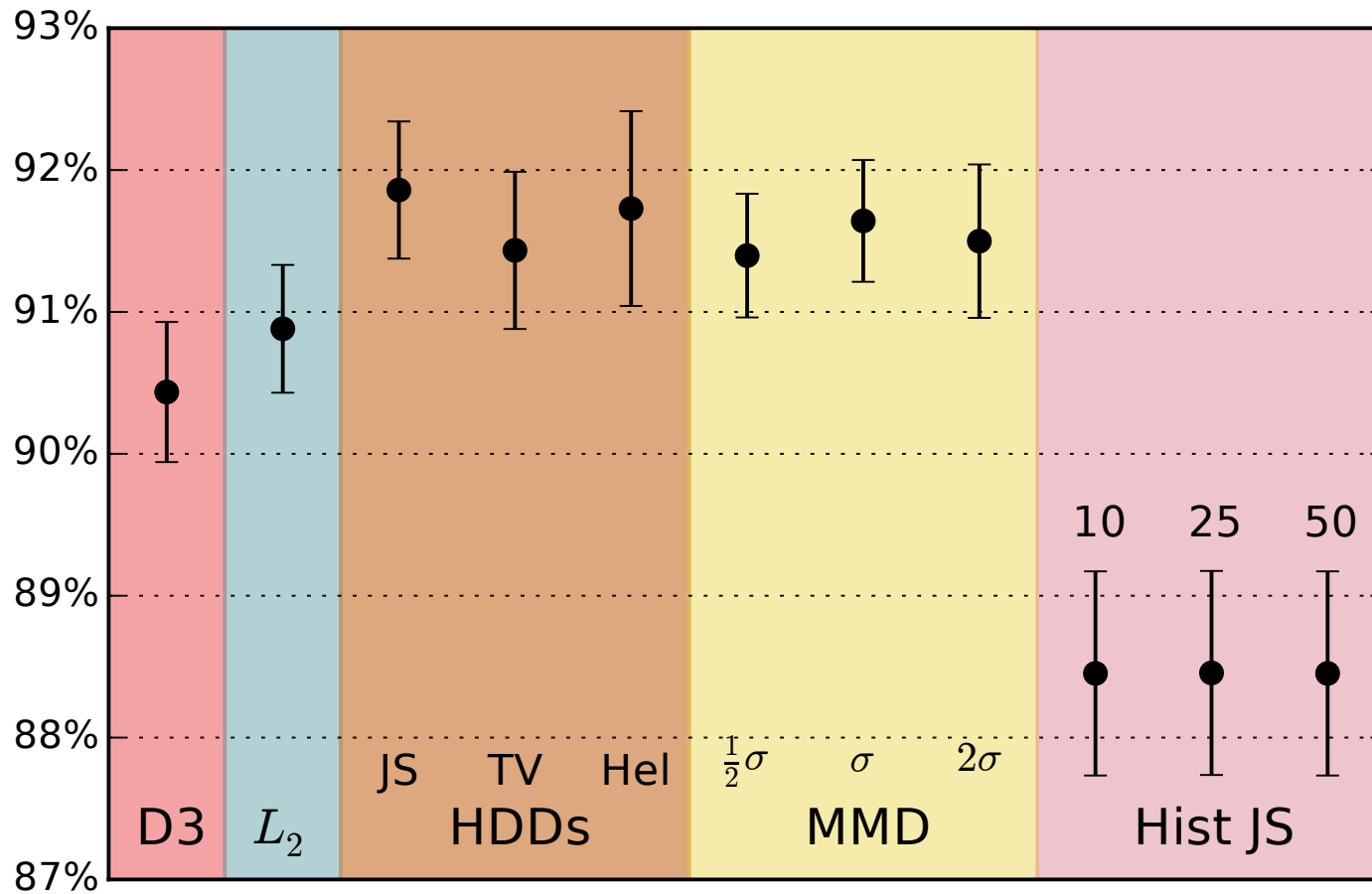
Scene-15 dataset (4 485 images)

Application: Scene Classification

Features from a deep convnet (16 layers):



Application: Scene Classification



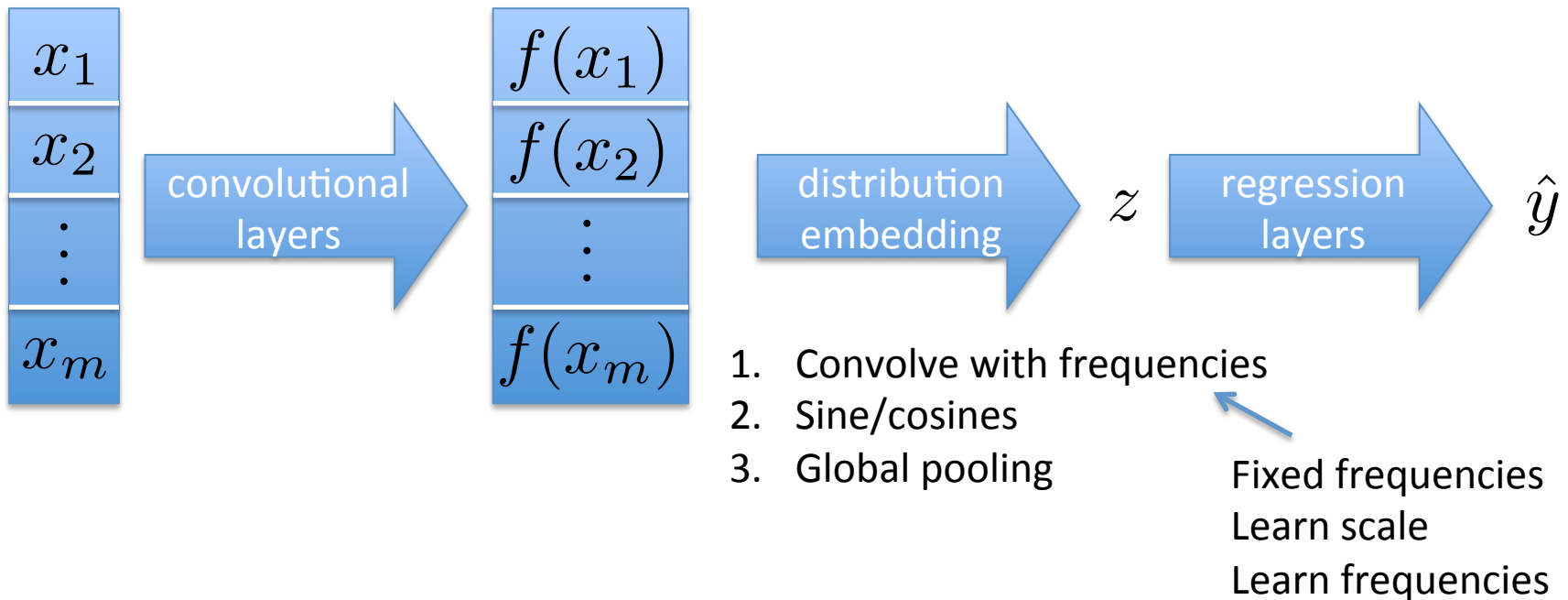
Using $A(p)$ features

Contributions

- **Learning on distributions** with nonparametric kernels
- **Scalable** approximate kernel embeddings
 - Random Fourier features analysis
 - New embeddings for distribution kernels
- **Flexible** distribution kernels
 - Deep mean maps in computer vision
 - MMD kernel learning for testing
- **Active** pointillistic pattern search

Deep distribution kernel learning

We can put distribution embeddings in a deep network:



Initial results for scene classification: small but consistent improvement.

