



Robust Kernel Density Estimation

Clayton Scott

EECS and Statistics
University of Michigan

Problem Statement

- $\mathbf{X}_1, \dots, \mathbf{X}_n \sim f(\mathbf{x}) = (1 - \epsilon)f_{\text{norm}}(\mathbf{x}) + \epsilon f_{\text{anom}}(\mathbf{x})$
- Tasks
 - Estimate $f_{\text{norm}}(\mathbf{x})$
 - Estimate $\{\mathbf{x} : f_{\text{norm}}(\mathbf{x}) > \lambda\}$

Kernel Density Estimate

-

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

- 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}_σ and a **feature map** $\Phi_\sigma : \mathbb{R}^d \rightarrow \mathcal{H}_\sigma$ such that

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

- Canonical feature map

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

- Reproducing property

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

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

KDE = mean in RKHS

-

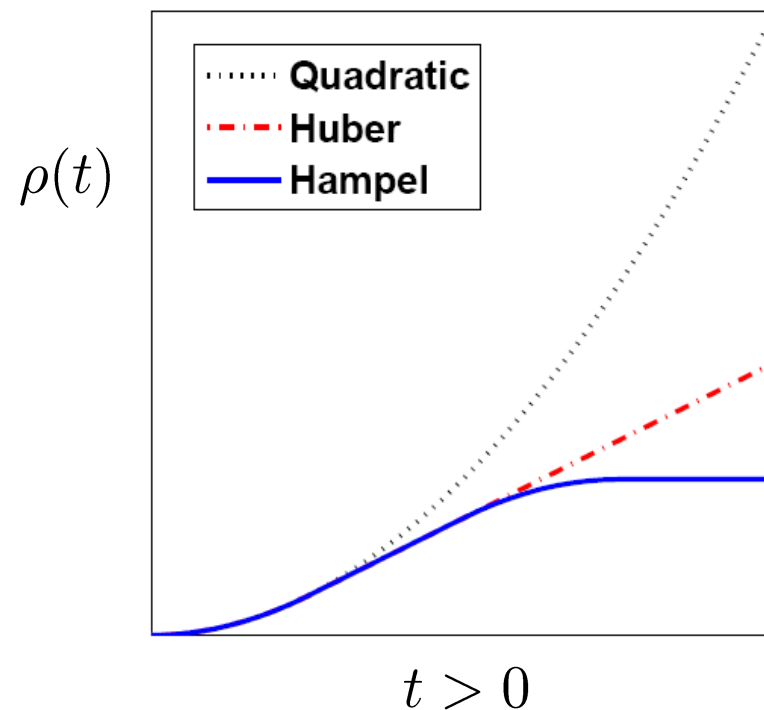
$$\begin{aligned}\hat{f}_{KDE} &= \frac{1}{n} \sum_{i=1}^n k_{\sigma}(\cdot, \mathbf{X}_i) \\ &= \frac{1}{n} \sum_{i=1}^n \Phi_{\sigma}(\mathbf{X}_i)\end{aligned}$$

- Idea: Estimate this mean **robustly**

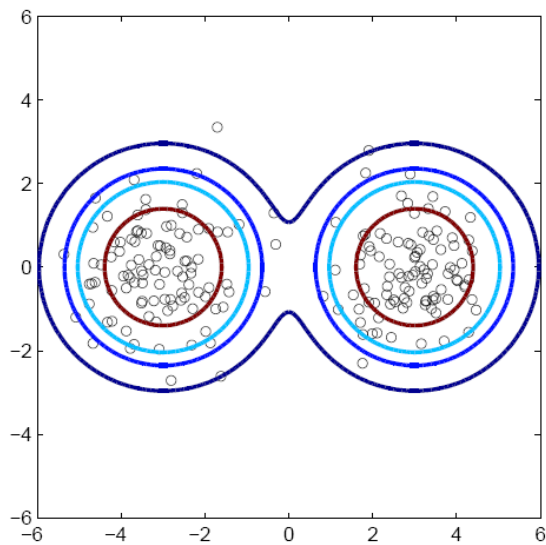
Robust Kernel Density Estimate

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

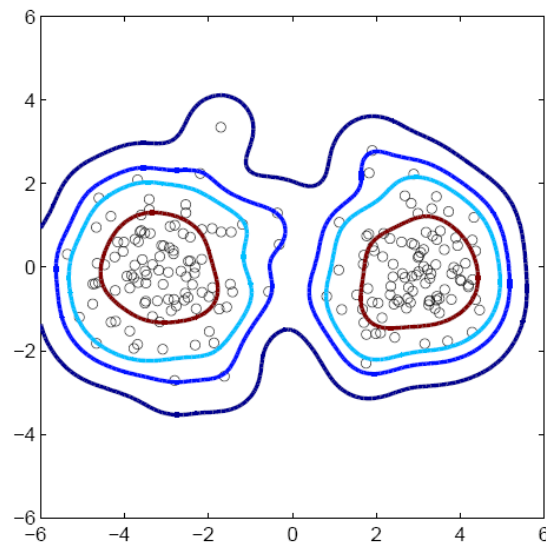
$$\hat{f}_{RKDE} = \arg \min_{g \in \mathcal{H}_\sigma} \sum_{i=1}^n \rho(\|\Phi_\sigma(\mathbf{X}_i) - g\|_{\mathcal{H}_\sigma})$$



