# Adventures in Kernel Density Estimation

## Clay Scott

joint with

## JooSeuk Kim, Robert Vandermeulen, and Efrén Cruz Cortés

ELECTRICAL &
COMPUTER ENGINEERING

UNIVERSITY OF MICHIGAN
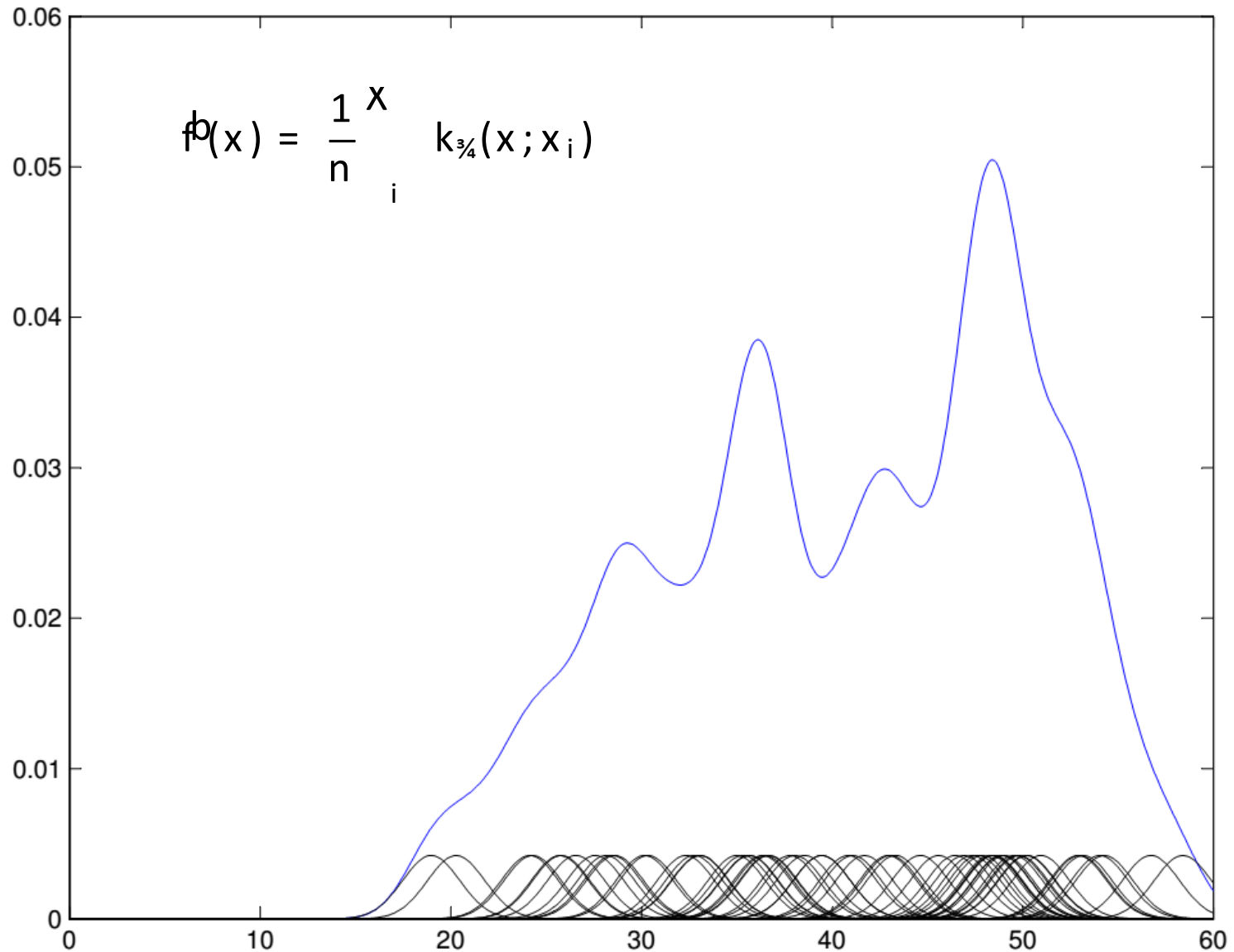
# Kernel Density Estimation

- $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{iid}{\sim} f$  (unknown density), $\mathbf{x}_i \in \mathbb{R}^d$

- Estimate $f$ via

$$\widehat{f}(\mathbf{x}) = \frac{1}{n} \sum_i k_\sigma(\mathbf{x}, \mathbf{x}_i)$$

- Example: Gaussian kernel

$$k_\sigma(\mathbf{x}, \mathbf{x}') = (2\pi\sigma^2)^{-d/2} \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

# Kernel Density Estimation



$$\hat{f}(x) = \frac{1}{n} \sum_i k_{\frac{3}{4}}(x\,;\,x_i)$$

# Applications of KDEs

Predictors based on KDEs: Classification, regression, anomaly detection

Estimates of information theoretic measures (entropy, KL divergence)

Clustering (e.g., image segmentation) via mean-shift algorithm

$(r; g; b; x; y) \; 2 \; \mathbb{R}^5$

gradient
ascent

mode of $\hat{f}$

# Conventional Analysis

² Estimation and approximation errors

$$\|\hat{f} - f\| \cdot \|\hat{f} - f * k_{\frac{3}{4}}\| + \|f * k_{\frac{3}{4}} - f\|$$