Contributions

- **Learning on distributions** with nonparametric kernels
- **Scalable** approximate kernel embeddings
 - Random Fourier features analysis
 - New embeddings for distribution kernels
- **Flexible** distribution kernels
 - Deep mean maps in computer vision
 - MMD kernel learning for testing
- **Active** pointillistic pattern search

Two-sample testing

Observed samples: $X \sim \mathbb{P}$ $Y \sim \mathbb{Q}$

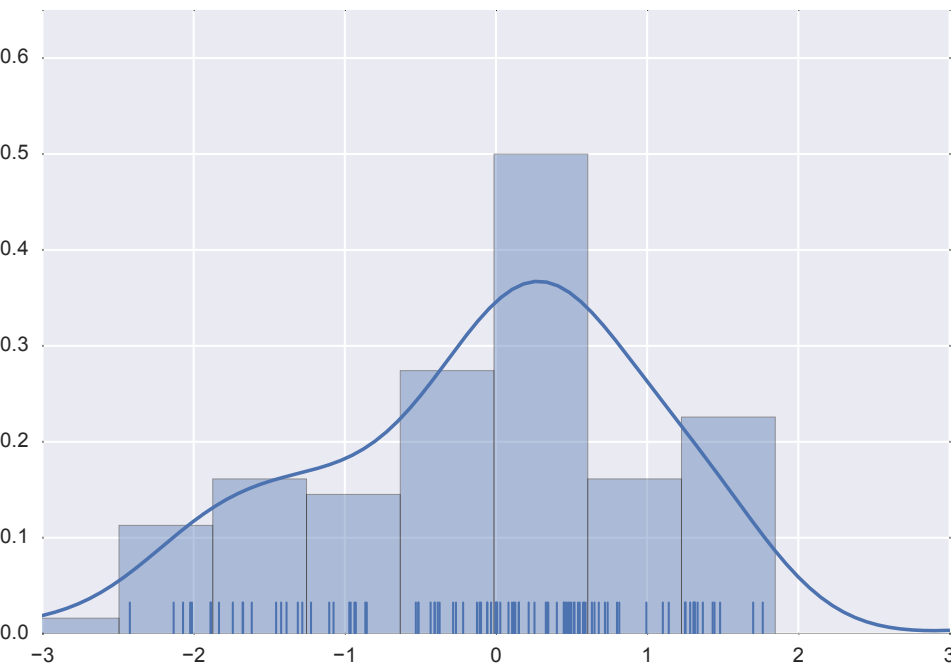
Hypothesis test on unobserved distributions: $\mathbb{P} \stackrel{?}{=} \mathbb{Q}$

Applications:

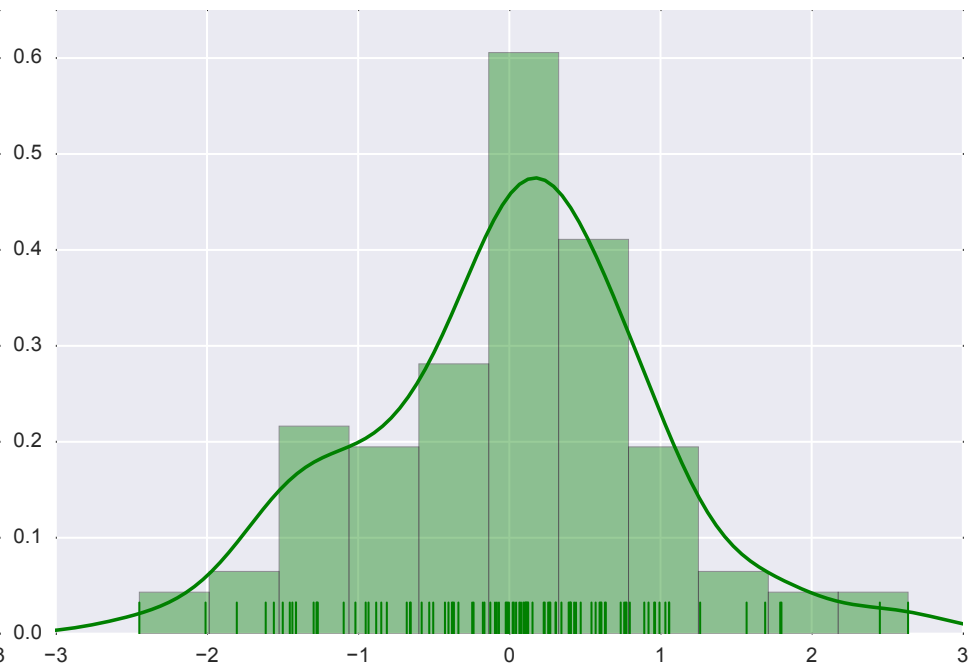
- Neuroscience: do these areas of the brain behave differently under different conditions?
- Schema alignment: do these two database columns mean the same thing?
- Many more...

Two-sample testing

Test with MMD:



$$X \sim \mathbb{P} = \mathcal{N}(0, 1)$$



$$Y \sim \mathbb{Q} = \text{Laplace}\left(0, \frac{1}{\sqrt{2}}\right)$$

Test: reject if $m \widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q}) > c_\alpha$.

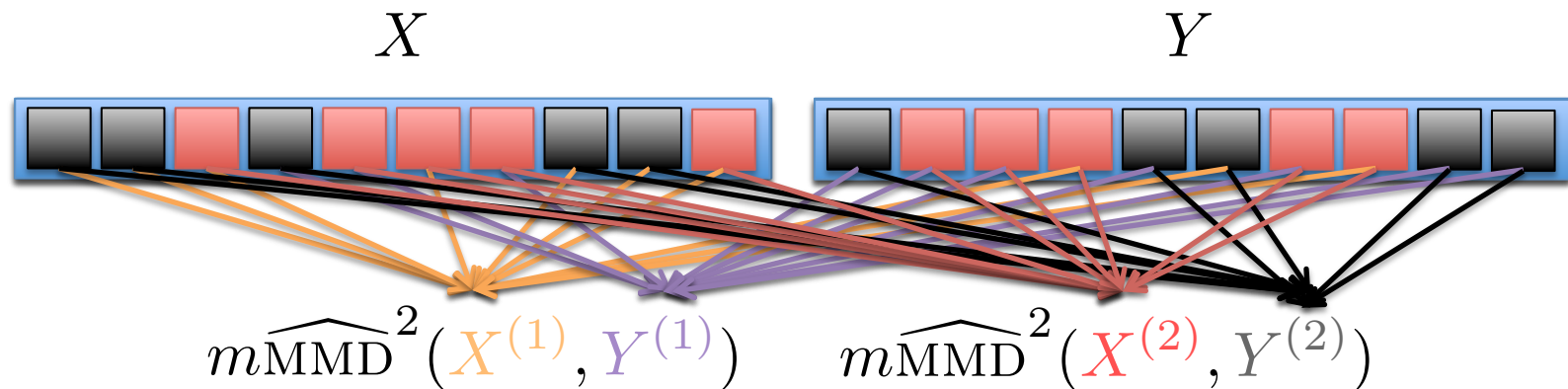
Permutation testing

When $\mathbb{P} = \mathbb{Q}$, MMD is asymptotically gross, so hard to find a threshold c_α that way.