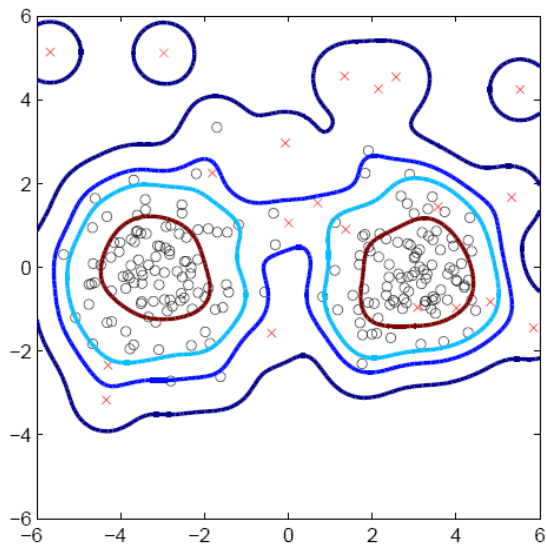
Example



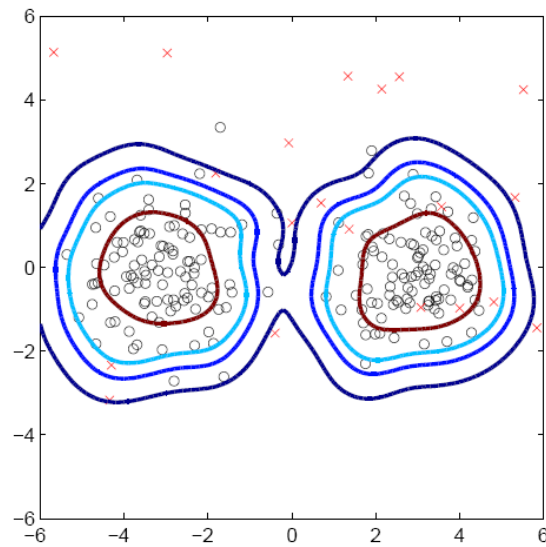
(a) True density



(b) KDE without outliers



(c) KDE with outliers



(d) RKDE with outliers

Outline

- Algorithm
- Influence Function
- Asymptotics
- Experiments: Anomaly Detection

Robust Multivariate Mean

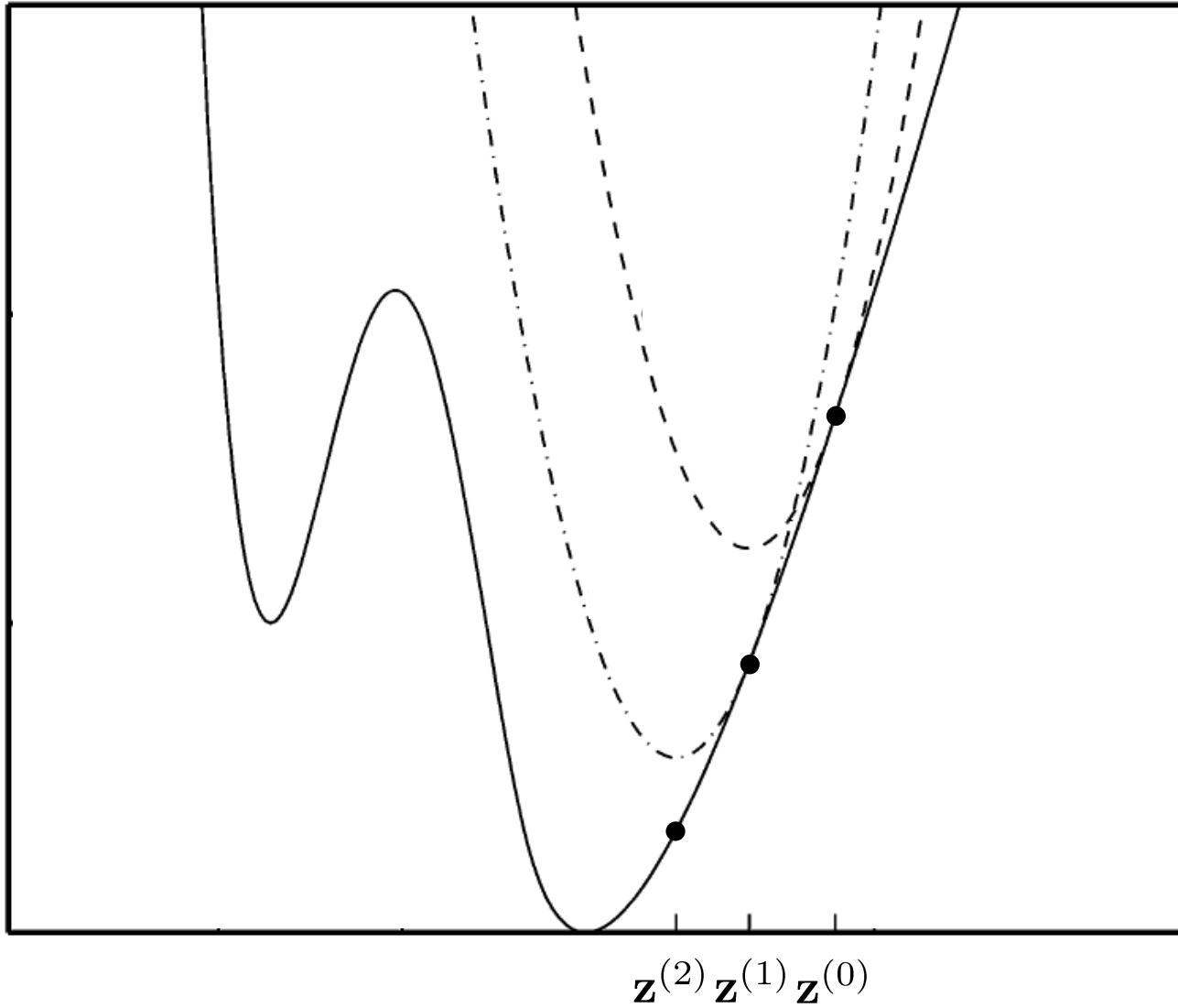
- $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^d$
- Sample mean

$$\bar{\mathbf{z}} = \arg \min_{\mathbf{z} \in \mathbb{R}^d} \sum_{i=1}^n \|\mathbf{z}_i - \mathbf{z}\|^2$$

- Robust sample mean

$$\bar{\mathbf{z}}_{\text{rob}} = \arg \min_{\mathbf{z} \in \mathbb{R}^d} \sum_{i=1}^n \rho(\|\mathbf{z}_i - \mathbf{z}\|)$$

Majorization / Minimization



Iterative Re-Weighted Least Squares

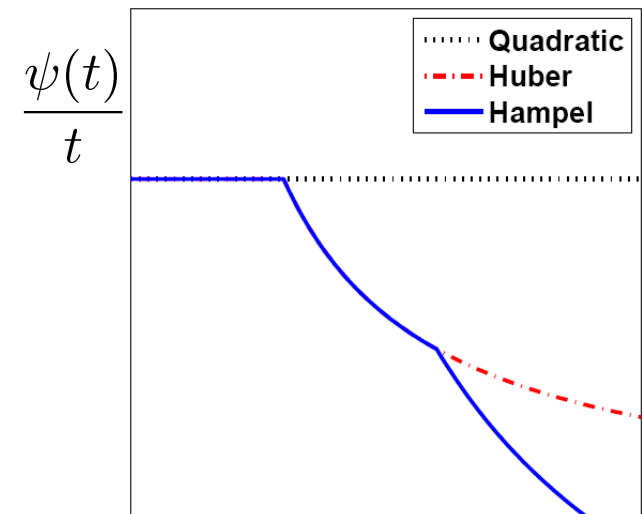
- Initialize

$$w_i^{(0)} = \frac{1}{n}$$

- Iterate

- $\bar{\mathbf{z}}_{\text{rob}}^{(k)} = \sum_{i=1}^n w_i^{(k-1)} \mathbf{z}_i,$
- $w_i^{(k)} \propto \frac{\psi(\|\mathbf{z}_i - \bar{\mathbf{z}}_{\text{rob}}^{(k)}\|)}{\|\mathbf{z}_i - \bar{\mathbf{z}}_{\text{rob}}^{(k)}\|}$
s.t. $\sum_{i=1}^n w_i^{(k)} = 1$

where $\psi = \rho'$



IRWLS Computation

-

$$\|\mathbf{z}_i - \bar{\mathbf{z}}_{\text{rob}}^{(k)}\|^2 = \|\mathbf{z}_i - \sum_{j=1}^n w_j^{(k-1)} \mathbf{z}_j\|^2$$

- Involves data only through inner products

$$\langle \mathbf{z}_i, \mathbf{z}_j \rangle$$

- Can **kernelize!** Substitute

$$\begin{aligned} \mathbf{z}_i &\rightarrow \Phi_\sigma(\mathbf{X}_i) \\ \langle \mathbf{z}_i, \mathbf{z}_j \rangle &\rightarrow \langle \Phi_\sigma(\mathbf{X}_i), \Phi_\sigma(\mathbf{X}_j) \rangle \\ &\quad \parallel \\ &\quad k_\sigma(\mathbf{X}_i, \mathbf{X}_j) \end{aligned}$$

Kernel IRWLS

- Initialize

$$w_i^{(0)} = \frac{1}{n}$$

- Iterate

- $f^{(k)} = \sum_{i=1}^n w_i^{(k-1)} \Phi_{\sigma}(\mathbf{X}_i)$
 $= \sum_{i=1}^n w_i^{(k-1)} k_{\sigma}(\cdot, \mathbf{X}_i)$
- $w_i^{(k)} \propto \frac{\psi(\|\Phi_{\sigma}(\mathbf{X}_i) - f^{(k)}\|)}{\|\Phi_{\sigma}(\mathbf{X}_i) - f^{(k)}\|}$
s.t. $\sum_{i=1}^n w_i^{(k)} = 1$

Kernel IRWLS Convergence

- **Theorem:** If ρ satisfies certain common assumptions, then $f^{(k)}$ converges to a stationary point.

If in addition ρ is convex and strictly increasing, then

$$f^{(k)} \rightarrow \hat{f}_{RKDE}$$

in \mathcal{H}_σ .

Representer Theorem

- Recall

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

- **Theorem:** If ρ satisfies certain common assumptions, then

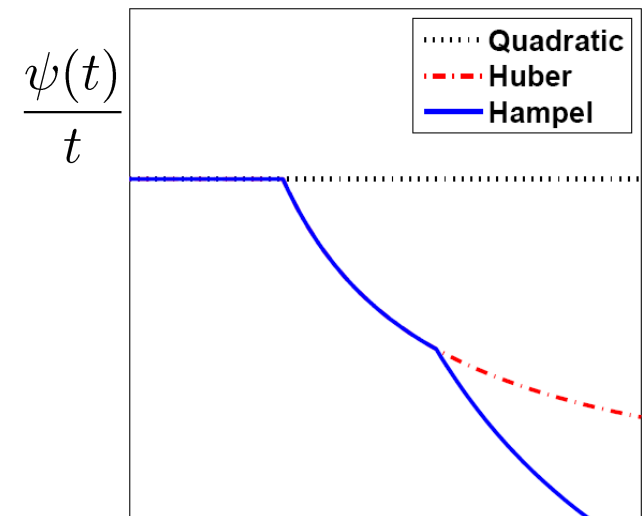
$$\hat{f}_{RKDE}(\mathbf{x}) = \sum_{i=1}^n w_i k_\sigma(\mathbf{x}, \mathbf{X}_i)$$

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

- Furthermore

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

where $\psi = \rho'$.