² Note

$$\hat{f} = \frac{1}{n} \sum_i k(\cdot, x_i)$$

= k convolved with empirical distribution

² Need $\frac{3}{4} \to 0$ for approximation error to vanish

² Need $n\frac{3}{4}^d \to 1$ for estimation error to vanish (by properties of convolutions and concentration inequalities)
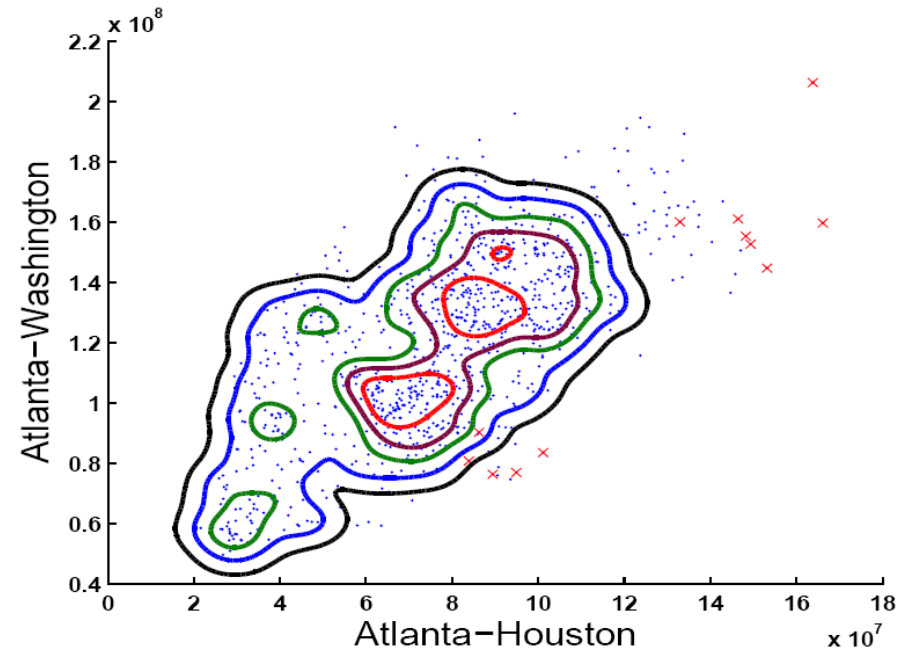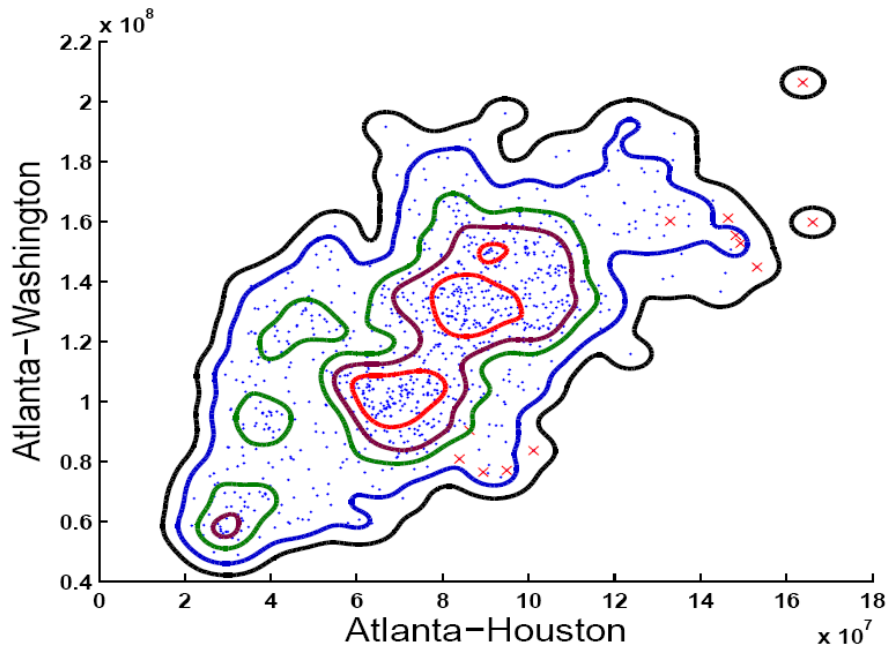
# Outline

Topics
- Robust KDEs
- Sparse KDEs
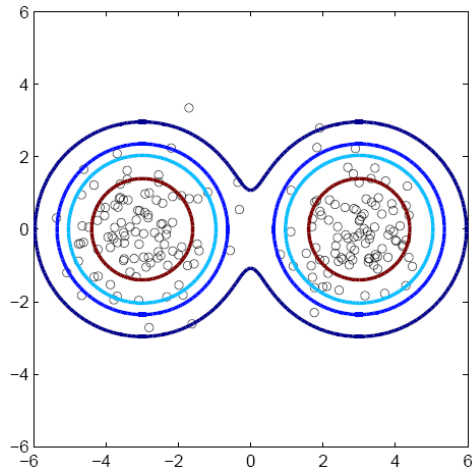- Consistency with fixed bandwidth

Themes
- KDE = mean in a function space
- Weighted KDEs to achieve above goals
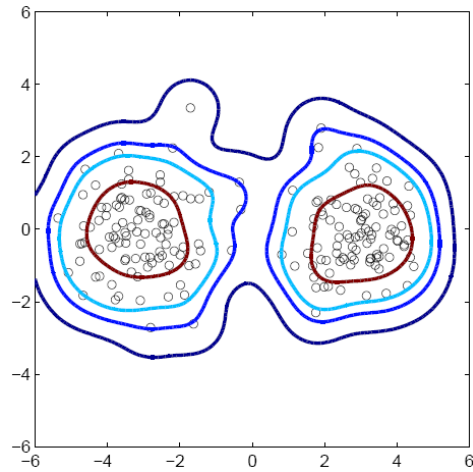
# Robust KDE for Anomaly Detection
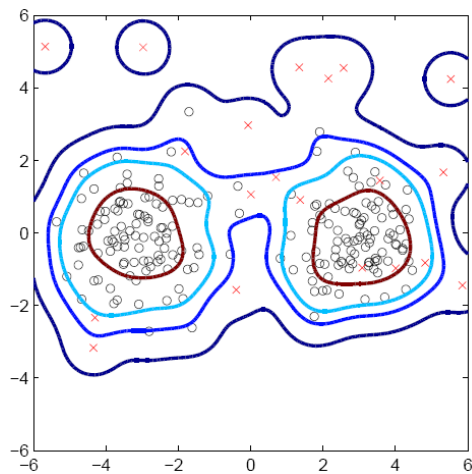
Abilene network traffic volumes
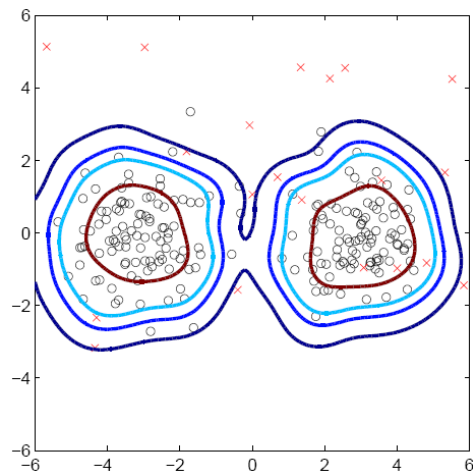
# Synthetic Example



(a) True density

(b) KDE without outliers

(c) KDE with outliers

(d) RKDE with outliers

# Problem Statement

- $\mathbf{X}_1, \ldots, \mathbf{X}_n \sim f(\mathbf{x}) = (1 - \epsilon)f_0(\mathbf{x}) + \epsilon f_1(\mathbf{x})$