Use a permutation test:

$$X = \{x_1, x_2, \dots, x_m\} \quad Y = \{y_1, y_2, \dots, y_m\}$$

Split randomly to estimate MMD when $\mathbb{P} = \mathbb{Q}$:



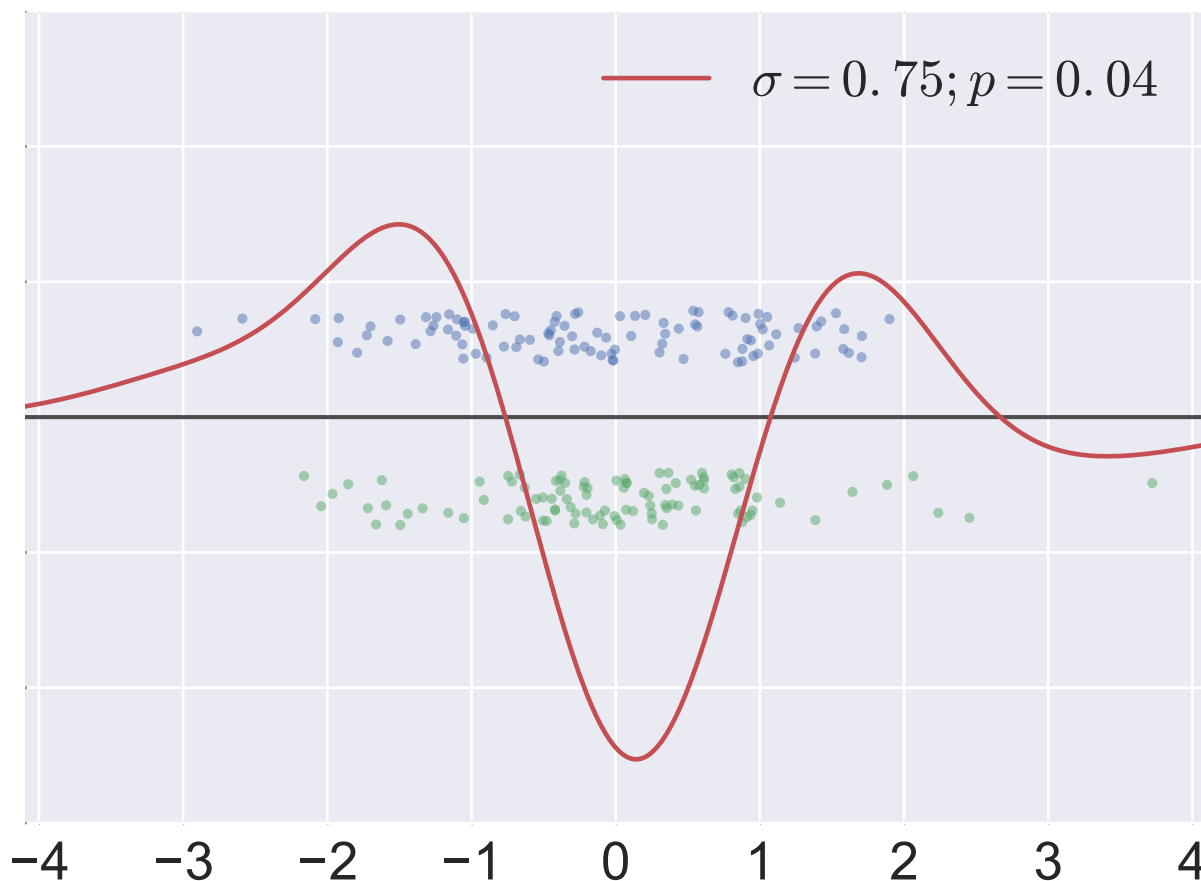
\hat{c}_α : $(1-\alpha)$ th quantile of $m\widehat{\text{MMD}}^2(X^{(i)}, Y^{(i)})$

The kernel matters!

Witness function f helps compare samples:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y)$$

$$f(x) = \mathbb{E}_{X \sim \mathbb{P}} k(x, X) - \mathbb{E}_{Y \sim \mathbb{Q}} k(x, Y)$$

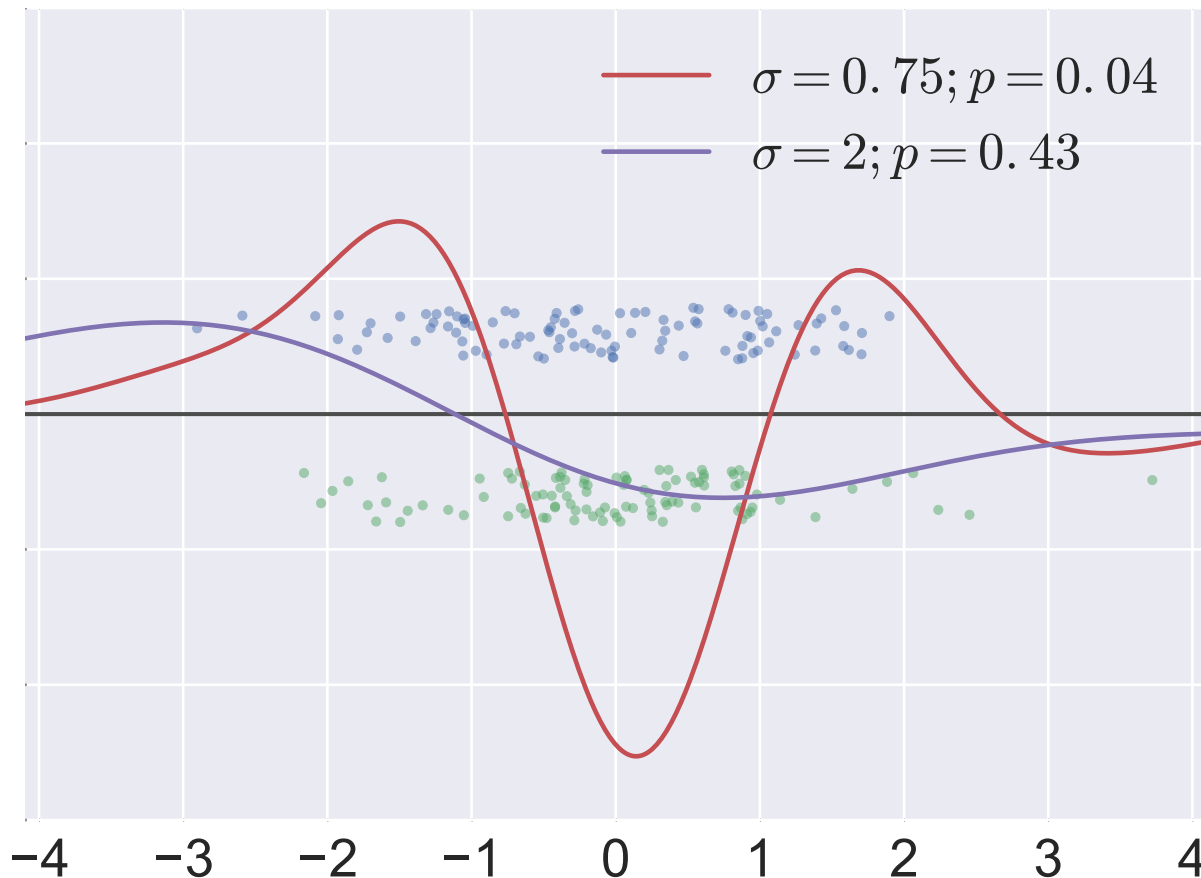


The kernel matters!

