

Exploring the Limits of Classification Accuracy

Carolyn Kim ¹ Lester Mackey ²

¹Computer Science Department, Stanford University

²Statistics Department, Stanford University

December 7, 2015

Classification

- Setup: random variable (X, Y) , where X describes the observations, and Y describes the class label
 - In our case, X takes values in \mathbb{R}^d (jet images), and Y takes values in ± 1 (“signal” W-jets or “background” QCD-jets).
- We can construct a classifier: $g : \mathbb{R}^d \rightarrow \{\pm 1\}$, with loss $L(g) := \mathbb{P}\{g(X) \neq Y\}$
- We want the optimal classifier (*Bayes Classifier*)

$$g^* = \operatorname{argmin}_{g: \mathbb{R}^d \rightarrow \{\pm 1\}} \mathbb{P}\{g(X) \neq Y\}$$

$$L^* := L(g^*)$$

- g^* is the classifier that outputs 1 if $\mathbb{P}\{1|x\} > \mathbb{P}\{-1|x\}$

k-Nearest Neighbors

The k-nearest neighbor classifier $g_{k,n}$ given n samples $(X_1, Y_1), \dots, (X_n, Y_n)$ with weights w_1, \dots, w_n is

$$g_{k,n}(x) = \begin{cases} 1 & \sum_{Y_i=1}^{X_i \in \{k\text{-nearest neighbors}(x)\}} w_i > \sum_{Y_i=-1}^{X_i \in \{k\text{-nearest neighbors}(x)\}} w_i \\ -1 & \text{otherwise} \end{cases}$$

Theorem (Universal Consistency of k-Nearest Neighbors, Devroye and Györfi, 1985, Zhao (1987))

For any distribution of (X, Y) , as $k \rightarrow \infty$, $k/n \rightarrow 0$, $n \rightarrow \infty$, i.i.d. samples, then $L(g_{k,n}) \rightarrow L^$.*

Theorem (Devroye, 1981)

For $k \geq 3$ and odd, $\lim_{n \rightarrow \infty} L(g_{1,n}) \leq L^(1 + \sqrt{\frac{2}{k}})$.*

Experimental setup

- Generate data: simulated signal and background events with $p_t \in [200, 400]$ GeV; each event is defined by a weight and 20-40 particles defined by $(\phi, \eta, \text{energy})$.
- Bin data, resulting in a jet image, a vector in \mathbb{R}^d .
- Optionally, whiten the data so the training covariance matrix is the identity.
- Compute the distances to the k -th nearest signal and background neighbors (this is enough information to do $2k - 1$ -nearest neighbor) in the “distance training set” (900K or 10M in size). In practice, this requires a lot of computational power!
- create a rejection versus efficiency curve

Step 1: Binning

Multiple possible binning strategies: equal size binning or equal weight (event weight only vs. event weight * energy bin bounds, energy only bin values vs energy density bin values)

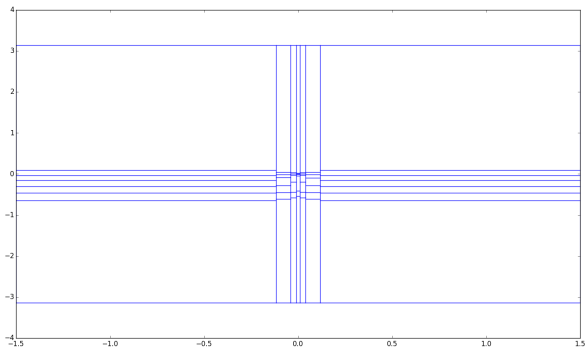
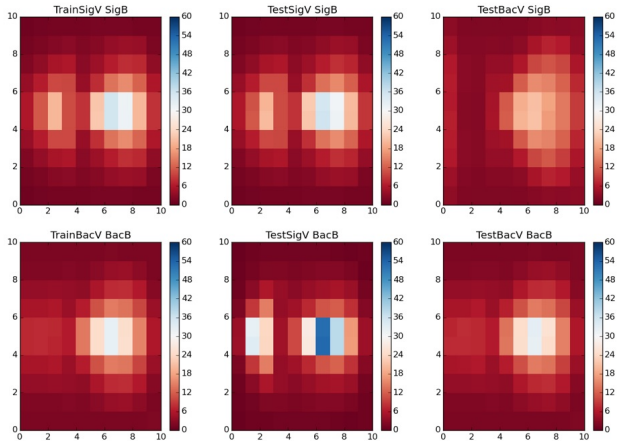