- Tasks

  - Estimate $f_0(\mathbf{x})$
  - Estimate $\{\mathbf{x} : f_0(\mathbf{x}) > \lambda\}$

# Kernel Density Estimate

- 

$$\widehat{f}_{KDE}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} k_\sigma(\mathbf{x}, \mathbf{X}_i)$$

- Gaussian kernel

$$k_\sigma(\mathbf{x}, \mathbf{x}') = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^d \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

# Gaussian RKHS

- There exists a **Hilbert space** $\mathcal{H}_\sigma$ and a **feature map** $\Phi_\sigma : \mathbb{R}^d \to \mathcal{H}_\sigma$ such that

$$k_\sigma\left(\mathbf{x}, \mathbf{x}'\right) = \left\langle \Phi_\sigma(\mathbf{x}), \Phi_\sigma(\mathbf{x}') \right\rangle_{\mathcal{H}_\sigma}$$

- Canonical feature map

$$\Phi_\sigma(\mathbf{x}) = k_\sigma(\cdot, \mathbf{x})$$

- Reproducing property

$$\forall\, g \in \mathcal{H}_\sigma, \quad g(\mathbf{x}) = \left\langle \Phi_\sigma(\mathbf{x}), g \right\rangle_{\mathcal{H}_\sigma}$$

- $\|\Phi_\sigma(\mathbf{x})\|^2 = k_\sigma(\mathbf{x}, \mathbf{x}) = (\sqrt{2\pi}\sigma)^{-d}$
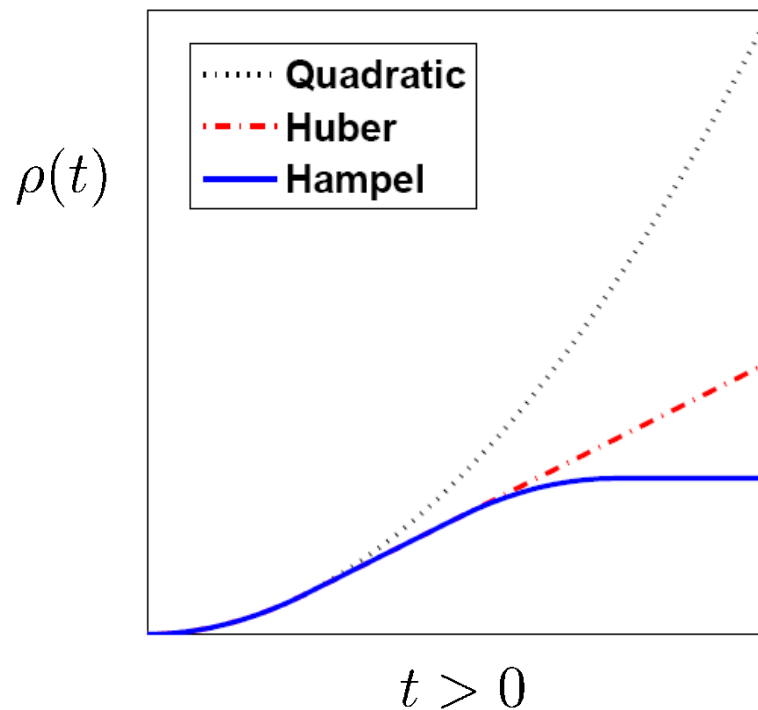
# KDE = mean in RKHS

●

$$\widehat{f}_{KDE} = \frac{1}{n} \sum_{i=1}^{n} k_\sigma\left(\cdot, \mathbf{X}_i\right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \Phi_\sigma(\mathbf{X}_i)$$

● Idea: Estimate this mean **robustly**

# Robust Kernel Density Estimate

$$\widehat{f}_{KDE} = \arg\min_{g \in \mathcal{H}_\sigma} \sum_{i=1}^{n} \|\Phi_\sigma(\mathbf{X}_i) - g\|_{\mathcal{H}_\sigma}^2$$

$$\widehat{f}_{RKDE} = \arg\min_{g \in \mathcal{H}_\sigma} \sum_{i=1}^{n} \rho\big(\|\Phi_\sigma(\mathbf{X}_i) - g\|_{\mathcal{H}_\sigma}\big)$$



$\rho(t)$

- Quadratic
- Huber
- Hampel

$t > 0$

# Representer Theorem

- Recall

$$\widehat{f}_{RKDE} = \arg\min_{g \in \mathcal{H}_\sigma} \sum_{i=1}^{n} \rho\big(\|\Phi_\sigma(\mathbf{X}_i) - g\|_{\mathcal{H}_\sigma}\big)$$

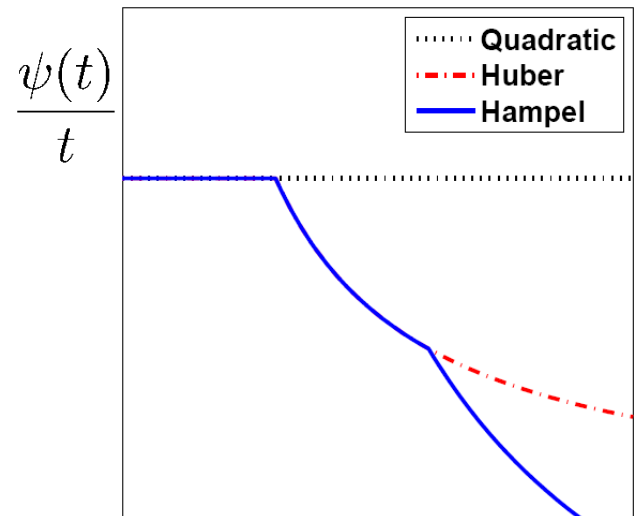- **Theorem**: If $\rho$ satisfies certain common assumptions, then

$$\widehat{f}_{RKDE}(\mathbf{x}) = \sum_{i=1}^{n} w_i k_\sigma(\mathbf{x}, \mathbf{X}_i)$$

for some $w_i \geq 0$, $\sum_{i=1}^{n} w_i = 1$.

- Furthermore

$$w_i \propto \frac{\psi(\|\Phi_\sigma(\mathbf{X}_i) - \widehat{f}_{RKDE}\|)}{\|\Phi_\sigma(\mathbf{X}_i) - \widehat{f}_{RKDE}\|}$$

where $\psi = \rho'$.