Witness function f helps compare samples:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y)$$

$$f(x) = \mathbb{E}_{X \sim \mathbb{P}} k(x, X) - \mathbb{E}_{Y \sim \mathbb{Q}} k(x, Y)$$

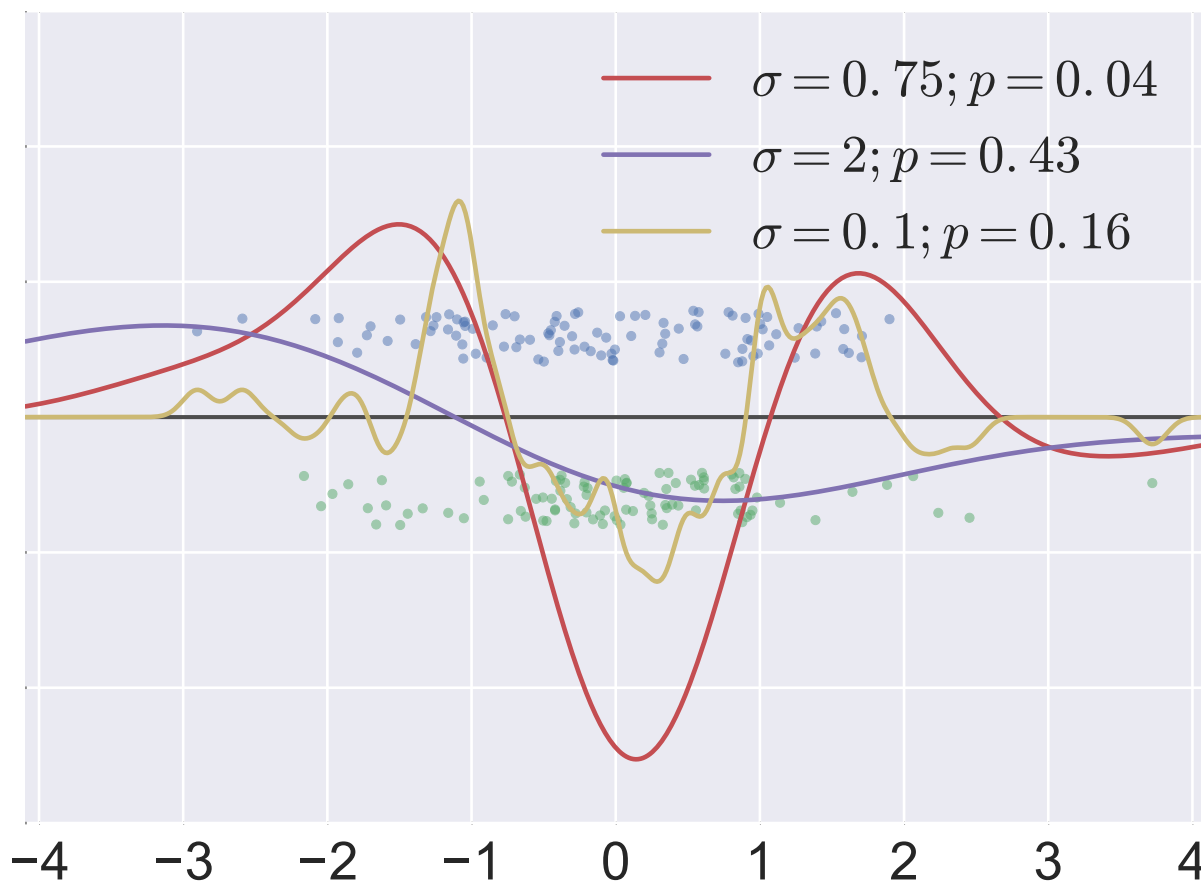


The kernel matters!

Witness function f helps compare samples:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y)$$

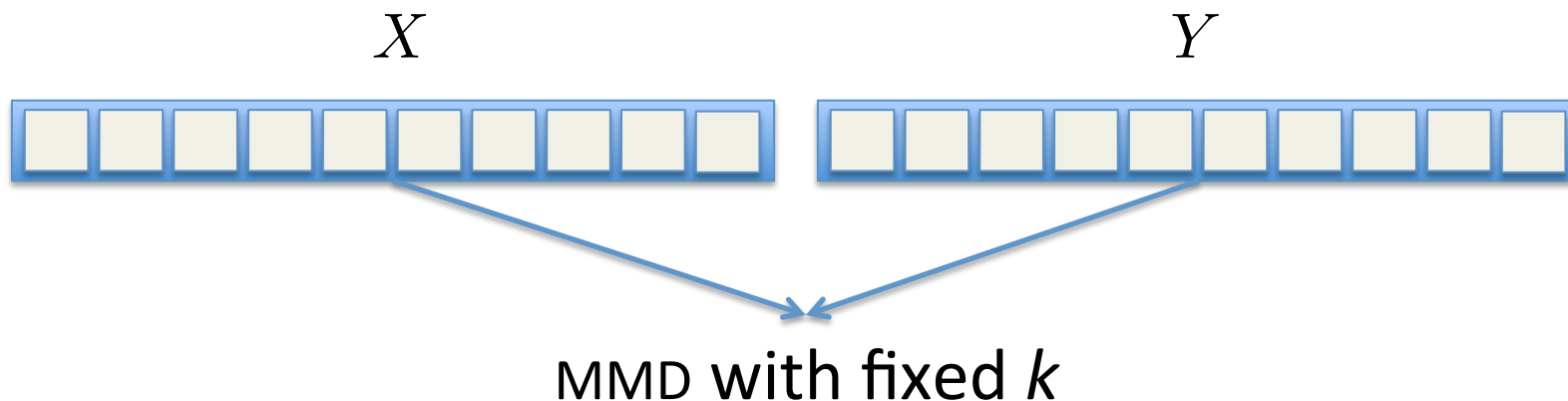
$$f(x) = \mathbb{E}_{X \sim \mathbb{P}} k(x, X) - \mathbb{E}_{Y \sim \mathbb{Q}} k(x, Y)$$



Choosing a kernel

So we need a way to pick a kernel to do the test.

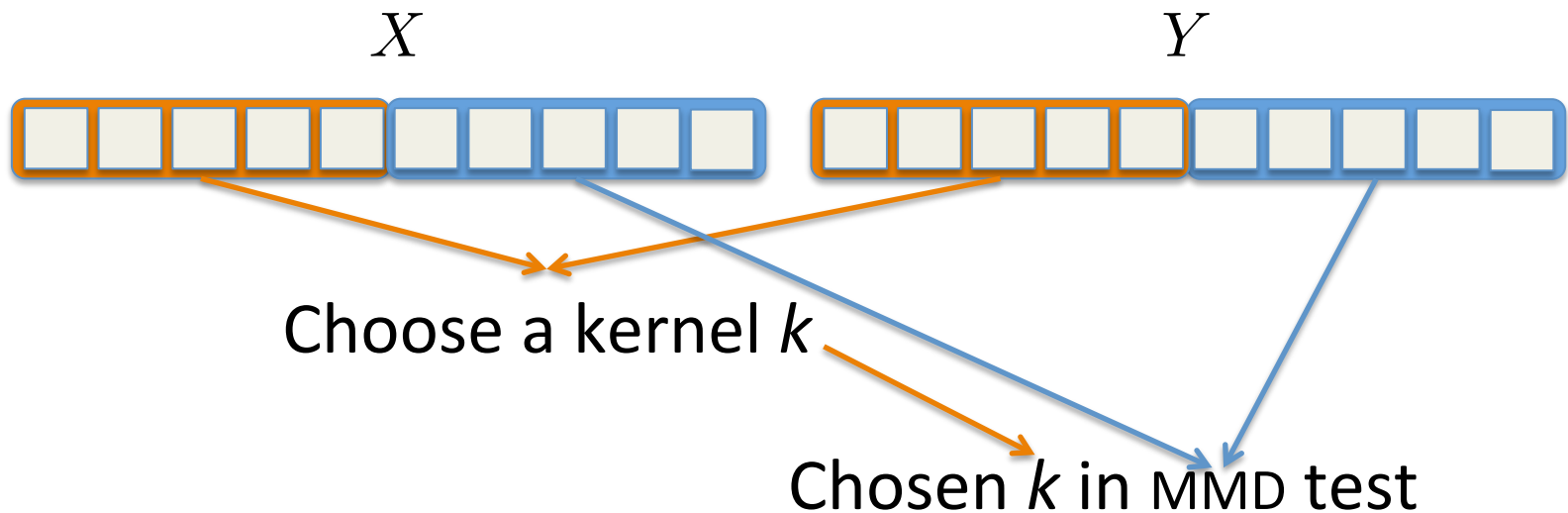
Before:



Choosing a kernel

So we need a way to pick a kernel to do the test.

