

Kernel Classifiers

## Kernel classifiers

A kernel classifier is a plug-in classifier obtained using a kernel density estimate for each class-conditional density.

## Kernel Density Estimation

Let  $x_1, \dots, x_n$  be a random sample from a distribution with density  $g(x)$ . A kernel density estimate of  $g$  is a nonparametric estimate given by

$$\hat{g}(x) = \frac{1}{n} \sum_{i=1}^n k(x-x_i)$$

where  $k(x)$  is a kernel function.

This estimate is also called the Parzen window estimate.

A kernel function should satisfy

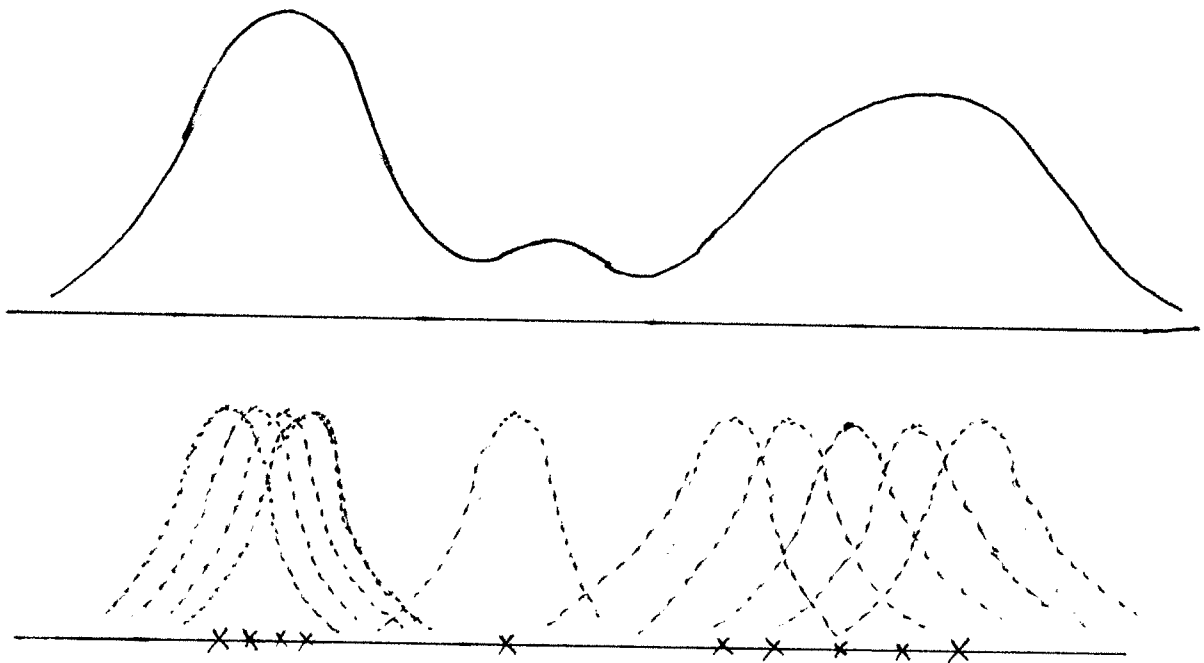
1.  $\int k(x) dx = 1$

2.  $k(x) \geq 0$

3.  $k(-x) = k(x)$

Primary example | Gaussian kernel with spherical cov.

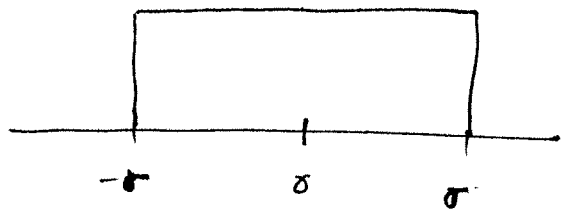
$$k(x) = k_{\sigma}(x) = (2\pi\sigma^2)^{-\frac{d}{2}} \exp\left\{-\frac{\|x\|^2}{2\sigma^2}\right\}$$



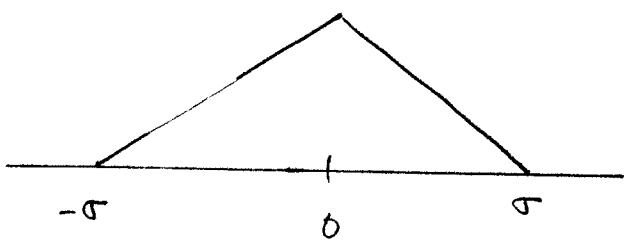
Other examples

Uniform kernel:

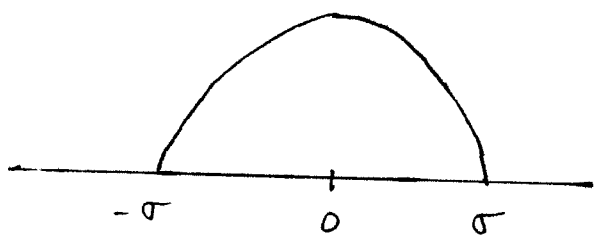
$$k_{\sigma}(x) = c_{\sigma} I_{\{\|x\| \leq \sigma\}}$$



Triangular kernel:



Epanechnikov kernel  
(parabolič)



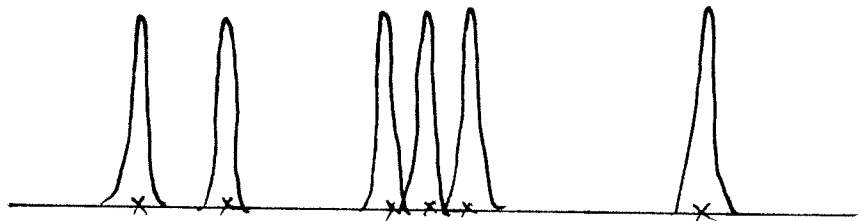
Thus, a KDE requires that we specify

1. a kernel
2. the width  $\sigma$  of the kernel.

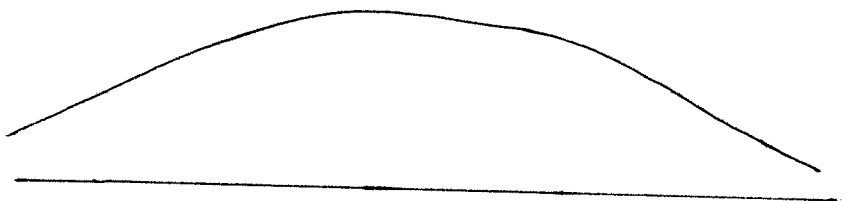
A great deal of research has been conducted on how to answer these questions. The consensus is that

- the kernel is not that critical
- the width  $\sigma$  is critical.

$\sigma$  too small



$\sigma$  too large



## Bandwidth Selection: Two Approaches

1. Estimate the probability of error of the kernel classifier as a function of  $\sigma$ , and choose the  $\sigma$  that minimizes this estimate.
2. Estimate the quality of  $\hat{g} = \hat{g}_\sigma$  as a density estimate and choose  $\sigma$  to optimize the quality.

### Bandwidth Selection by Error Estimation

Recall: The probability of error of a classifier  $f$  is

$$R(f) = P(f(X) \neq Y).$$

Let  $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$  be a training dataset.

Let  $\mathcal{A}_n^\sigma: D_n \mapsto \hat{f}_n^\sigma$  be the kernel discrimination rule/classifier based on a (fixed) kernel having bandwidth  $\sigma$ .

How should we estimate  $R(\hat{f}_n^\sigma)$ ?

## Training error estimate

$$\hat{R}(\hat{f}_n^\sigma) = \frac{1}{n} \sum_{i=1}^n I_{\{\hat{f}_n^\sigma(x_i) \neq y_i\}}$$

where

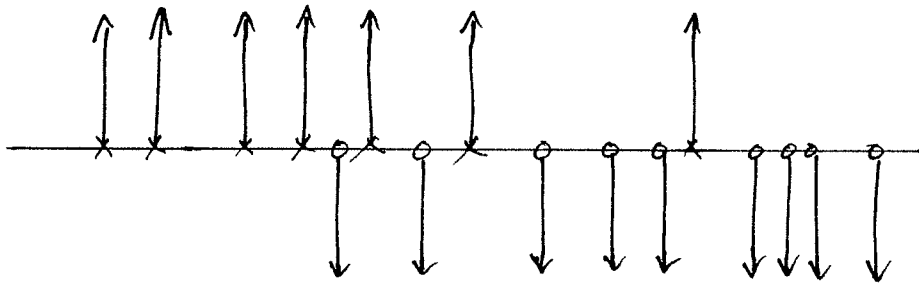
$$I_{\{p\}} = \begin{cases} 1 & \text{if } p \text{ is true} \\ 0 & \text{if } p \text{ is false} \end{cases}$$

How would we choose  $\sigma$  to minimize this error estimate?