# Robustness Interpretation

- Notice that

$$\|\Phi_\sigma(\mathbf{x}) - \widehat{f}\|_{\mathcal{H}_\sigma}^2 = \langle \Phi_\sigma(\mathbf{x}) - \widehat{f}, \Phi_\sigma(\mathbf{x}) - \widehat{f} \rangle_{\mathcal{H}_\sigma}$$
$$= \|\Phi_\sigma(\mathbf{x})\|_{\mathcal{H}_\sigma}^2 - 2\langle \Phi_\sigma(\mathbf{x}), \widehat{f} \rangle_{\mathcal{H}_\sigma} + \|\widehat{f}\|_{\mathcal{H}_\sigma}^2$$
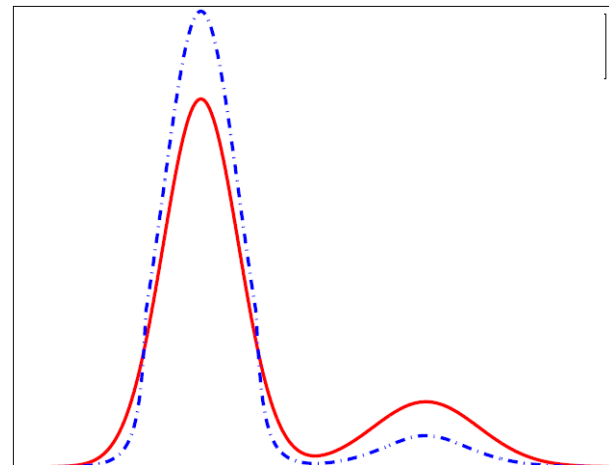$$= (\sqrt{2\pi}\sigma)^{-d} - 2\widehat{f}(\mathbf{x}) + \|\widehat{f}\|_{\mathcal{H}_\sigma}^2$$

- **Conclusion**:

$$w_i \text{ is small} \iff \|\Phi_\sigma(\mathbf{X}_i) - \widehat{f}_{RKDE}\| \text{ is large}$$
$$\iff \widehat{f}_{RKDE}(\mathbf{X}_i) \text{ is small}$$

- RKDE **down-weights** outlying points

# Other Results

- Efficient algorithm: iterative reweighted least squares (converges to global or local optimum depending on whether ρ is convex)

- Influence function: also reveals robustness to outliers relative to standard KDE

- Consistency
  - Convex ρ: converges to density of data (same as KDE)
  - Nonconvex ρ: converges to transformed version of density of data; this equals the uncontaminated density under certain assumptions on the contamination

- Experimental validation

# Sparse Approximation of Kernel Means

² Sparse kernel mean:

$$\sum_i \beta_i \Phi(\cdot, x_i) \approx \frac{1}{n} \sum_i \Phi(\cdot, x_i)$$

where $|f i : \beta_i \neq 0g| \leq k$

² Major complexity gains

± Evaluation of kernel density estimate: $O(n)$ ¡! $O(k)$

± Mean shift clustering: $O(n^2)$ ¡! $O(nk)$

² Set $z_i := \Phi(\cdot, x_i) \in V$

± $V = L^2(\mathbb{R}^d)$ for density estimation, or

± $V = \mathrm{RKHS}$ of $\Phi$, if $\Phi$ is a reproducing kernel

# Sparse Approximation of a Sample Mean

Setting

- inner product space $(V; \langle \cdot, \cdot \rangle)$

- $z_1, \ldots, z_n \in V$, $\quad \bar{z} = \frac{1}{n} \sum_i z_i$

- desired sparsity $k \ll n$

- $\alpha = (\alpha_1, \ldots, \alpha_n)$

Goal: (approximately) solve

$$\min_{\alpha \in \mathbb{R}^n} \left\| \bar{z} - \sum_i \alpha_i z_i \right\|_V$$

$$\text{s.t. } \|\alpha\|_0 \leq k$$

using an algorithm with $O(nk)$ time and space complexity

# Existing SA Methods Too Slow

Matching pursuit:

- Requires computing $\langle \hat{z}; z_i \rangle$ for each $i$

- Each $\langle \hat{z}; z_i \rangle = \frac{1}{n} \sum_j \langle \hat{z}_j; z_i \rangle$ requires $O(n)$ operations

- $\Rightarrow$ $O(n^2)$ complexity

# Problem Simplification

- Denote $I \subseteq [n] := \{1, 2, \ldots, n\}$

- Equivalent formulations:

$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2}\left\| z - \sum_i \beta_i z_i \right\|^2$$

$$\text{s.t. } \|\beta\|_0 \le k$$

$$\min_{\substack{I \subseteq [n] \\ |I| = k}} \min_{(\beta_i)_{i \in I}} \frac{1}{2}\left\| z - \sum_{i \in I} \beta_i z_i \right\|^2$$

$$\min_{\substack{I \subseteq [n] \\ |I| = k}} \frac{1}{2}\left\| z - z_I \right\|^2 :$$

- where $z_I$ is the projection of $z$ onto span$\{ z_i \mid i \in I \}$

# Incoherence-Based Bound

Theorem: Assume $\langle z_i, z_i \rangle = C$ for all $i$. For any $I \subseteq [n]$,

$$\|\hat{z}_i - z_I\| \le \left(1 - \frac{|I|}{n}\right)^r \sqrt{\frac{1}{C}(C^2 - \mu_I^2)} .$$

where

$$\mu_I := \min_{j \ge I} \max_{i \in I} \langle z_i, z_j \rangle$$

measures the incoherence of $\{z_i : i \in I\}$

# Bound Minimization

- For most kernels of interest,

$$\langle z_i, z_j \rangle = \langle \phi(\cdot, x_i), \phi(\cdot, x_j) \rangle$$
$$= g(\|x_i - x_j\|)$$

for some strictly decreasing $g$

- Example: Gaussian kernel, $V = $ RKHS

$$\langle \phi(\cdot, x_i), \phi(\cdot, x_j) \rangle = \phi(x_i, x_j)$$
$$= (2\pi\sigma^2)^{-d=2} \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