Figure 1: Sample bin bounds for an equal weighting scheme

Mean heatmap of one binning strategy

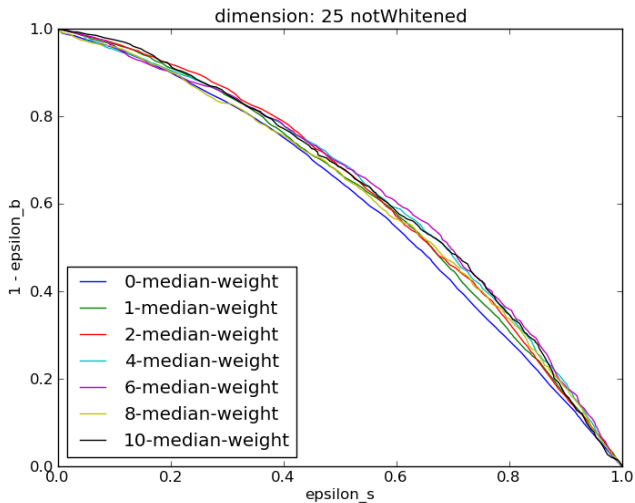
Binned Values (Energies) for dim 100



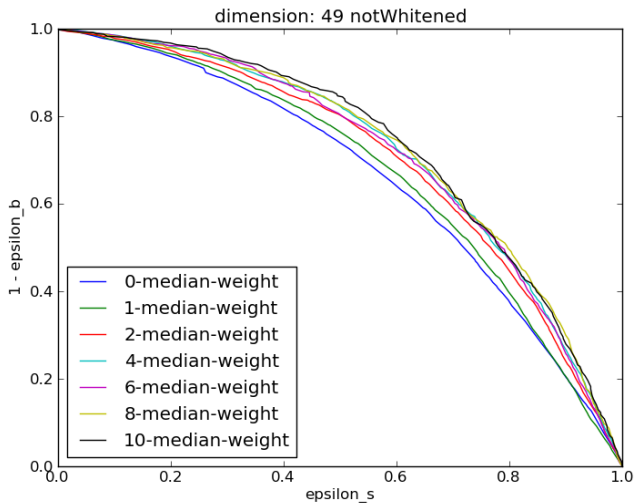
Plotting rejection versus efficiency curve

- x -axis is signal efficiency (proportion of signal classified as signal)
- y -axis is $1 - \text{background efficiency}$
- The 1-D discriminant is the ratio between the probability densities of distances to the k -th nearest signal and background neighbor. (2D-likelihood without taking the ratio has empirically not been better.)
- use one set of distances as a “curve training” set to estimate the densities, and another set of distances as the “curve testing” set to plot the curve

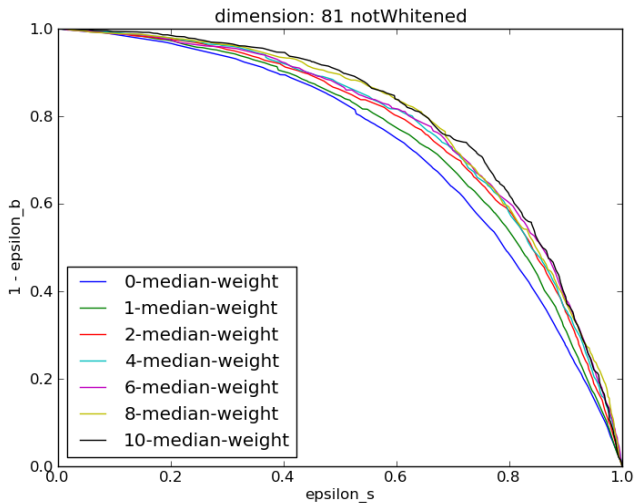
curve training, testing: 100K; distance training: 900K



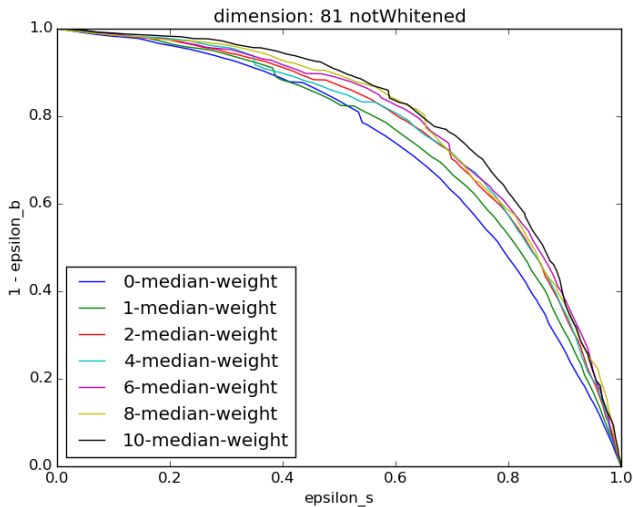
curve training, testing: 100K; distance training: 900K