Robustness Interpretation # 1

- Notice that

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

- **Conclusion:**

$$\begin{aligned}w_i \text{ is small} &\iff \|\Phi_\sigma(\mathbf{X}_i) - \hat{f}_{RKDE}\| \text{ is large} \\ &\iff \hat{f}_{RKDE}(\mathbf{X}_i) \text{ is small}\end{aligned}$$

- **RKDE down-weights** outlying points

Connection to Data Depth

- Just argued that

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

- Therefore

$$\|\Phi_\sigma(\mathbf{x}) - \hat{f}\|_{\mathcal{H}_\sigma} < a \iff \hat{f}(\mathbf{x}) > \lambda$$

where $\lambda = \frac{1}{2}((\sqrt{2\pi}\sigma)^{-d} + \|\hat{f}\|_{\mathcal{H}_\sigma}^2 - a^2)$

- Likelihood depth w.r.t. $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$



“centroid” depth w.r.t. $\{\Phi_\sigma(\mathbf{X}_1), \dots, \Phi_\sigma(\mathbf{X}_n)\}$

Influence Function

- Scalar parameter estimator

$$\hat{\theta}(F)$$

- e.g., $\theta = \text{mean of } F$

$$\hat{\theta}(F) = \int x dF(x)$$

$$\hat{\theta}(F_n) = \frac{1}{n} \sum_{i=1}^n X_i$$

- Influence function

$$IF(x'; \hat{\theta}, F) = \lim_{s \rightarrow 0} \frac{\hat{\theta}((1-s)F + s\delta_{x'}) - \hat{\theta}(F)}{s}$$

Influence Function

- Density estimator

$$\hat{f}(\mathbf{x}; F)$$

- e.g., KDE

$$\hat{f}_{KDE}(\mathbf{x}; F) = \int k_{\sigma}(\mathbf{x}, \mathbf{y}) dF(\mathbf{y})$$

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

- Influence function

$$IF(\mathbf{x}, \mathbf{x}'; \hat{f}, F) = \lim_{s \rightarrow 0} \frac{\hat{f}(\mathbf{x}; (1-s)F + s\delta_{\mathbf{x}'}) - \hat{f}(\mathbf{x}; F)}{s}$$

Influence Function

- KDE

$$IF(\mathbf{x}, \mathbf{x}'; \hat{f}_{KDE}, F_n) = -\frac{1}{n} \sum_{i=1}^n k_{\sigma}(\mathbf{x}, \mathbf{X}_i) + k_{\sigma}(\mathbf{x}, \mathbf{x}')$$

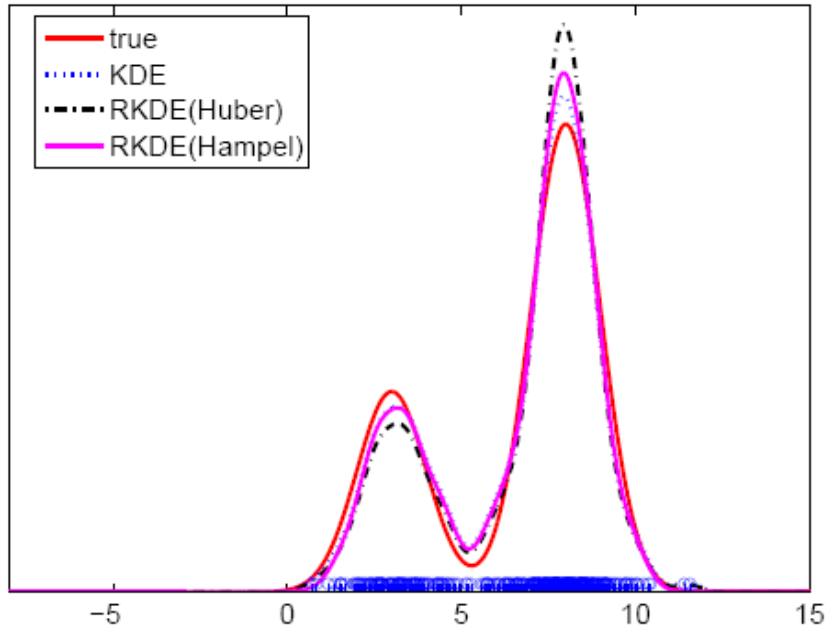
- RKDE: **Theorem:**

$$IF(\mathbf{x}, \mathbf{x}'; \hat{f}_{RKDE}, F_n) = \sum_{i=1}^n \alpha_i k_{\sigma}(\mathbf{x}, \mathbf{X}_i) + \alpha' k_{\sigma}(\mathbf{x}, \mathbf{x}')$$

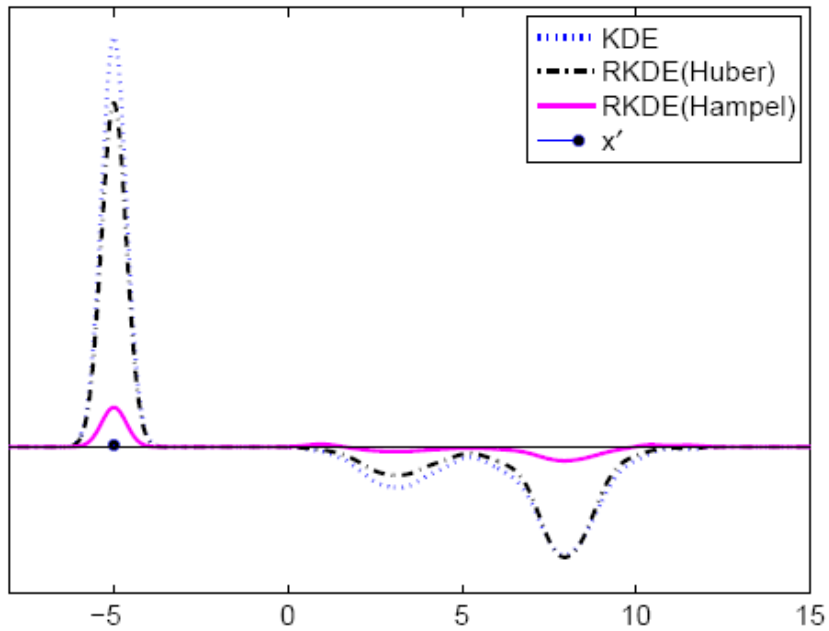
where $(\alpha_1, \dots, \alpha_n, \alpha')$ is the solution of a system of linear equations.

Example

$$\hat{f}(\mathbf{x}; F_n)$$



$$IF(\mathbf{x}, \mathbf{x}'; \hat{f}, F_n)$$




Robustness Interpretation # 2

- Exact formula

$$\alpha' = \frac{\frac{\psi(\|\Phi_\sigma(\mathbf{x}') - \hat{f}_{RKDE}\|)}{\|\Phi_\sigma(\mathbf{x}') - \hat{f}_{RKDE}\|}}{\frac{1}{n} \sum_{i=1}^n \frac{\psi(\|\Phi_\sigma(\mathbf{X}_i) - \hat{f}_{RKDE}\|)}{\|\Phi_\sigma(\mathbf{X}_i) - \hat{f}_{RKDE}\|}}$$

Measure of
outlyingness



- KDE

$$\alpha' = 1$$

- RKDE

$$\alpha' < 1$$

if \mathbf{x}' is an outlier

Asymptotics: σ fixed

$$f_\sigma^n = \arg \min_{g \in \mathcal{H}_\sigma} \underbrace{\frac{1}{n} \sum_{i=1}^n \rho(\|\Phi_\sigma(\mathbf{X}_i) - g\|_{\mathcal{H}_\sigma})}_{J_n(g)}$$

$$f_\sigma = \arg \min_{g \in \mathcal{H}_\sigma} \underbrace{\int \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) f(x) dx}_{J(g)}$$

Theorem: If f has compact support, then

$$J(f_\sigma^n) \xrightarrow{i.p.} J(f_\sigma) \quad \text{as } n \rightarrow \infty$$

Asymptotics: σ fixed

- **Corollary:** If f has compact support, and ρ is convex and strictly increasing, then

$$\|f_\sigma^n - f_\sigma\|_{\mathcal{H}_\sigma} \xrightarrow{i.p.} 0 \quad \text{as } n \rightarrow \infty$$

- **Corollary:** Under the same assumptions

$$\begin{aligned} \|f_\sigma^n - f_\sigma\|_\infty &= \sup_{\mathbf{x}} |\langle \Phi_\sigma(\mathbf{x}), f_\sigma^n - f_\sigma \rangle_{\mathcal{H}_\sigma}| \\ &\leq \underbrace{\|\Phi_\sigma(\mathbf{x})\|_{\mathcal{H}_\sigma}}_{\|(\sqrt{2\pi}\sigma)^{-d/2}} \|f_\sigma^n - f_\sigma\|_{\mathcal{H}_\sigma} \xrightarrow{i.p.} 0 \quad \text{as } n \rightarrow \infty \end{aligned}$$

Robustness Interpretation # 3

- Representer theorem (infinite sample):

$$f_\sigma = k_\sigma * p_\sigma$$

where p_σ is a density, $\text{supp}(p_\sigma) \subseteq \text{supp}(f)$

- Quadratic loss (KDE)

$$p_\sigma = f$$

- Robust loss (RKDE)

$$p_\sigma(\mathbf{x}) = \frac{w_\sigma(\mathbf{x})f(\mathbf{x})}{\int w_\sigma(\mathbf{y})f(\mathbf{y})d\mathbf{y}}$$

Tails are
down-weighted

where

$$w_\sigma(\mathbf{x}) = \frac{\psi(\|\Phi_\sigma(\mathbf{x}) - f_\sigma\|_{\mathcal{H}_\sigma})}{\|\Phi_\sigma(\mathbf{x}) - f_\sigma\|_{\mathcal{H}_\sigma}}$$

Asymptotics: $\sigma \rightarrow 0$

- Does

$$\lim_{\sigma \rightarrow 0} f_\sigma$$

exist?