# Bound Minimization = k-center problem

² For most kernels of interest,

$$\langle z_i ; z_j \rangle = \langle \acute{A}(\cdot, x_i); \acute{A}(\cdot, x_j) \rangle$$
$$= g(\| x_i - x_j \|)$$

for some strictly decreasing g

² Bound minimization

$$I_k^{\bowtie} = \arg\max_{\substack{I \subseteq [n] \\ |I| = k}} \underset{l \in I}{\overset{\circ}{}}$$

$$= \arg\max_{\substack{I \subseteq [n] \\ |I| = k}} \min_{j \geq I} \max_{i \geq I} g(\| x_i - x_j \|)$$

$$= \arg\min_{\substack{I \subseteq [n] \\ |I| = k}} \max_{j \geq I} \min_{i \geq I} \| x_i - x_j \|$$
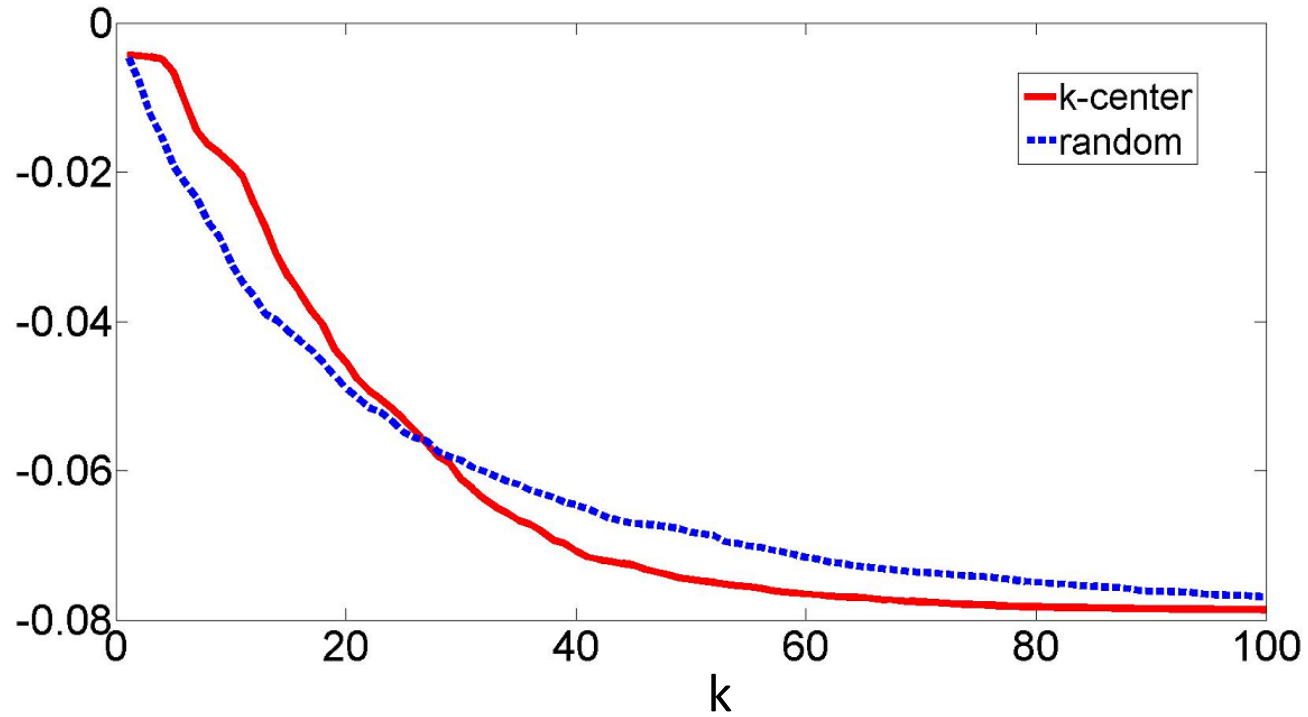
=) k-center problem

# k-center problem and algorithm

² Given n cities, place warehouses in k cities to minimize the maximum distance of a city to the nearest warehouse

² NP-Complete

² Greedy O(nk) 2-approximation algorithm



k = 6

# Example

# Weighted KDEs for Density Estimation

²  Estimate f via

$$\hat{f}_\circledR(x) = \sum_i \circledR_i \, k_{\frac{3}{4}}(x; x_i)$$

²  Theoretical question: Can we establish consistency (and rates) with ¾ ¯xed as n ! 0?

²  Empirical question: can we get better estimates in ¯nite sample settings?

# Estimation Method

$$\|\hat{f}_\rho - f\|_{L_2}^2 = \int (\hat{f}_\rho(x) - f(x))^2 dx$$

$$= \int \hat{f}_\rho(x)^2 dx - 2 \int \hat{f}_\rho(x) f(x) dx + \int f(x)^2 dx$$

First term:

$$\int \hat{f}_\rho(x)^2 dx = \sum_i \sum_j \rho_i \rho_j \int k_\sigma(x; x_i) k_\sigma(x; x_j) dx$$

$$= \rho^T Q \rho$$

Second term:

$$\int \hat{f}_\rho(x) f(x) dx = \sum_i \rho_i \left[ \int k_\sigma(x; x_i) f(x) dx \right]$$

$$\approx \sum_i \rho_i \left[ \frac{1}{n-1} \sum_{j \neq i} k_\sigma(x_j; x_i) dx \right]$$