curve training, testing: 100K; distance training: 900K

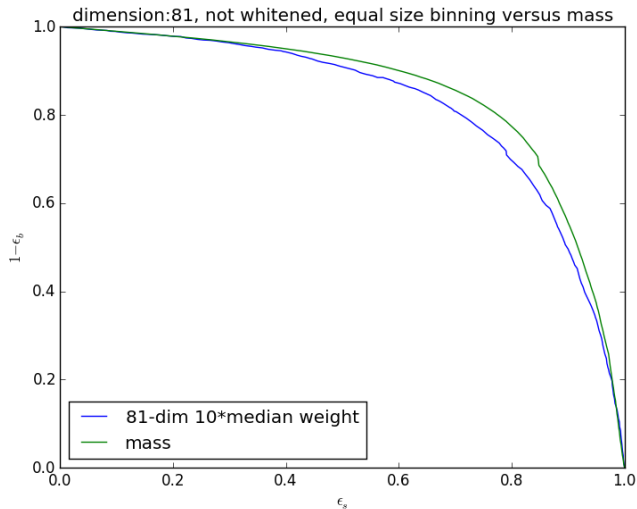


curve training, testing: 1M; distance training: 10M



curve training, testing: 1M; distance training: 10M

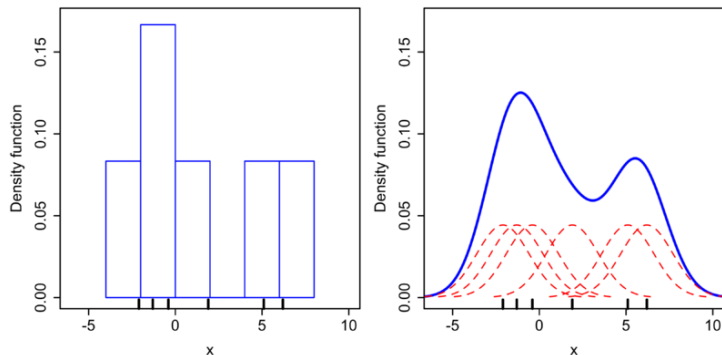
How well are we doing? Unfortunately, worse than mass...



Kernels

A kernel function $K : \mathbb{R}^d \rightarrow \mathbb{R}$ intuitively creates “bumps” around 0 (e.g. Gaussian kernel $K(x) = e^{-\|x\|^2}$).

We can estimate the probability density function by summing up kernel functions centered at the data points: $P(y_j|x) \propto \sum_{Y_i=y_j} w_i K(x - X_i)$



Credit: <http://en.wikipedia.org>

The kernel classifier $g_{K,n}$ for a kernel function K given n samples $(X_1, Y_1), \dots, (X_n, Y_n)$ with weights w_1, \dots, w_n is

$$g_{K,n}(x) = \begin{cases} 1 & \sum_{Y_i=1} w_i K\left(\frac{x-X_i}{h}\right) > \sum_{Y_i=-1} w_i K\left(\frac{x-X_i}{h}\right) \\ -1 & \text{otherwise} \end{cases}$$

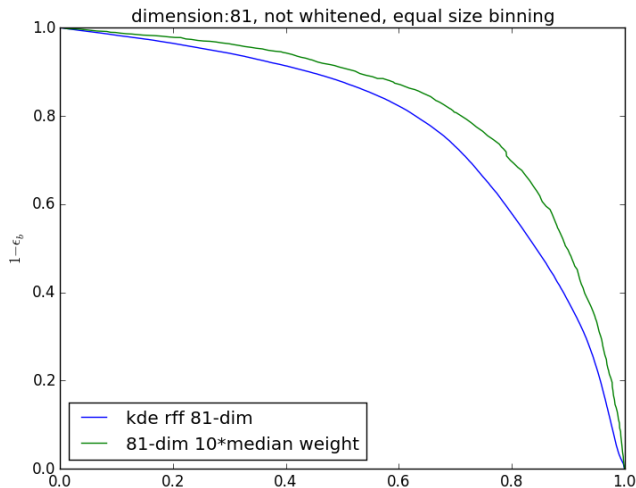
Theorem (Devroye and Krzyzak, 1989)

For any distribution of (X, Y) , if $h \rightarrow 0$ and $nh^d \rightarrow \infty$ as $n \rightarrow \infty$, i.i.d. samples, then $L(g_{\text{Gaussian},n}) \rightarrow L^$.*

This classifier can converge faster than the k-NN estimator if the conditional densities are smooth.

Random Fourier Feature Kernel Density Estimation

Randomized algorithm to approximate the Gaussian kernel, which makes it more efficient (at least a 10x speedup.)



Next Steps

- Use FLANN, a library for fast approximate nearest neighbors.
- Scale to higher dimensions: currently it takes 10 hours to run 81-dimensional data; use more data!
- Tune random Fourier Feature parameters
- Other strategies: e.g. use independent component analysis