- If $f_\sigma \rightarrow f_0$, and $\sigma = \sigma_n \rightarrow 0$, under what conditions does

$$f_{\sigma_n}^n \rightarrow f_0 ?$$

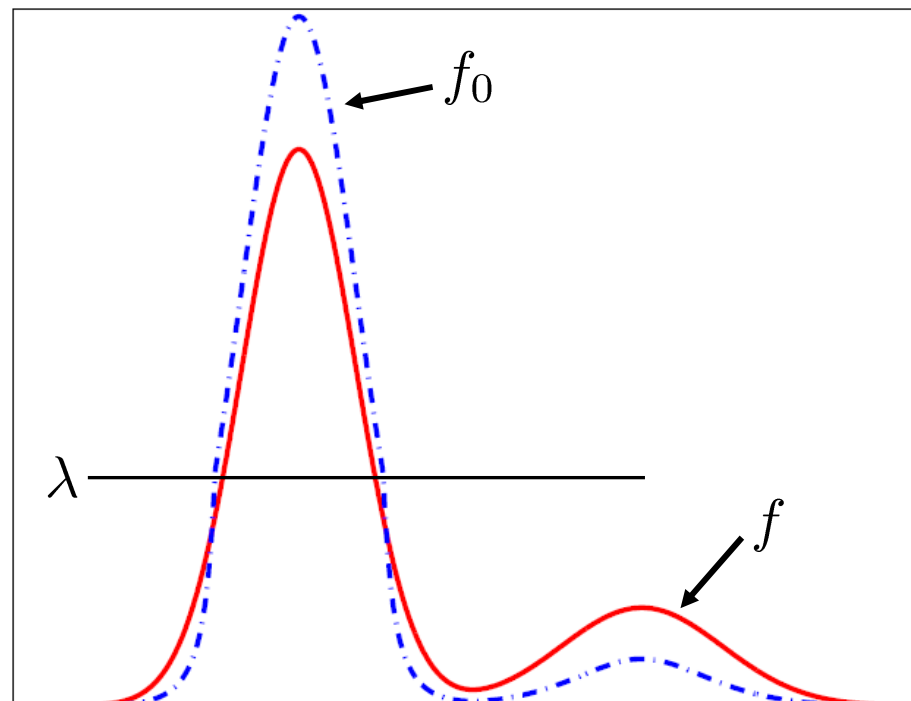
- How does f_0 depend on ρ ?

Hampel Loss

- Conjecture:

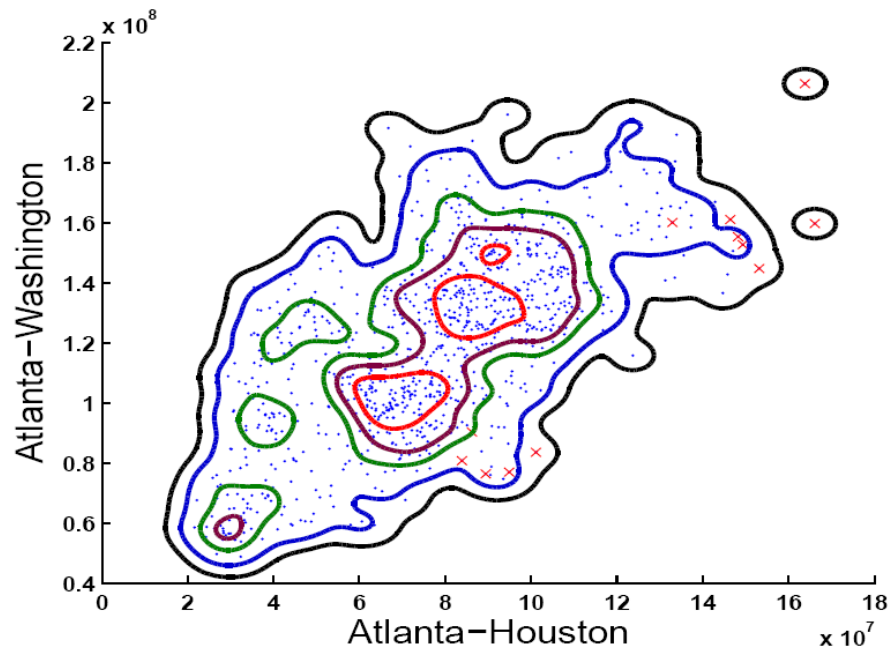
$$f_0(\mathbf{x}) = \begin{cases} af(\mathbf{x}), & f(\mathbf{x}) \geq \lambda \\ \frac{f(\mathbf{x})}{b-cf(\mathbf{x})}, & f(\mathbf{x}) < \lambda \end{cases}$$

for some $a > 1, b, c, \lambda > 0$

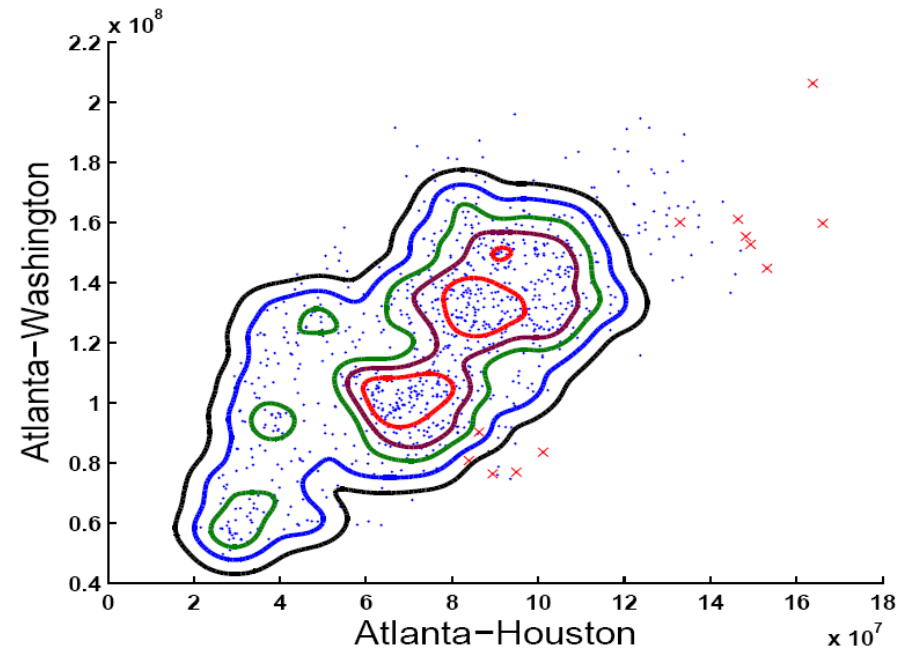


Example: Hampel Loss

KDE



RKDE



Anomaly Detection

- Training data

$$\mathbf{X}_1, \dots, \mathbf{X}_n \sim f(\mathbf{x}) = (1 - \epsilon)f_{\text{norm}}(\mathbf{x}) + \epsilon f_{\text{anom}}(\mathbf{x})$$

- Task: Classify normal examples versus anomalies
- Anomaly detector:

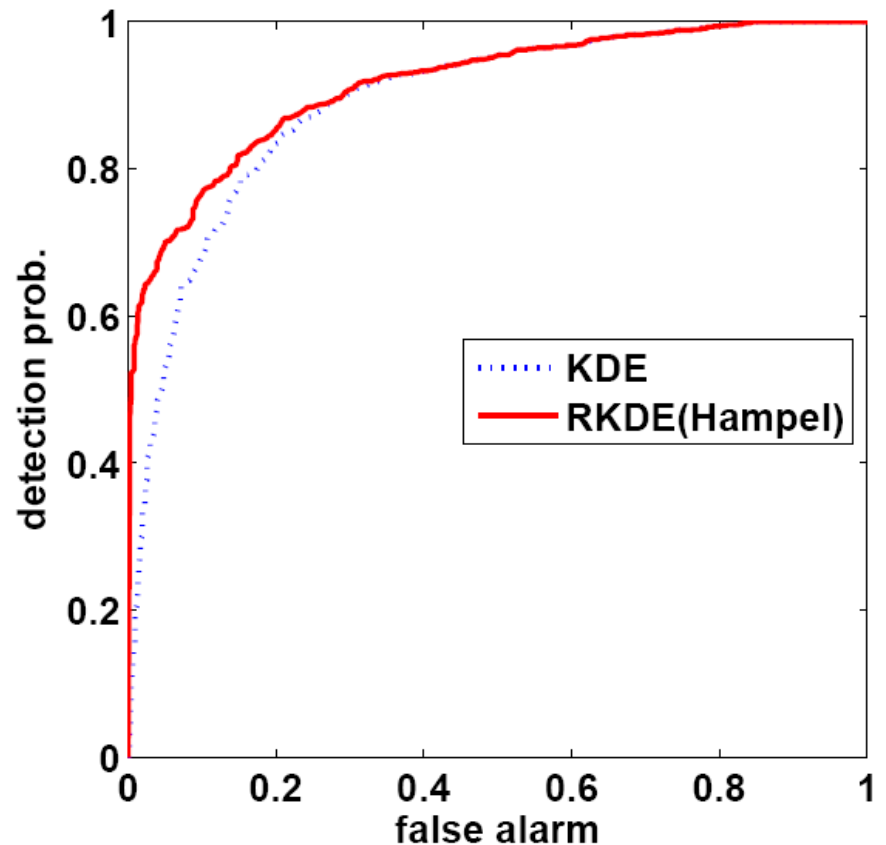
$$\{\mathbf{x} : \hat{f}(\mathbf{x}) > \lambda\}$$

for all $\lambda \geq 0$

- Performance:
Area under ROC (AUC)
- Testing data

$$\mathbf{X}'_1, \dots, \mathbf{X}'_{n'} \sim f_{\text{norm}}(\mathbf{x})$$

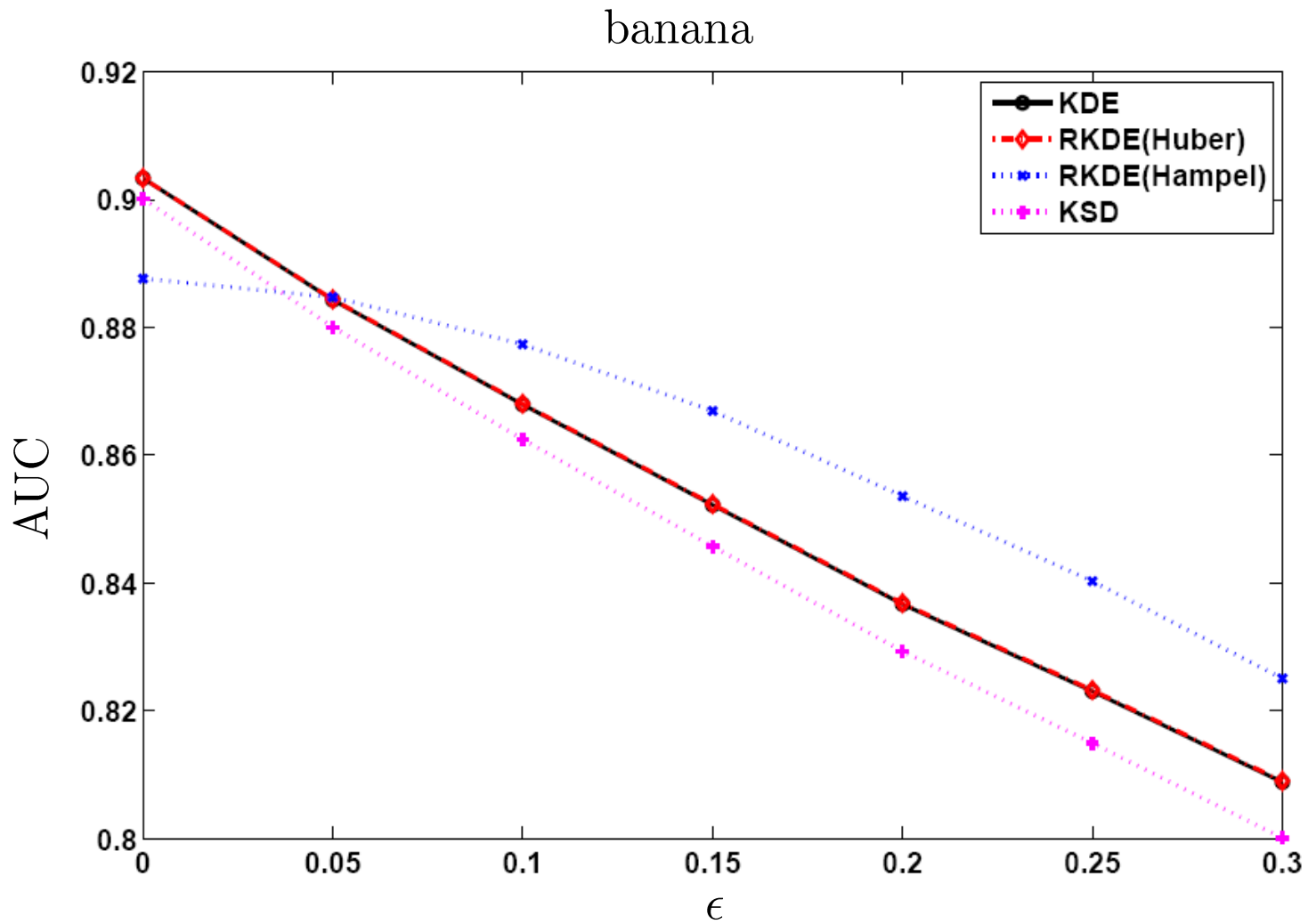
$$\mathbf{X}''_1, \dots, \mathbf{X}''_{n''} \sim f_{\text{anom}}(\mathbf{x})$$



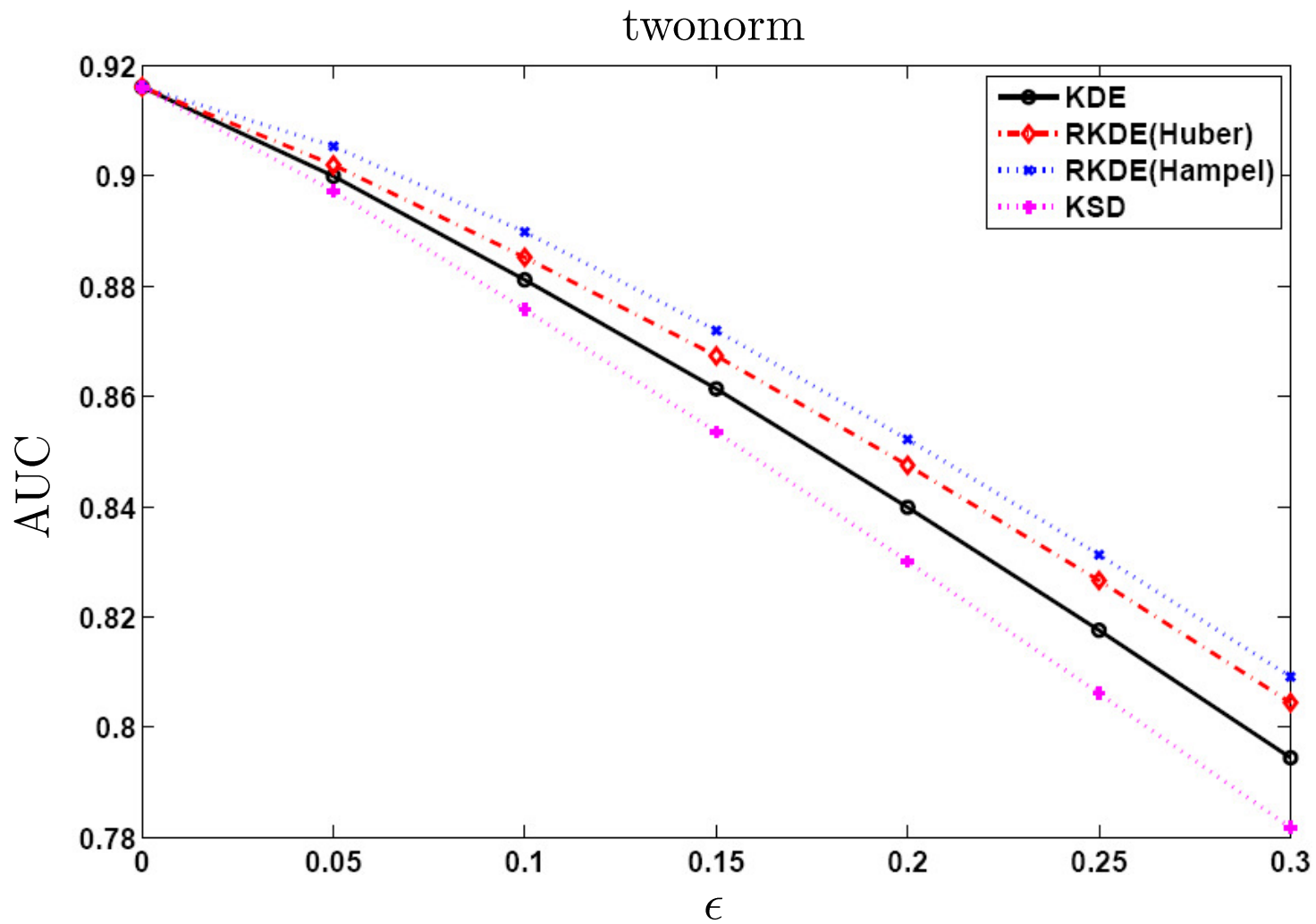
Anomaly Detection

- 15 benchmark datasets for binary classification
- $\epsilon = 0, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3$
- Algorithms compared
 - KDE
 - RKDE (Huber)
 - RKDE (Hampel)
 - Kernelized spatial depth (Chen et al., 2009)