$$= r^T \rho$$

# Fixed-Bandwidth KDE

Objective function:
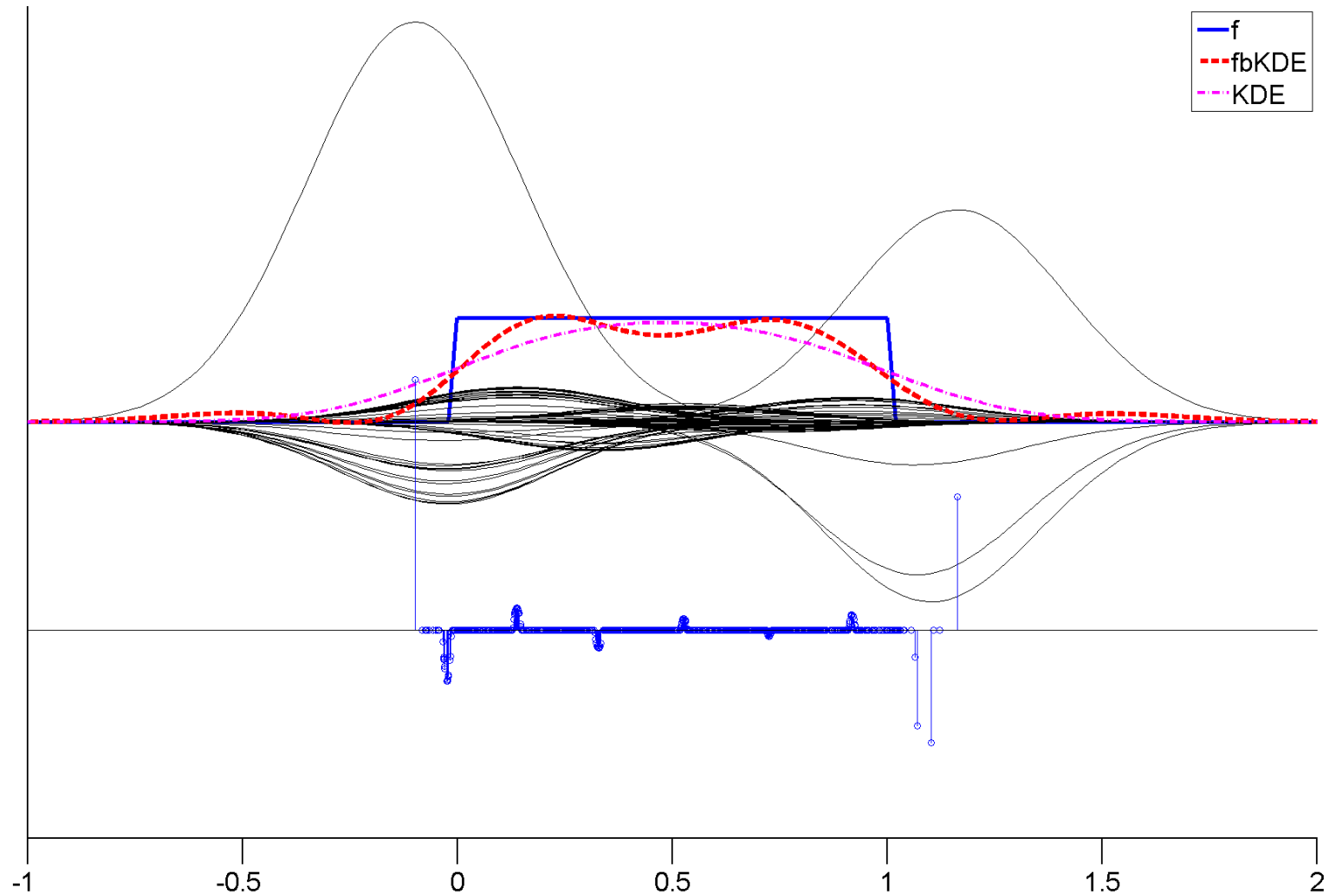
$$J(\alpha) = \alpha^T Q \alpha - 2\mathbf{r}^T \alpha$$

Optimization problem:

$$\widehat{\alpha} \longleftarrow \underset{\alpha \in \mathbb{R}^n}{\arg\min} \quad J(\alpha)$$

$$\text{s.t.} \quad \|\alpha\|_1 \leq R$$

Density estimator:

$$\widehat{f}(\mathbf{x}) = \sum_i \widehat{\alpha}_i k_\sigma(\mathbf{x}, \mathbf{x}_i)$$