Split data:



How to pick k ? Typically: maximize MMD.

But we want the (asymptotically) *most powerful test*. 54

Asymptotic power of MMD

When $\mathbb{P} \neq \mathbb{Q}$, MMD is asymptotically normal:

$$\frac{\widehat{\text{MMD}}^2 - \text{MMD}^2}{\sqrt{V_m}} \xrightarrow{D} \mathcal{N}(0, 1) \quad V_m = \text{Var}_{\substack{X \sim P^m \\ Y \sim Q^m}} \left[\widehat{\text{MMD}}^2(X, Y) \right]$$

and we can analyze the power:

$$\Pr_{H_1} \left(m \widehat{\text{MMD}}^2 > \hat{c}_\alpha \right)$$

MMD t -statistic

$$\Pr_{H_1} \left(m \widehat{\text{MMD}}^2 > \hat{c}_\alpha \right) \rightarrow 1 - \Phi \left(\frac{c_\alpha}{m\sqrt{V_m}} - \frac{\text{MMD}^2}{\sqrt{V_m}} \right)$$

So we can maximize the power by maximizing

$$\tau_U = \frac{\text{MMD}^2}{\sqrt{V_m}} - \frac{c_\alpha}{m\sqrt{V_m}} \qquad \hat{\tau}_U = \frac{\widehat{\text{MMD}}^2}{\sqrt{\hat{V}_m}} - \frac{\hat{c}_\alpha}{m\sqrt{\hat{V}_m}}$$

But V_m is $O(1/m)$, so the first term dominates for large m , and we should be able to get away with maximizing

$$t_U = \frac{\text{MMD}^2}{\sqrt{V_m}} \qquad \hat{t}_U = \frac{\widehat{\text{MMD}}^2}{\sqrt{\hat{V}_m}}$$

t -statistic estimator

$$\hat{\tau}_U = \frac{\widehat{\text{MMD}}^2}{\sqrt{\hat{V}_m}} - \frac{\hat{c}_\alpha}{m\sqrt{\hat{V}_m}}$$

$$\widehat{\text{MMD}}^2 := \frac{1}{\binom{m}{2}} \sum_{i \neq j} k(X_i, X_j) + k(Y_i, Y_j) - k(X_i, Y_j) - k(X_j, Y_i)$$

\hat{c}_α is from a permutation test, so the average of a bunch of MMD estimates

t -statistic estimator

$$\begin{aligned}
 \widehat{V}_m &= \frac{4(m-2)}{m(m-1)} \hat{\zeta}_1 + \frac{2}{m(m-1)} \hat{\zeta}_2 \\
 \hat{\zeta}_1 &= \frac{1}{m(m-1)(m-2)} \left(\mathbf{1}^\top \tilde{K}_{XX} \tilde{K}_{XX} \mathbf{1} - \|\tilde{K}_{XX}\|_F^2 \right) - \left(\frac{1}{m(m-1)} \mathbf{1}^\top \tilde{K}_{XX} \mathbf{1} \right)^2 \\
 &\quad - \frac{2}{m^2(m-1)} \mathbf{1}^\top \tilde{K}_{XX} K_{XY} \mathbf{1} + \frac{2}{m^3(m-1)} \mathbf{1}^\top \tilde{K}_{XX} \mathbf{1} \mathbf{1}^\top K_{XY} \mathbf{1} \\
 &\quad + \frac{1}{m(m-1)(m-2)} \left(\mathbf{1}^\top \tilde{K}_{YY} \tilde{K}_{YY} \mathbf{1} - \|\tilde{K}_{YY}\|_F^2 \right) - \left(\frac{1}{m(m-1)} \mathbf{1}^\top \tilde{K}_{YY} \mathbf{1} \right)^2 \\
 &\quad - \frac{2}{m^2(m-1)} \mathbf{1}^\top \tilde{K}_{YY} K_{XY}^\top \mathbf{1} + \frac{2}{m^3(m-1)} \mathbf{1}^\top \tilde{K}_{YY} \mathbf{1} \mathbf{1}^\top K_{XY} \mathbf{1} \\
 &\quad + \frac{1}{m^2(m-1)} \left(\mathbf{1}^\top K_{XY}^\top K_{XY} \mathbf{1} - \|K_{XY}\|_F^2 \right) - 2 \left(\frac{1}{m^2} \mathbf{1}^\top K_{XY} \mathbf{1} \right)^2 \\
 &\quad + \frac{1}{m^2(m-1)} \left(\mathbf{1}^\top K_{XY} K_{XY}^\top \mathbf{1} - \|K_{XY}\|_F^2 \right) \\
 \hat{\zeta}_2 &= \frac{1}{m(m-1)} \left\| \tilde{K}_{XX} + \tilde{K}_{YY} - \tilde{K}_{XY} - \tilde{K}_{XY}^\top \right\|_F^2
 \end{aligned}$$

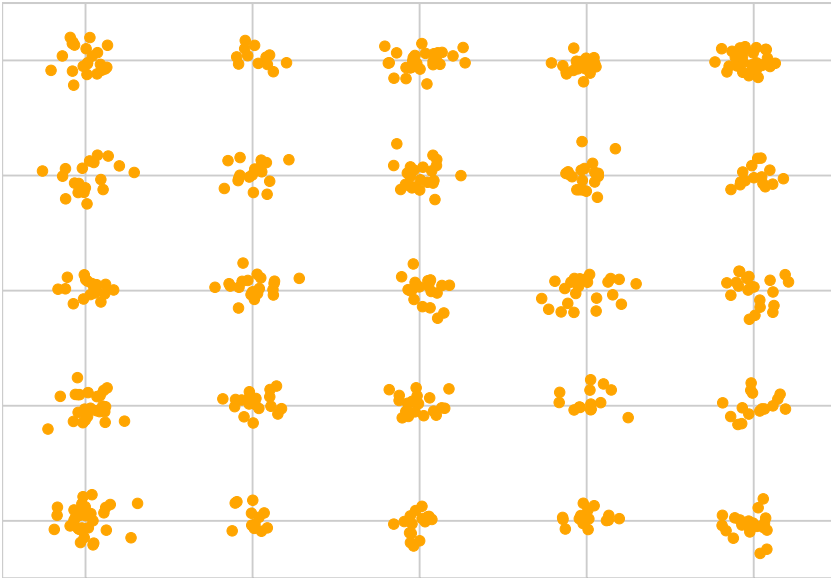
t -statistic estimator

Can even get gradients of t_U and (with some more effort) τ_U , to help maximize it.