Write the plug-in classifier as

$$\hat{\pi}_i \hat{g}_i^\sigma(x) - \hat{\pi}_0 \hat{g}_0^\sigma(x) \stackrel{!}{\geq} 0$$

Letting  $\sigma \rightarrow 0$  we get this picture:



That is,  $\hat{R}(\hat{f}_\sigma) = 0$  for  $\sigma$  sufficiently small

$\Rightarrow$  overfitting

The training error estimate is also called the apparent error rate or resubstitution error estimate.



## Holdout Error Estimate

Partition the training data into two batches

$$D_n \begin{cases} \rightarrow D_m = \{(x_1, y_1), \dots, (x_m, y_m)\} \\ \rightarrow D_n \setminus D_m = \{(x_{m+1}, y_{m+1}), \dots, (x_n, y_n)\} \end{cases}$$

Basic Idea: Run  $f_m^\sigma$  on  $D_m$  and evaluate  $\hat{f}_m^\sigma$  on  $D_n \setminus D_m$ .

Formally, the holdout error estimate is

$$\hat{R}_{HO}(\hat{f}_m^\sigma) = \frac{1}{n-m} \sum_{i=m+1}^n \mathbb{I}_{\{\hat{f}_m^\sigma(x_i) \neq y_i\}}$$

What is the expected value of  $\hat{R}_{HO}(\hat{f}_m^\sigma)$ ?

$$E\{\hat{R}_{Ho}(\hat{f}_m^*)\} = \frac{1}{n-m} \sum_{i=m+1}^n E\{I_{\{\hat{f}_m^\sigma(x_i) \neq y_i\}}\}$$

since  $D_m$ ,  
 $D_n \setminus D_m$  are  
independent

$$= \frac{1}{n-m} \sum_{i=m+1}^n P(\hat{f}_m^\sigma(x_i) \neq y_i)$$

$$= \frac{1}{n-m} \sum_{i=m+1}^n \mathcal{R}(\hat{f}_m^\sigma)$$

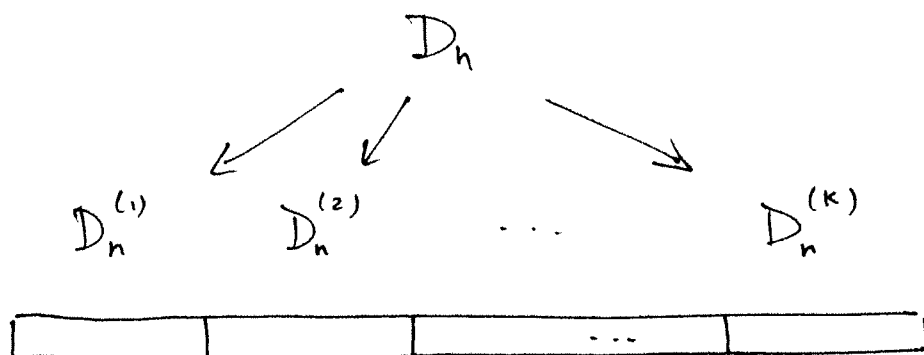
$$= \mathcal{R}(\hat{f}_m^\sigma)$$

$\Rightarrow \hat{R}_{Ho}$  is an unbiased error estimate.

Unfortunately, by holding out some of our data, our classifier is based on a fraction of the available data, which is undesirable, especially if  $n$  is small.

# Cross-Validation

Partition  $D_n$  into  $K$  equally sized groups



## K-fold Cross-Validation

1. For each  $k=1, \dots, K$ , hold out  $D_n^{(k)}$ , run  $f^\sigma$  on  $D_n \setminus D_n^{(k)}$  to get  $\hat{f}_{(k)}^\sigma$ .
2. Estimate the error probability of  $\hat{f}_{(k)}^\sigma$ :

$$\hat{R}^{(k)}(\hat{f}_{(k)}^\sigma) = \frac{1}{|D_n^{(k)}|} \sum_{(x_i, y_i) \in D_n^{(k)}} \mathbb{I}_{\{\hat{f}_{(k)}^\sigma(x_i) \neq y_i\}}$$

3. Average:

$$\hat{R}_{cr}(\hat{f}_n^\sigma) = \frac{1}{K} \sum_{k=1}^K \hat{R}^{(k)}(\hat{f}_{(k)}^\sigma)$$

The most popular choice is  $K=n$ , which is called leave-one-out cross validation (LOOCV).