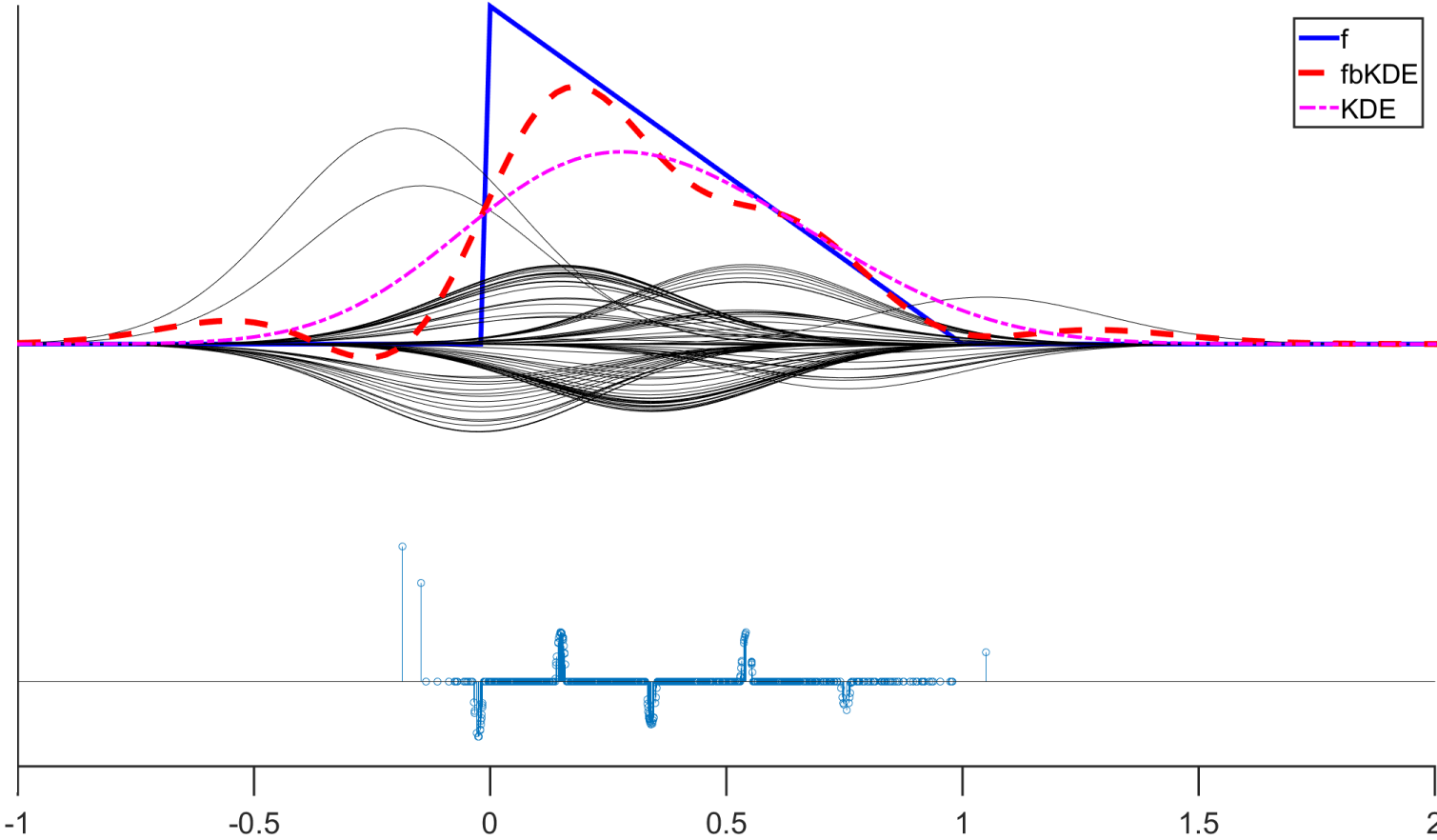# Example – Uniform Density

# Oracle Inequality

²  There is a constant $C$ such that with probability at least $1 - \delta$

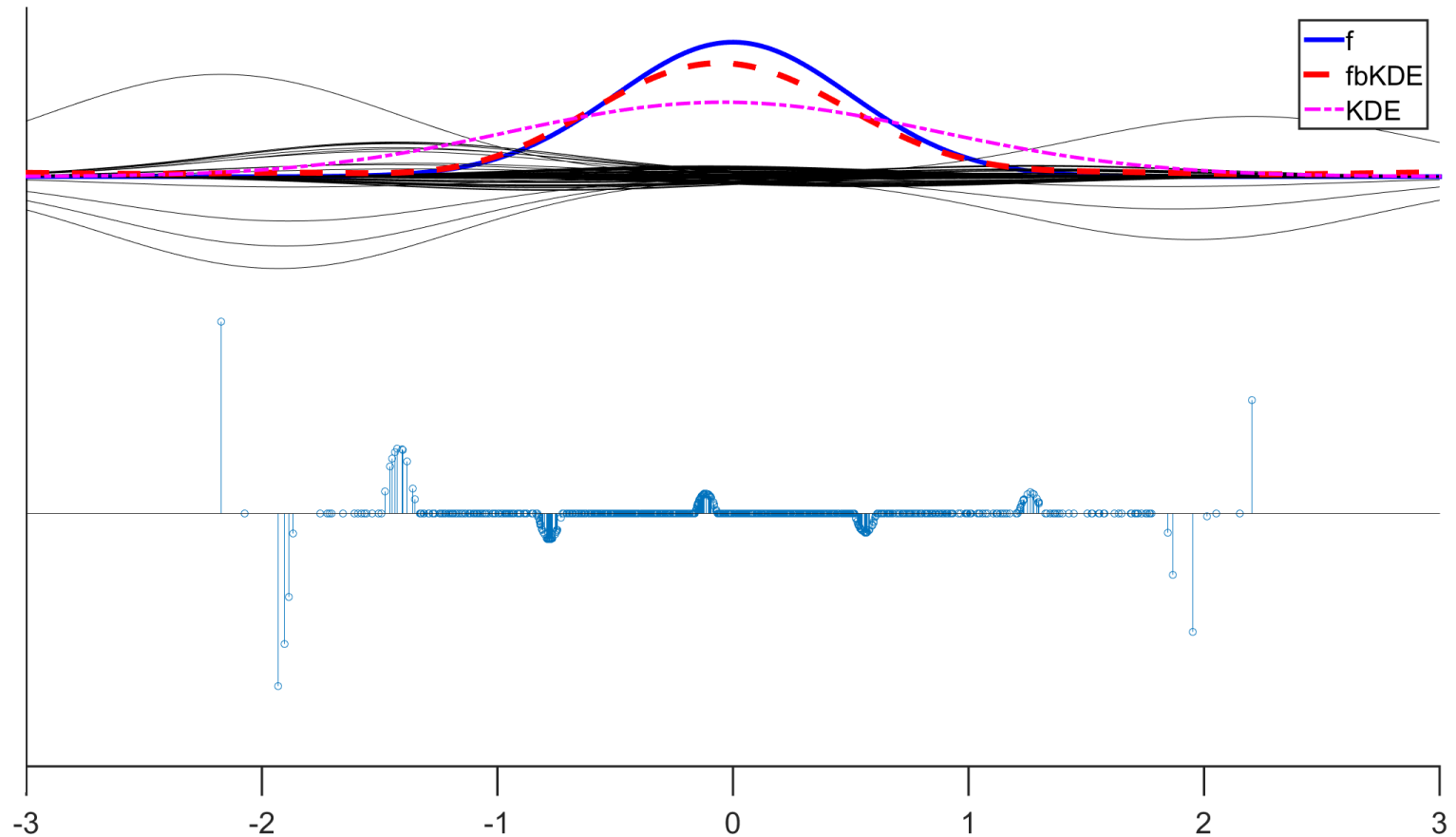$$\|\hat{f} - f\|_{L_2}^2 \cdot \inf_{\phi : \|\phi\|_{1} \cdot R} \|f_\phi - f\|_{L_2}^2 + C R \sqrt{\frac{\log(1/\delta)}{n}}$$

(ignoring $\log n$ terms).