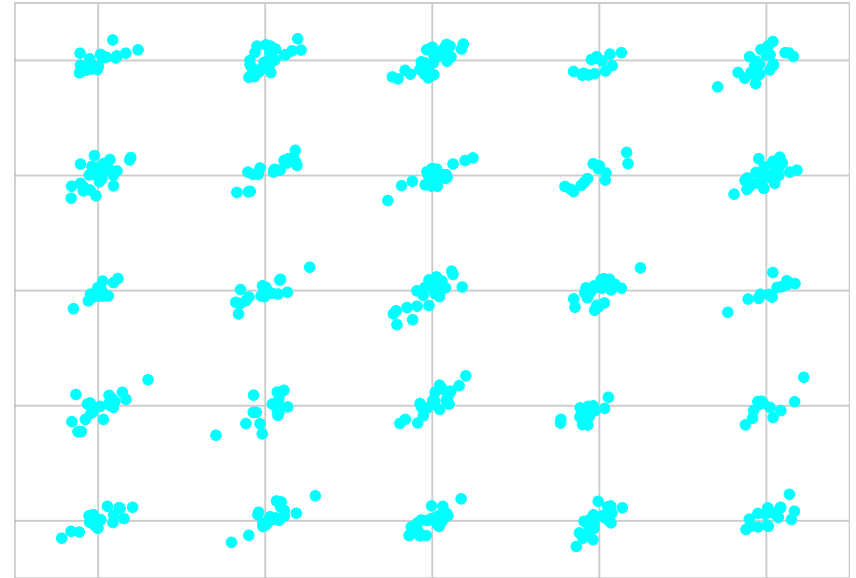
(automatic differentiation is your friend)

Kernel choice on Blobs

Blobs dataset:



vs

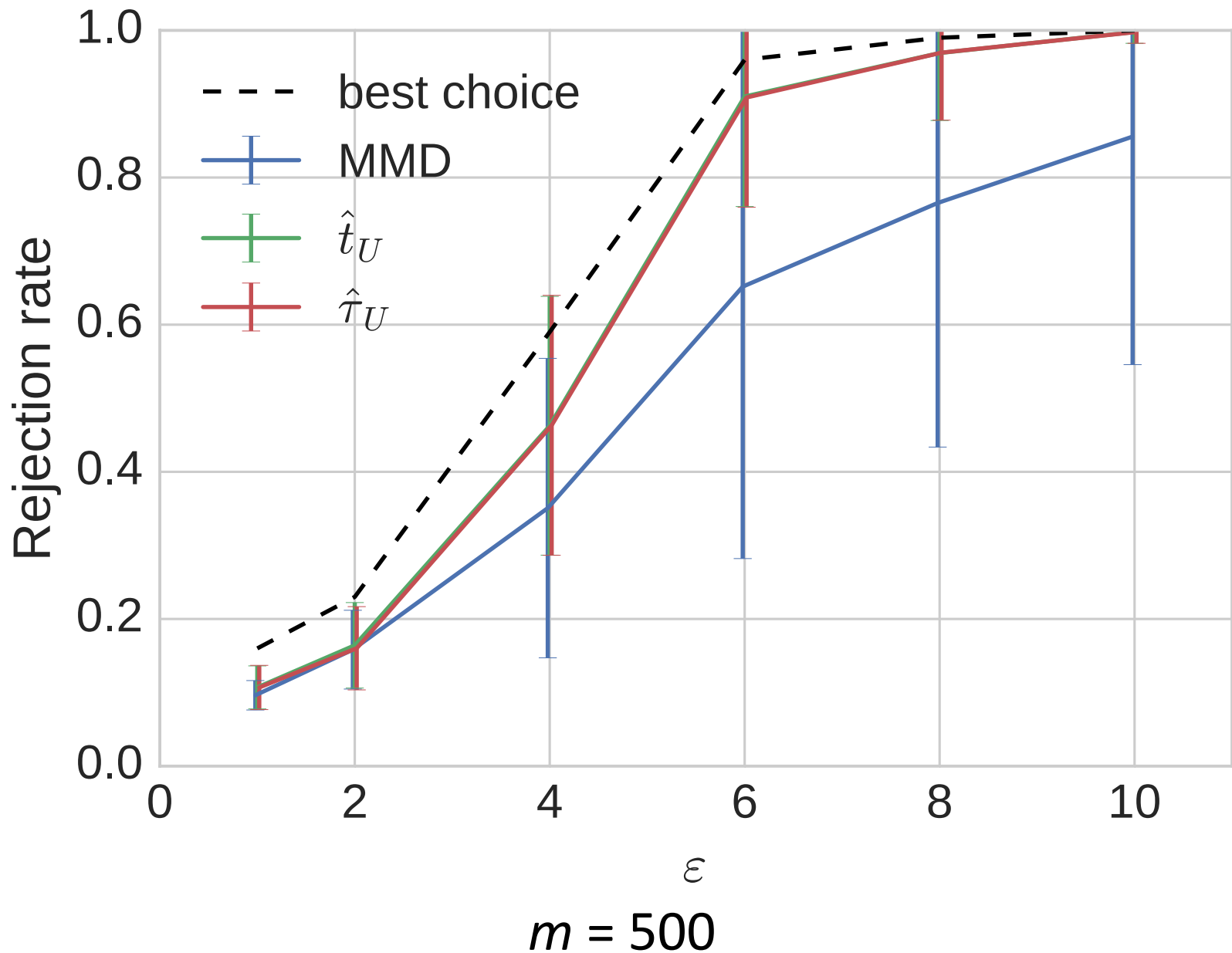


Mixture of $\mathcal{N}\left(\mu_{ij}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$

Mixture of $\mathcal{N}\left(\mu_{ij}, \begin{bmatrix} 1 & \frac{\varepsilon-1}{\varepsilon+1} \\ \frac{\varepsilon-1}{\varepsilon+1} & 1 \end{bmatrix}\right)$

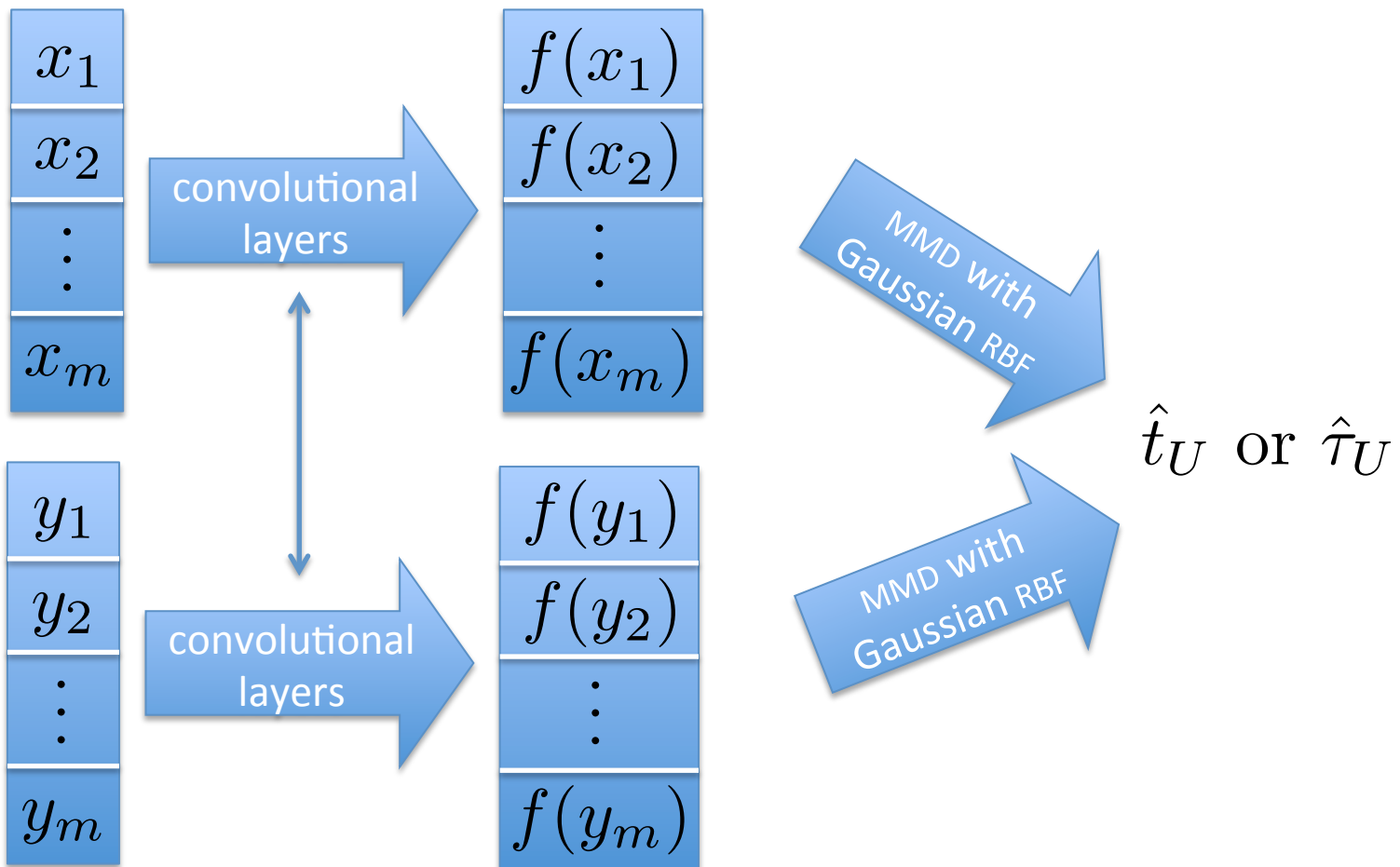
When $\varepsilon=1$, $P = Q$; this picture has $\varepsilon=6$.