How to choose  $K$ ?

- Large  $K \Rightarrow$  low bias, high variance, high computational cost
- Small  $K \Rightarrow$  higher bias, lower variance

$K=5, 10, n$  are the most common choices.

For small datasets (small  $n$ ),  $K=5$  can lead to unacceptably small training sets  $D_n \setminus D_n^{(k)}$ .

Ten Dollar Question

Why is the bias not 0?

## Bandwidth Selection by Least-Squares Cross-Validation

Let us measure the performance of the KDE  $\hat{g}_\sigma$  using the integrated squared error:

$$\text{ISE}(\sigma) = \int (\hat{g}_\sigma(x) - g(x))^2 dx$$

Idea: Choose  $\sigma$  to minimize this expression

But it depends on  $g$ , which is unknown.

What can we do?

Observe:

$$\text{ISE}(\sigma) = \int \hat{g}_\sigma(x)^2 dx - 2 \int g(x) \hat{g}_\sigma(x) dx$$

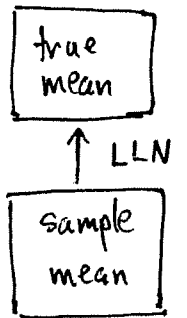
expectation w.r.t.  $g(x)$  +  $\int g(x)^2 dx$

independent of  $\sigma$

We can approximate

$$\int g(x) \hat{g}_\sigma(x) dx = E_{g(x)} \{ \hat{g}_\sigma(x) \}$$

$$\approx \frac{1}{n} \sum_{i=1}^n \hat{g}_\sigma(x_i)$$



This suggests that to select  $\sigma$ , we should minimize

$$\int \hat{g}_\sigma(x)^2 dx - \frac{2}{n} \sum_{i=1}^n \hat{g}_\sigma(x_i)$$

However, by using the sample mean we are prone to overfitting. Instead, we can employ a cross-validation estimate of  $g(x_i)$ . Define

$$\hat{g}_\sigma^{-i}(x) = \frac{1}{n-1} \sum_{\substack{j=1 \\ j \neq i}}^n K_\sigma(x-x_j)$$

Then  $\hat{g}_\sigma^{-i}(x_i)$  is a LOOCV estimate of  $g(x_i)$

Exercise | It can be shown that

$$\int K_{\sigma}(x-x_i) \cdot K_{\sigma}(x-x_j) dx = K_{\sqrt{2}\sigma}(x_i-x_j).$$

Use this fact to derive a computable expression

for

$$\text{LSCV}(\sigma) := \int \hat{g}_{\sigma}(x)^2 dx - \frac{2}{h} \sum_{i=1}^n \hat{g}_{\sigma}^{-i}(x_i)$$

least-squares cross-validation

Solution

$$\text{LSCV}(\sigma) = \int \left( \frac{1}{n} \sum_{i=1}^n K_{\sigma}(x-x_i) \right) \left( \frac{1}{n} \sum_{j=1}^n K_{\sigma}(x-x_j) \right) dx$$

$$- \frac{2}{n} \sum_{i=1}^n \frac{1}{n-1} \sum_{j \neq i} K_{\sigma}(x_i - x_j)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \int K_{\sigma}(x-x_i) K_{\sigma}(x-x_j) dx$$

$$- \frac{2}{n(n-1)} \sum_{\substack{i \neq j \\ 1 \leq i, j \leq n}} K_{\sigma}(x_i - x_j)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K_{\sqrt{2}\sigma}(x_i - x_j)$$

$$- \frac{2}{n(n-1)} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} K_{\sigma}(x_i - x_j)$$