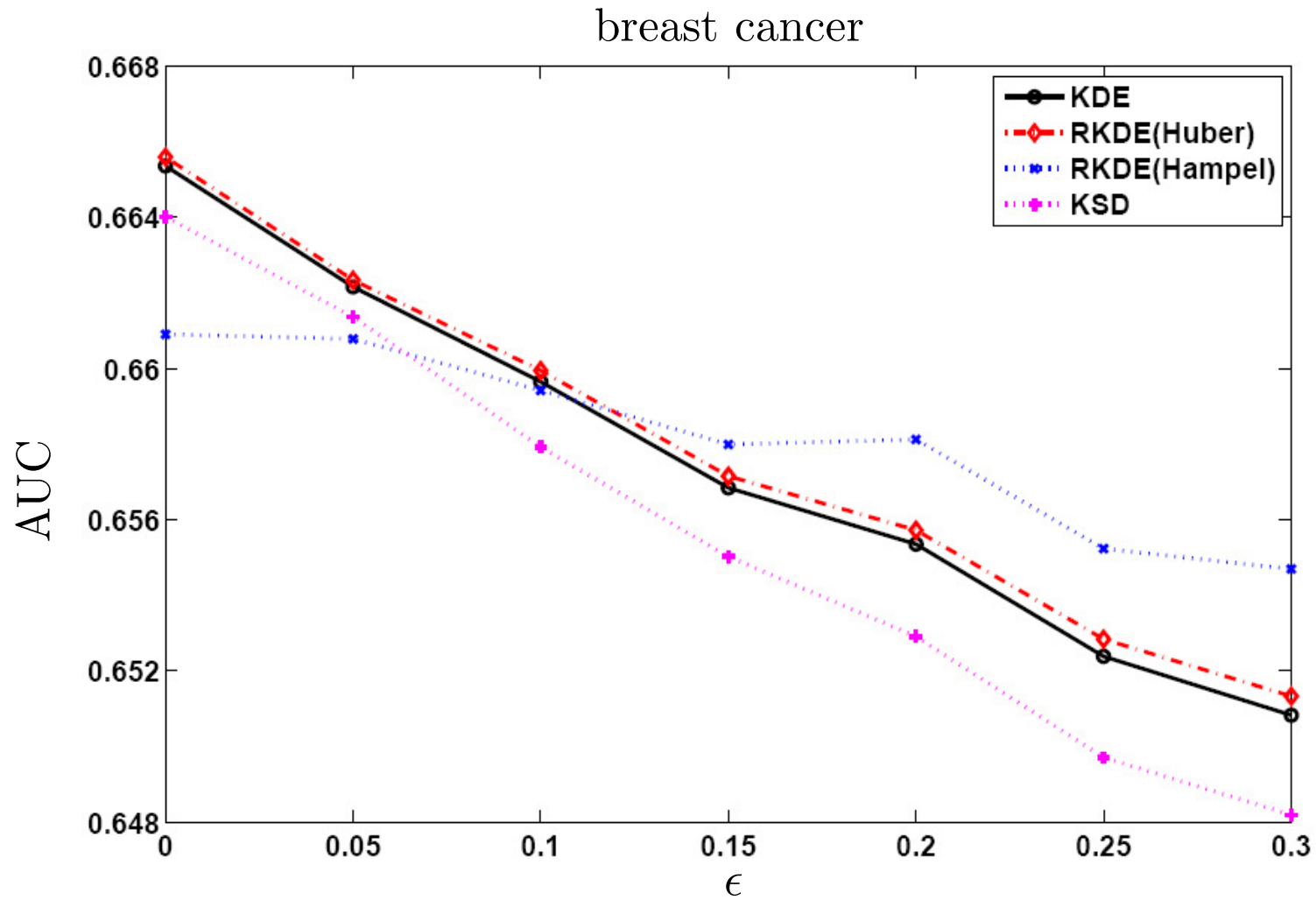
Anomaly Detection: AUC versus ϵ



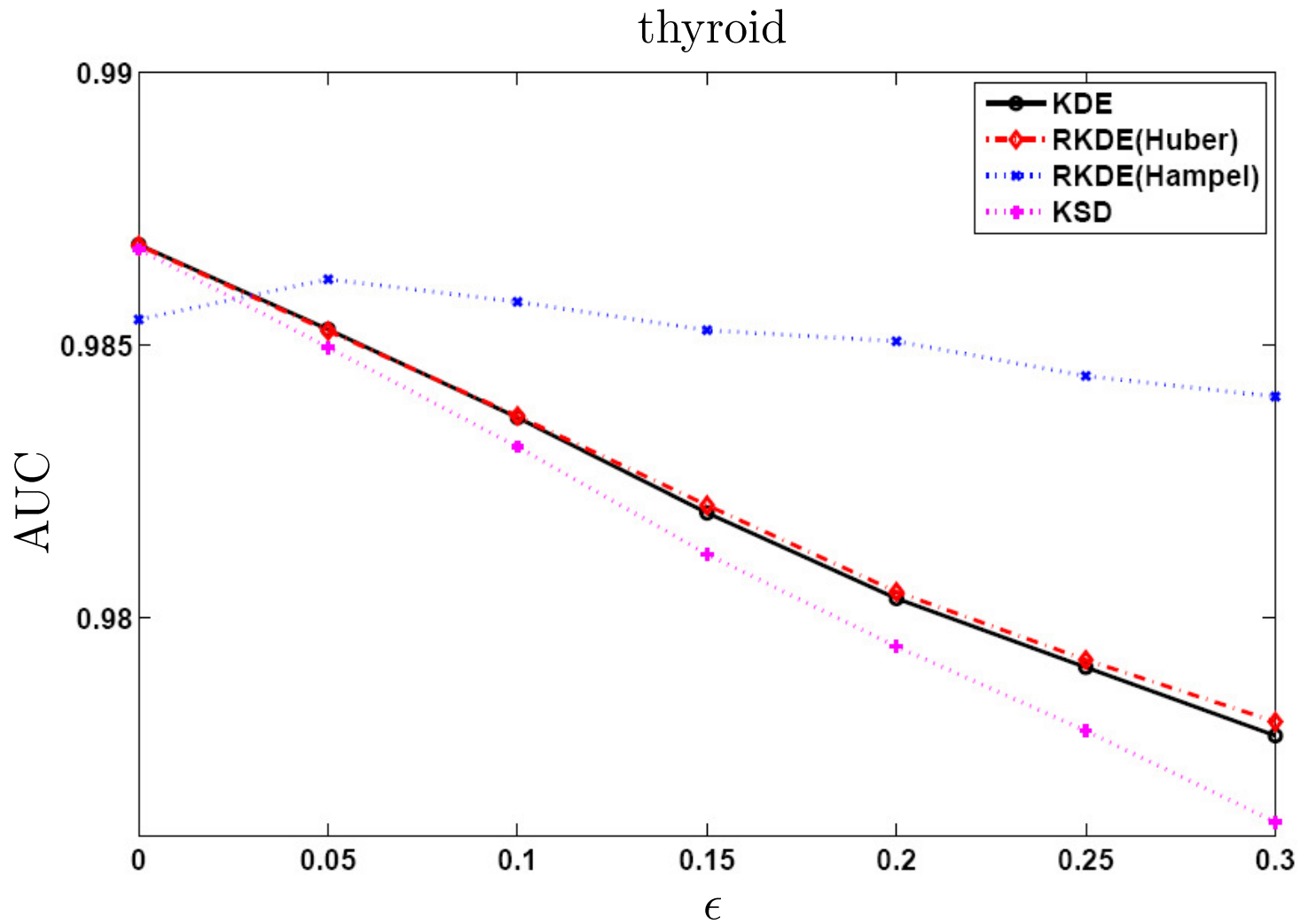
Anomaly Detection: AUC versus ϵ



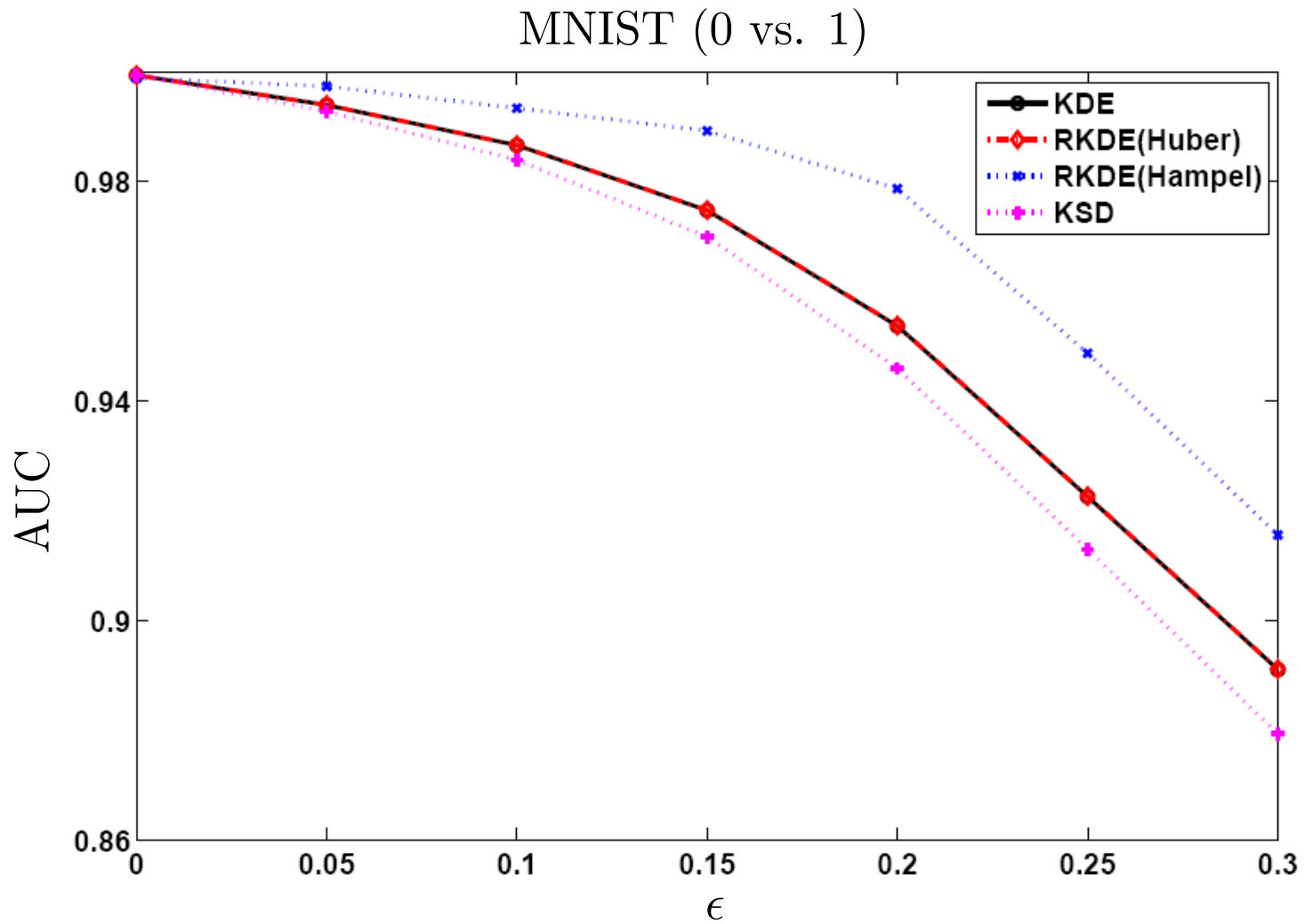
Anomaly Detection: AUC versus ϵ



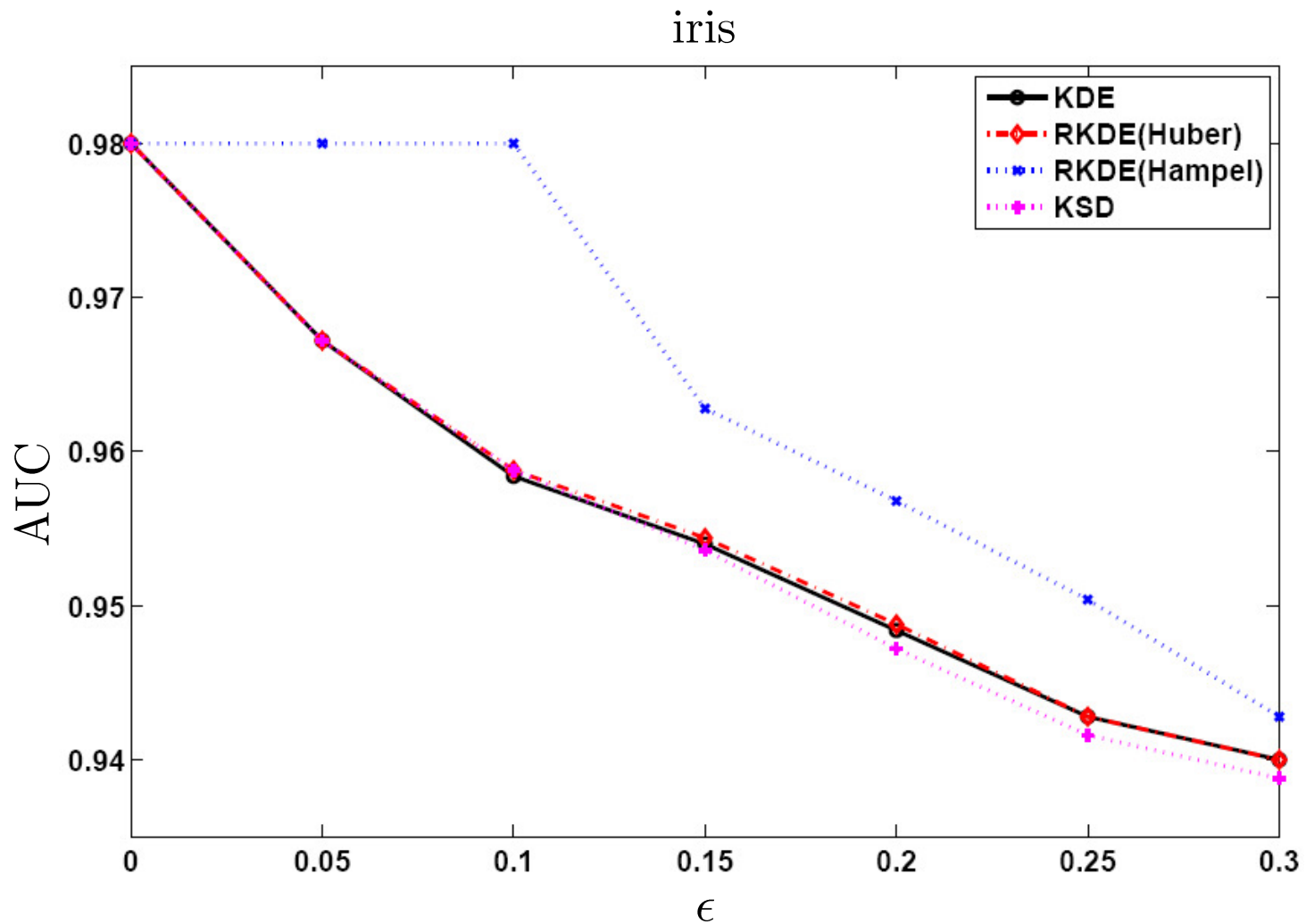
Anomaly Detection: AUC versus ϵ



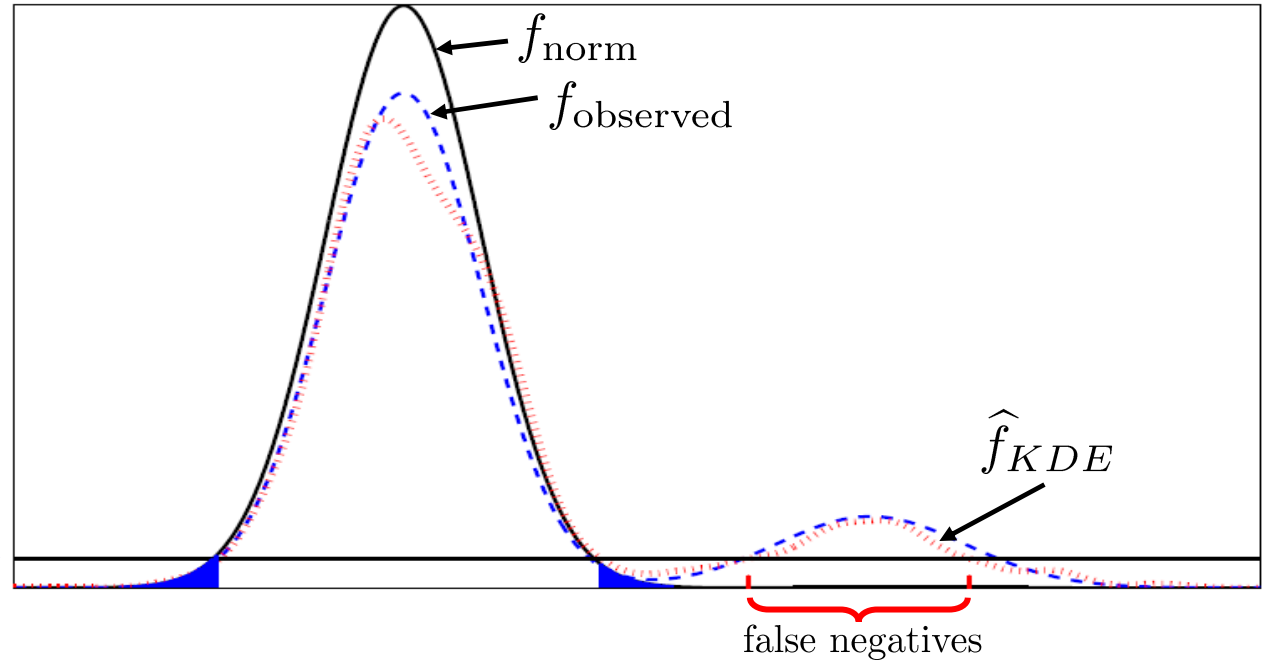
Anomaly Detection: AUC versus ϵ



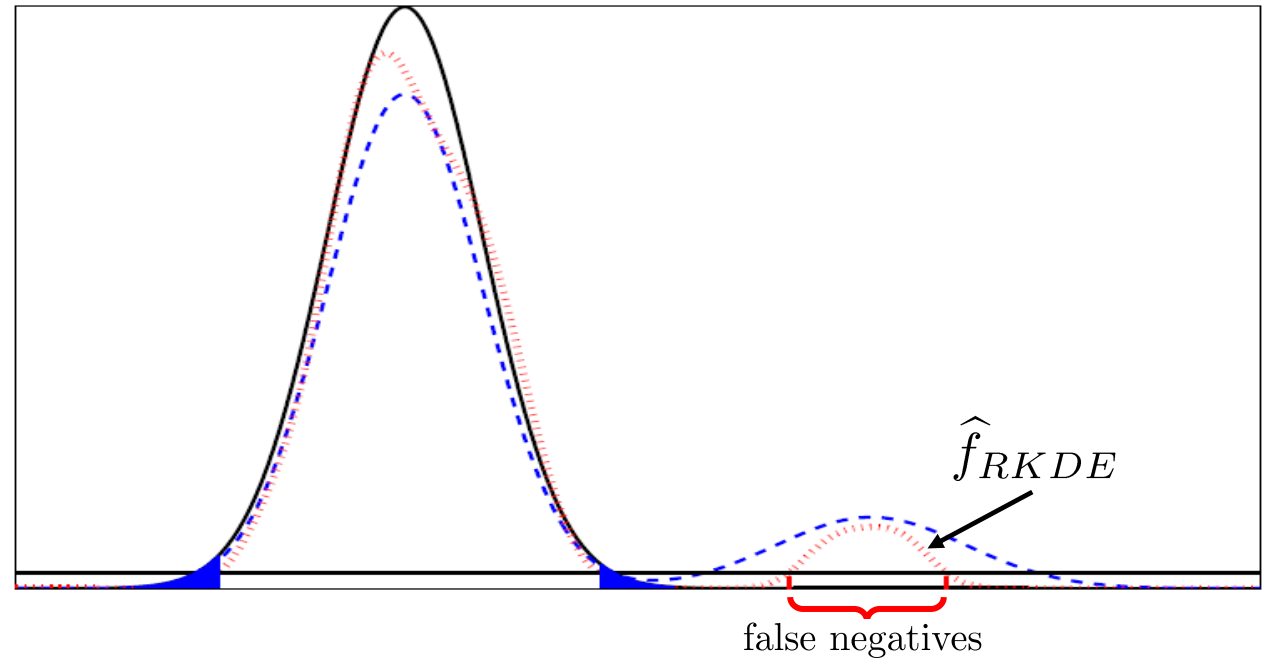
Anomaly Detection: AUC versus ϵ



Explanation



levels chosen to ensure same false alarm rate (shaded area)



Robustness Interpretation # 4:

- Define

$\alpha(\hat{f}) =$ smallest α such that some level set of \hat{f}
is a consistent estimate of the
 $(1 - \alpha)$ -**generalized quantile set** of f_{norm}

- **Conjecture:** If

- ϵ is sufficiently small
- f_{anom} and f_{norm} are sufficiently distinct

then

$$\alpha(\hat{f}_{RKDE}) < \alpha(\hat{f}_{KDE})$$

for the RKDE based on Hampel's loss

Conclusions

- Robust nonparametric density estimation is possible, if
 - ϵ is small
 - f_{anom} is outlying with respect to f_{norm}
- Easy implementation:
 - Training: Kernel IRWLS: $O(n^2)$ steps per iteration