Kernel choice on Blobs



Deep Kernels

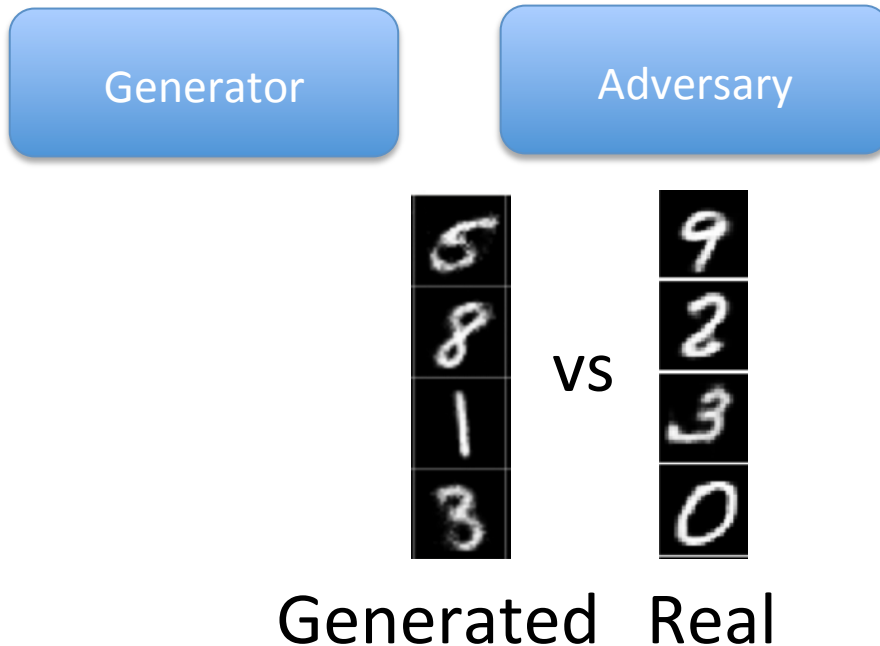
Map through layers of a deep network:



Generative Models

Generative adversarial networks:

- *Generator* comes up with samples; trained to trick the adversary.
- *Adversary* tries to distinguish between generator sample and true data; trained to beat the generator.



But adversary is really just a two-sample test.

Kernel learning helps prevent the generator's tricks.

Contributions

- **Learning on distributions** with nonparametric kernels
- **Scalable** approximate kernel embeddings
 - Random Fourier features analysis
 - New embeddings for distribution kernels
- **Flexible** distribution kernels
 - Deep mean maps in computer vision
 - MMD kernel learning for testing
- **Active** pointillistic pattern search

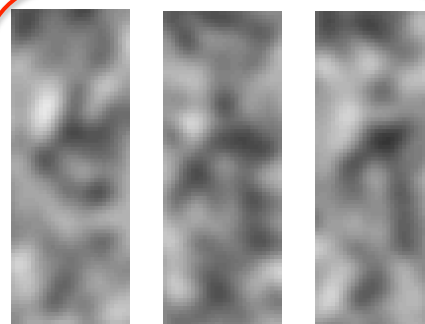
Active Pointillistic Pattern Search



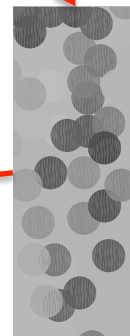
Search for *region patterns* with *point observations*.