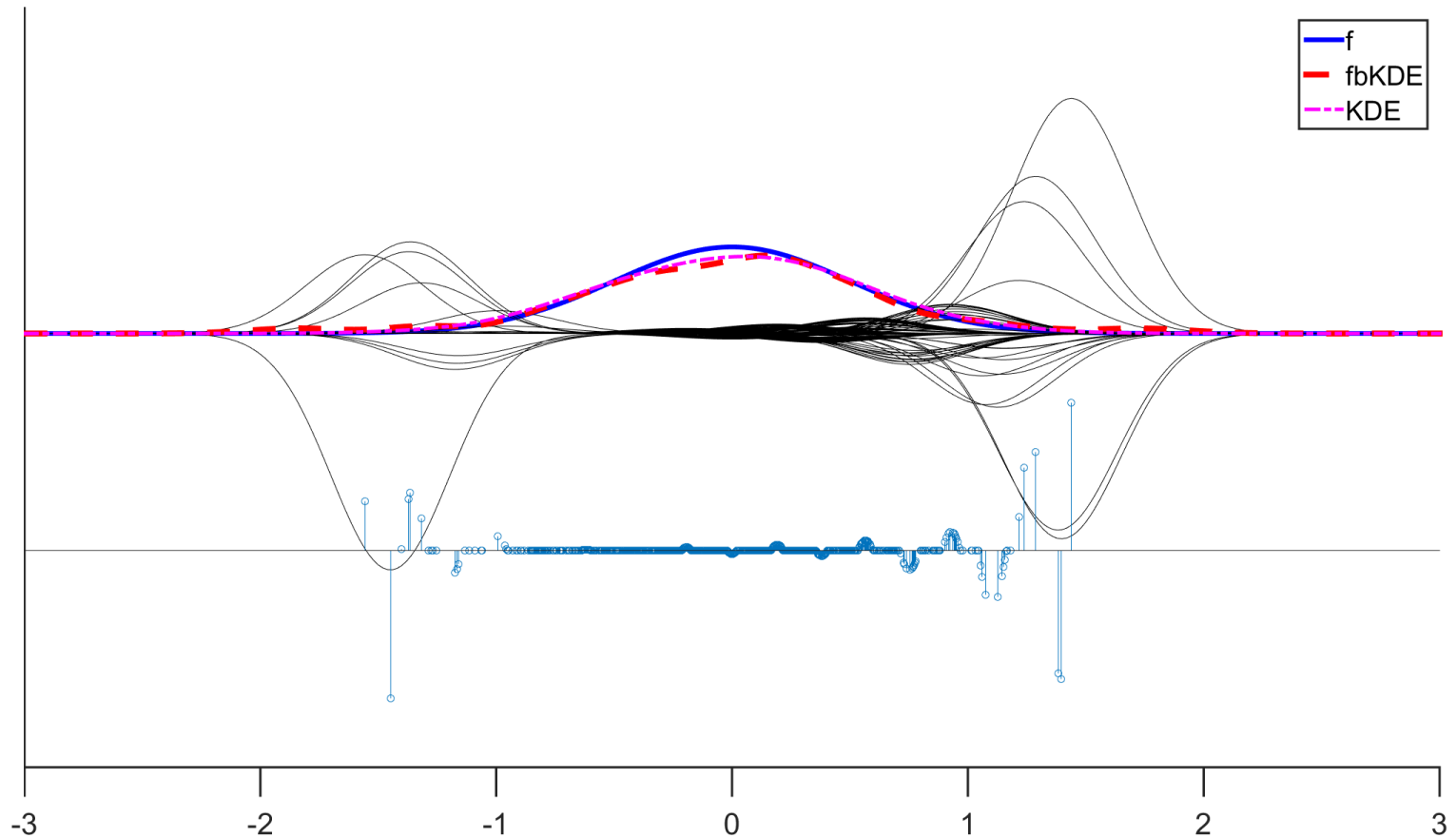²  Can be used to deduce consistency, rates of convergence with fixed $\lambda$.

# Triangular Density
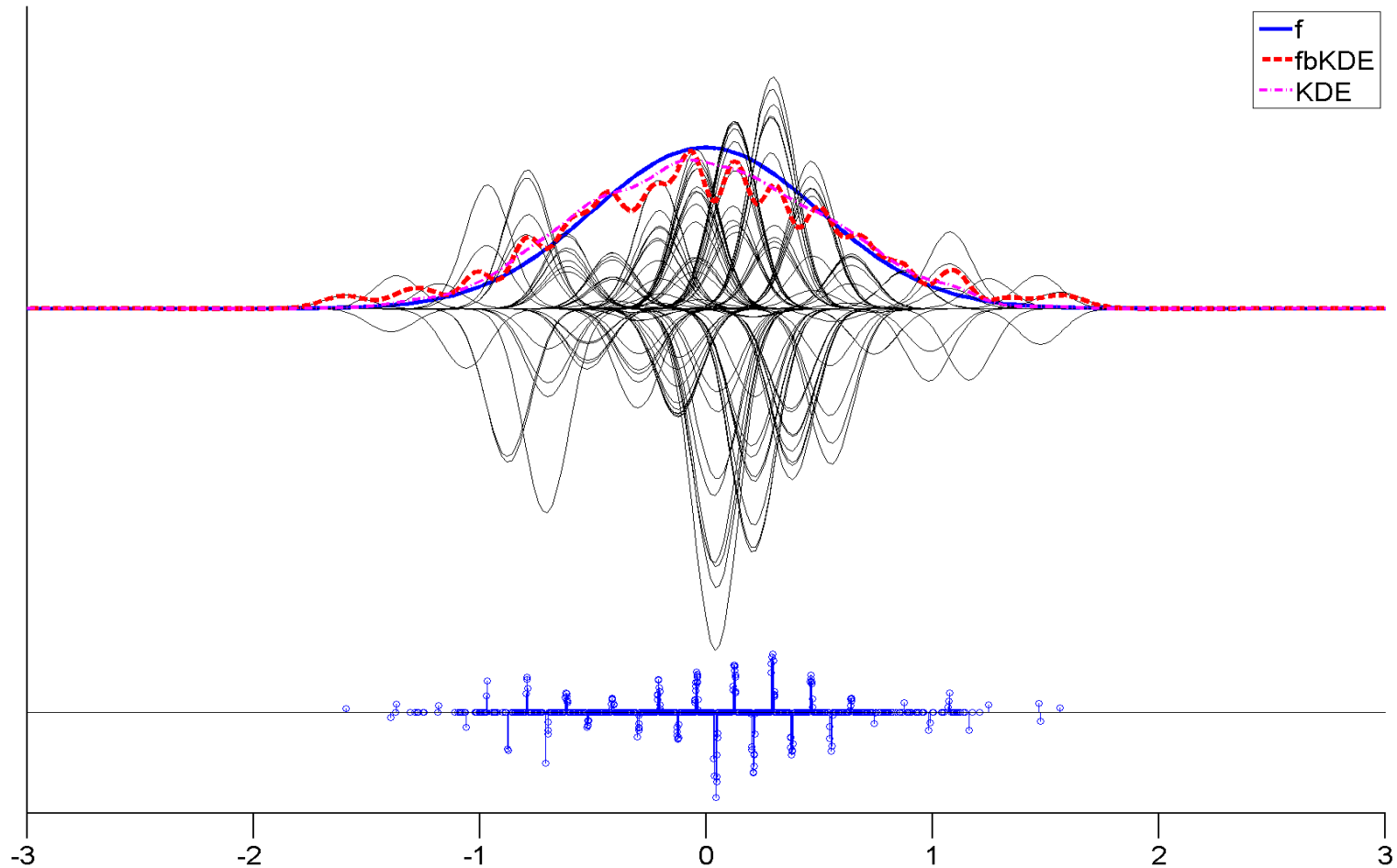


Legend:
- f
- fbKDE
- KDE

# Gaussian Density – Bandwidth Too Large

# Gaussian Density – Bandwidth Just Right

# Gaussian Density – Not Enough Regularization

# **Summary**

Topics
- Robust KDEs
- Sparse KDEs
- Consistency with fixed bandwidth

Themes
- KDE = mean in a function space
- Weighted KDEs to achieve above goals