 - Testing: $O(n)$ steps
- Interesting theory and works well in practice
- Data depth can be simple if it is kernelized
- Many open questions

Extensions / Open Questions

- Other kernels?
 - for anomaly detection, kernel need not integrate to 1
- Consistency, rates, limiting distribution?
- Re-analysis of KDE?
- Breakdown point?
- Other kernelizable depth functions / multivariate medians?

e.g., “spherical data depth,” Elmore et al., 2006.



Acknowledgments

- Joint work with JooSeuk Kim, Ph.D. candidate at the University of Michigan
- Supported by NSF Awards 0830490 and 0953135.

Parameter Tuning

- **Kernel bandwidth:** $\sigma =$ median distance to nearest neighbor (among training data)
- **Hampel loss parameters:**
 - Train RKDE with $\rho(t) = t$ to get \hat{f}_1
 - Set $d_i = \|\Phi_\sigma(\mathbf{X}_i) - \hat{f}_1\|$
 - Define a , b , and c to be the 50th, 95th, and 100th percentiles of $\{d_i\}$

Proof Idea

- Representer theorem (infinite sample):

$$f_\sigma = k_\sigma * p_\sigma$$

where p_σ is a density, $\text{supp}(p_\sigma) \subseteq \text{supp}(f)$

- Define