template

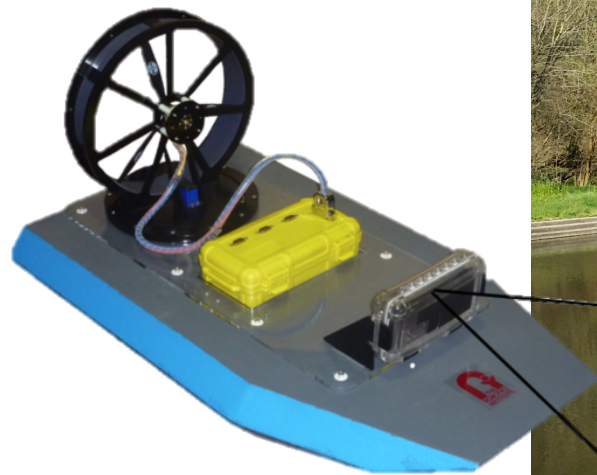


posterior draws

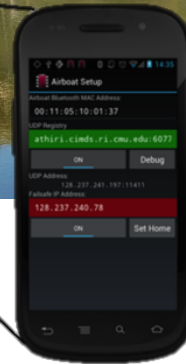


observations

Active Pointillistic Pattern Search

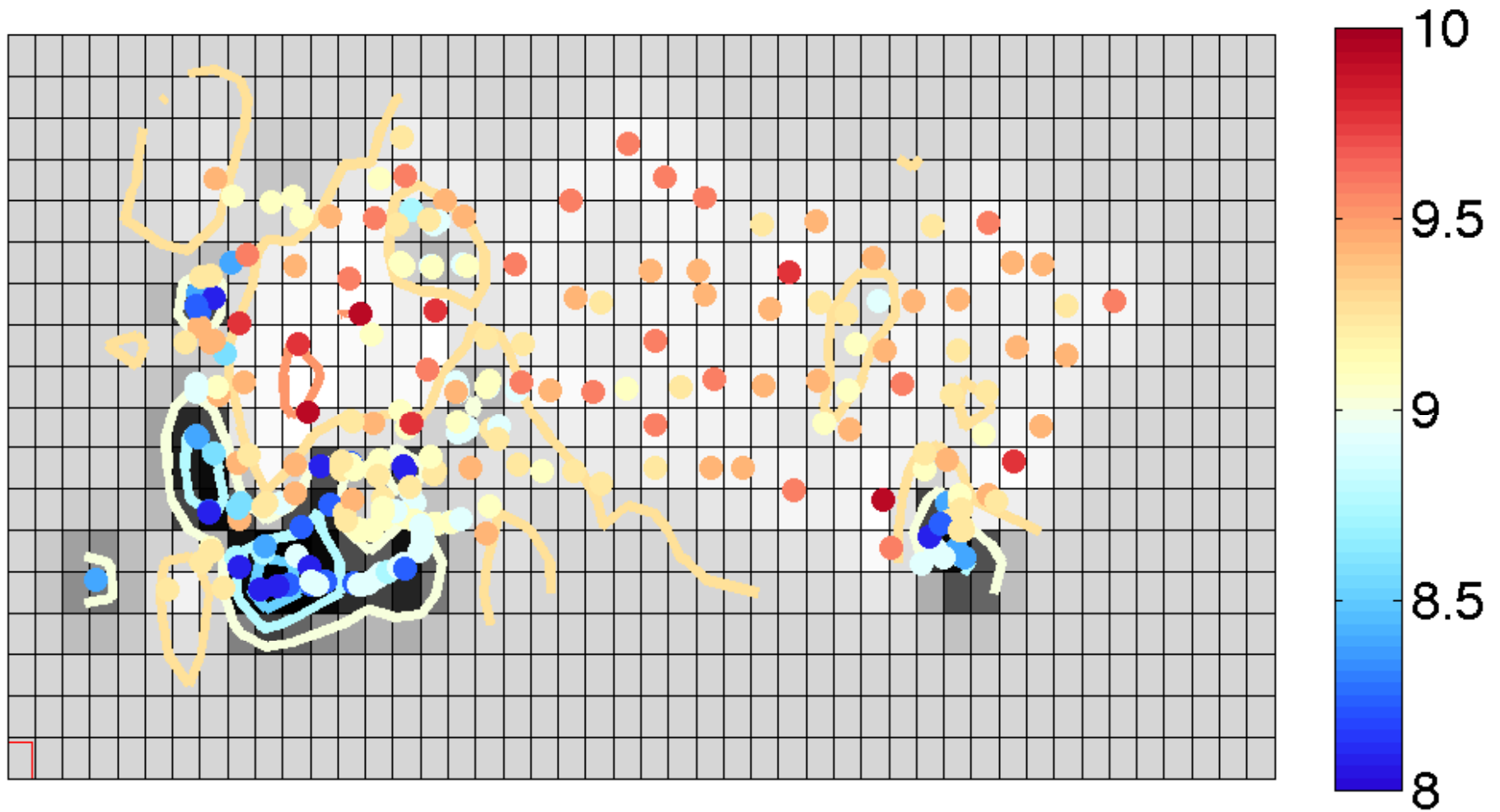


Watercraft

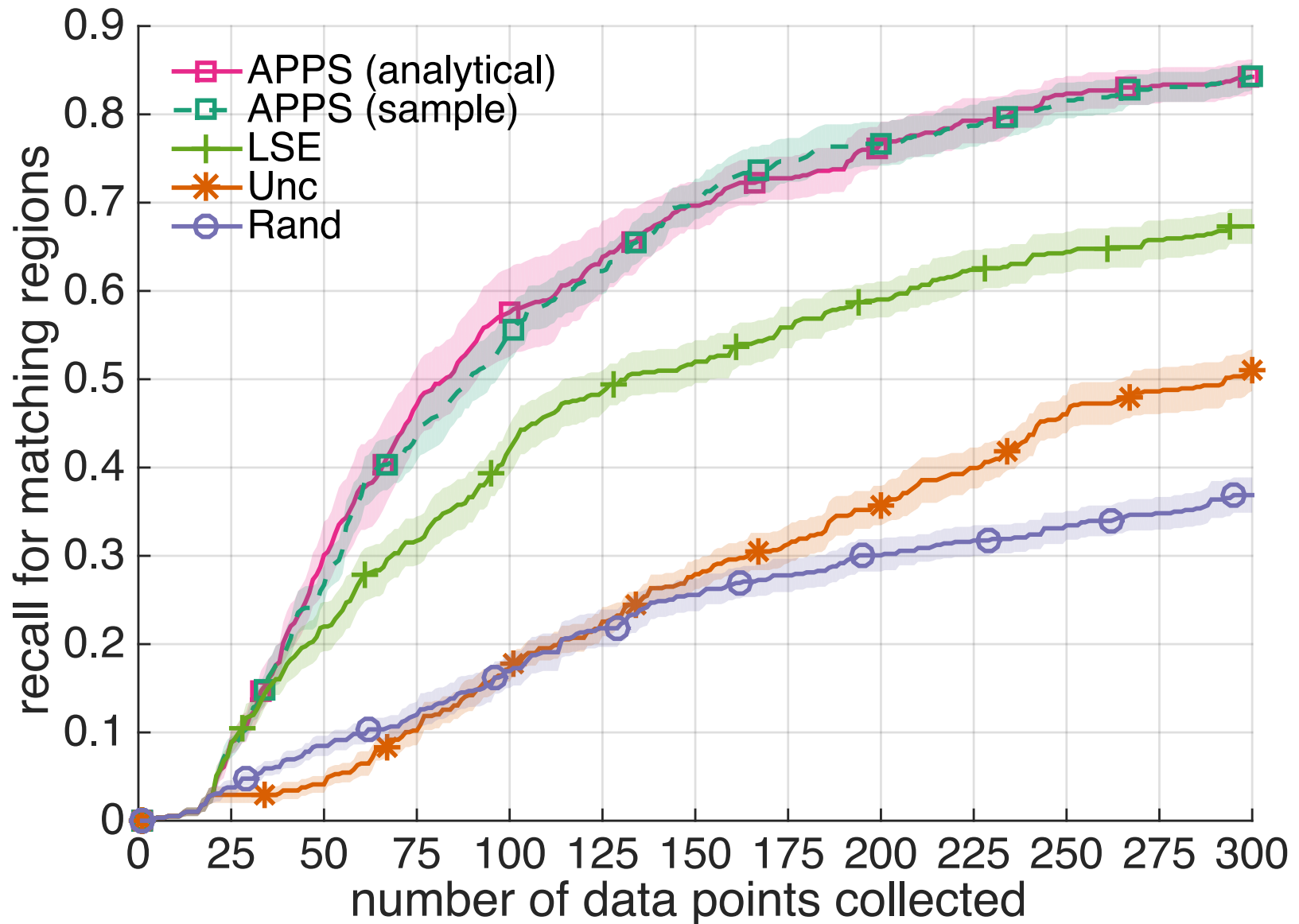


Android smartphone

Active Pointillistic Pattern Search



Active Pointillistic Pattern Search



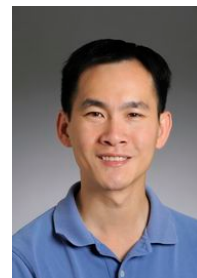
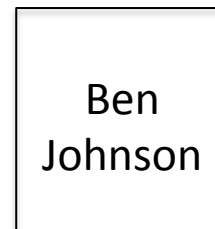
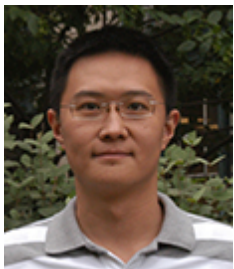
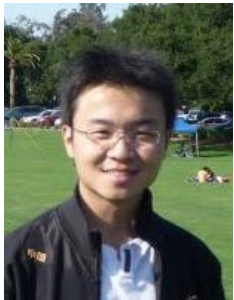
Take-Home Messages

- Think about how you model your data.
- Distributions and sets can work pretty well.
 - Cosmology, nuclear threat detection, scene classification, parametric statistical inference, polling, autonomous sensing...
- Random embeddings can help with scalability...
 - if you use random Fourier features, use the right one
- ...and with flexibility
 - Plug the MMD embedding into deep learning and go crazy

Things Still to Do

- Deep kernel learning
 - Different parameterizations of kernels
- More applications!
 - Word and document embeddings
 - Kernel-learning two-sample test as adversary in a GAN
- Active learning on distributions

Thanks!



Contributions

- **Learning on distributions** with nonparametric kernels
- **Scalable** approximate kernel embeddings
 - Random Fourier features analysis
 - New embeddings for distribution kernels
- **Flexible** distribution kernels
 - Deep mean maps in computer vision
 - MMD kernel learning for testing
- **Active** pointillistic pattern search