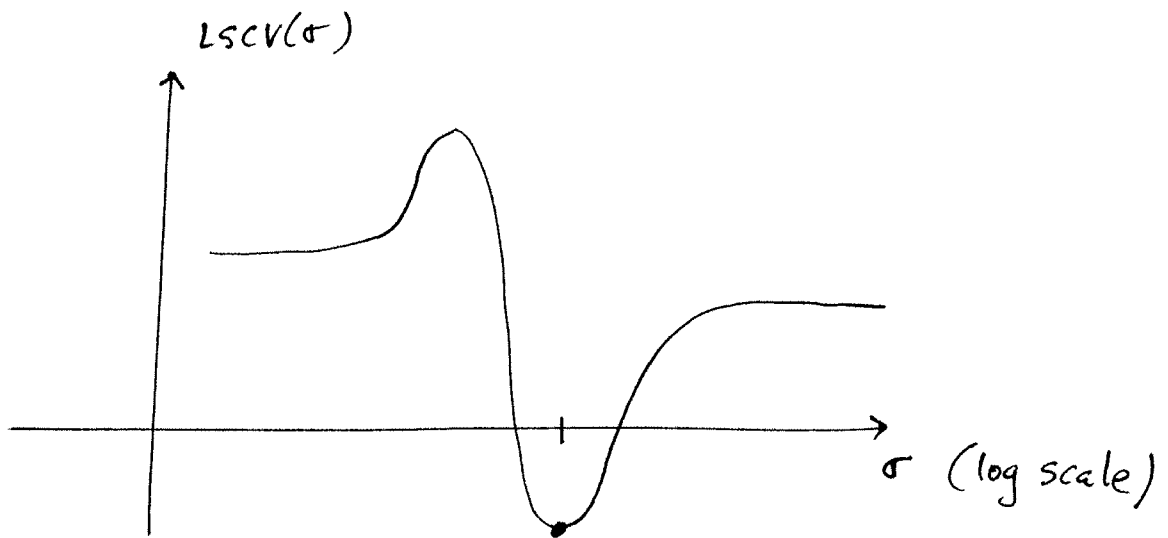
$$= \frac{1}{n^2} \cdot \sum_{i=1}^n K_{\sqrt{2}\sigma}(0)$$

$$+ \frac{1}{n} \sum_{i \neq j} \left( \frac{1}{n} K_{\sqrt{2}\sigma}(x_i - x_j) - \frac{2}{n-1} K_{\sigma}(x_i - x_j) \right)$$

$$= \boxed{\frac{(4\pi\sigma^2)^{-\frac{d}{2}}}{n} + \frac{1}{n} \sum_{i \neq j} \left( \frac{1}{n} K_{\sqrt{2}\sigma}(x_i - x_j) - \frac{2}{n-1} K_{\sigma}(x_i - x_j) \right)}$$



In practice, one evaluates  $LSCV(\sigma)$  on a predetermined grid of points and selects the minimizer



The graph is usually smooth with a clearly visible global minimizer. Local minima may also be present, however.

# Naive Bayes Classifier

The naive Bayes classifier makes the following unrealistic assumption:

The density  $g(x)$  ( $= g_0(x)$  or  $g_1(x)$ ) for either class factors as

$$g(x) = \prod_{i=1}^d g_i(x^{(i)})$$

Where  $x^{(i)}$  is the  $i^{\text{th}}$  coordinate of  $x$ .

In other words, the coordinates of  $X|Y=y$  are independent.

Under this assumption,  $g$  can be estimated by estimating each coordinate's density separately:

For each  $i=1, \dots, d$ , use  $x_1^{(i)}, \dots, x_n^{(i)}$  to estimate  $g_i(x)$  using a kernel density estimator.

The final plug-in classifier is

$$\frac{\hat{g}_1(x)}{\hat{g}_0(x)} = \frac{\prod_{i=1}^d \hat{g}_{1,i}(x^{(i)})}{\prod_{i=1}^d \hat{g}_{0,i}(x^{(i)})} \begin{matrix} > & \frac{\hat{\pi}_0}{\hat{\pi}_1} \\ < & 0 \end{matrix}$$

The obvious advantage of the Naive Bayes classifier is its simplicity: It is much easier to estimate a univariate density than a multivariate density.

Does this tradeoff outweigh the clearly incorrect modeling assumption? In many cases, yes. Keep in mind that in classification, the ultimate goal is not to learn the densities  $g_0, g_1$ , but to learn the set  $\{x: f^*(x) = 1\}$ . In many cases the deficiencies of  $\hat{g}_0, \hat{g}_1$  occur away from the decision boundary and therefore do not influence the classifier.