$$D_\sigma = \left\{ k_\sigma * p \mid \begin{array}{l} p \text{ is a density with} \\ \text{supp}(p) \subseteq \text{supp}(f) \end{array} \right\}$$

- Uniform error bound

$$P \left\{ \sup_{g \in D_\sigma} |J_n(g) - J(g)| > \epsilon \right\} \leq \underbrace{N(D_\sigma, \|\cdot\|_{\mathcal{H}_\sigma}, \epsilon)}_{\text{covering number}} \exp \left\{ \frac{-cn\epsilon^2}{\rho(\sigma^{-d/2})^2} \right\}$$

Anomaly Detection: Average Ranks

	epsilon						
	0	0.05	0.1	0.15	0.2	0.25	0.3
KDE	2.43	2.63	2.67	2.73	2.80	2.73	2.70
RKDE_Huber	2.10	2.33	2.13	2.20	2.13	2.33	2.23
RKDE_Hampel	3.10	2.13	1.87	1.67	1.67	1.47	1.60
KSD	2.37	2.90	3.33	3.40	3.40	3.47	3.47

CD at 0.05 1.21

CD at 0.10 1.08

Spatial Depth

- Spatial depth:

$$\text{outlyingness}(\mathbf{x}) \propto \left\| \underbrace{\sum_{i=1}^n \frac{\psi(\|\mathbf{x}_i - \mathbf{x}\|)}{\|\mathbf{x}_i - \mathbf{x}\|} (\mathbf{x}_i - \mathbf{x})}_{= 0 \text{ at multivariate median}} \right\|$$

with $\psi \equiv 1$

= 0 at multivariate median

- If ψ is decreasing (ρ nonconvex), then multiple stationary points exist (besides multivariate median), and “outlyingness” interpretation may be lost