

# 732A99/TDDE01 Machine Learning

## Lecture 3a Block 1: Kernel Methods

Jose M. Peña

IDA, Linköping University, Sweden

# Contents

- ▶ Histogram, Moving Window, and Kernel Classification
- ▶ Histogram, Moving Window, and Kernel Regression
- ▶ Histogram, Moving Window, and Kernel Density Estimation
- ▶ Kernel Selection
- ▶ Kernel Trick
- ▶ Summary

# Literature

- ▶ Main source
  - ▶ Bishop, C. M. *Pattern Recognition and Machine Learning*. Springer, 2006. Sections 2.5 and 6.1-6.2.
- ▶ Additional source
  - ▶ Devroye, L., Györfi, L. and Lugosi, G. *A Probabilistic Theory of Pattern Recognition*. Springer, 1996. Sections 6.4 and 10.0.
  - ▶ Hastie, T., Tibshirani, R. and Friedman, J. *The Elements of Statistical Learning*. Springer, 2009. Chapter 6.

## Histogram Classification

- ▶ Consider binary classification with input space  $\mathbb{R}^D$ .
- ▶ The best classifier under the 0-1 loss function is  $y^*(\mathbf{x}) = \arg \max_y p(y|\mathbf{x})$ .
- ▶ Since  $\mathbf{x}$  may not appear in the finite training set  $\{(\mathbf{x}_n, t_n)\}$  available, then
  - ▶ divide the input space into  $D$ -dimensional cubes of side  $h$ , and
  - ▶ classify according to majority vote in the cube  $C(\mathbf{x}, h)$  that contains  $\mathbf{x}$ .

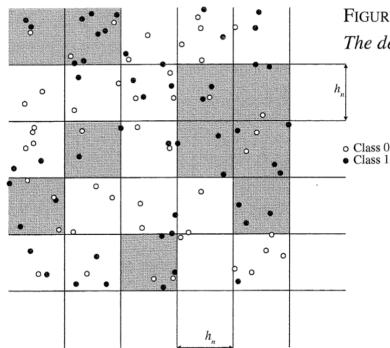


FIGURE 6.1. A cubic histogram rule:  
The decision is 1 in the shaded area.

- ▶ In other words,

$$y_C(\mathbf{x}) = \begin{cases} 0 & \text{if } \sum_n \mathbf{1}_{\{t_n=1, \mathbf{x}_n \in C(\mathbf{x}, h)\}} \leq \sum_n \mathbf{1}_{\{t_n=0, \mathbf{x}_n \in C(\mathbf{x}, h)\}} \\ 1 & \text{otherwise} \end{cases}$$

## Moving Window Classification

- ▶ The histogram rule is less accurate at the borders of the cube, because those points are not as well represented by the cube as the ones near the center. Then,
  - ▶ consider the points within a certain distance to the point to classify, and
  - ▶ classify the point according to majority vote.

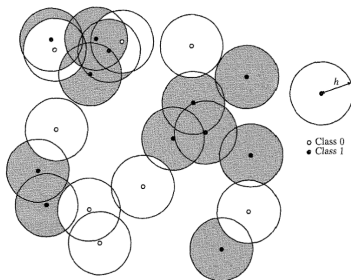


FIGURE 10.1. *The moving window rule in  $\mathcal{R}^2$ . The decision is 1 in the shaded area.*

- ▶ In other words,

$$y_S(\mathbf{x}) = \begin{cases} 0 & \text{if } \sum_n \mathbf{1}_{\{t_n=1, \mathbf{x}_n \in S(\mathbf{x}, h)\}} \leq \sum_n \mathbf{1}_{\{t_n=0, \mathbf{x}_n \in S(\mathbf{x}, h)\}} \\ 1 & \text{otherwise} \end{cases}$$

where  $S(\mathbf{x}, h)$  is a  $D$ -dimensional closed ball of radius  $h$  centered at  $\mathbf{x}$ .

## Kernel Classification

- ▶ The moving window rule gives equal weight to all the points in the ball, which may be counterintuitive. Then,

$$y_k(\mathbf{x}) = \begin{cases} 0 & \text{if } \sum_n \mathbf{1}_{\{t_n=1\}} k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) \leq \sum_n \mathbf{1}_{\{t_n=0\}} k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) \\ 1 & \text{otherwise} \end{cases}$$

where  $k : \mathbb{R}^D \rightarrow \mathbb{R}$  is a kernel function, which is usually non-negative and monotone decreasing along rays starting from the origin. The parameter  $h$  is called smoothing factor or width.

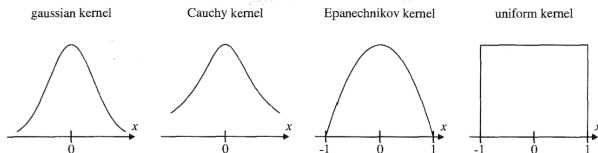


FIGURE 10.3. Various kernels on  $\mathcal{R}$ .

- ▶ Gaussian kernel:  $k(u) = \exp(-\|u\|^2)$  where  $\|\cdot\|$  is the Euclidean norm.
- ▶ Cauchy kernel:  $k(u) = 1/(1 + \|u\|^{D+1})$
- ▶ Epanechnikov kernel:  $k(u) = (1 - \|u\|^2)\mathbf{1}_{\{\|u\| \leq 1\}}$
- ▶ Moving window kernel:  $k(u) = \mathbf{1}_{\{u \in S(0,1)\}}$

# Kernel Classification

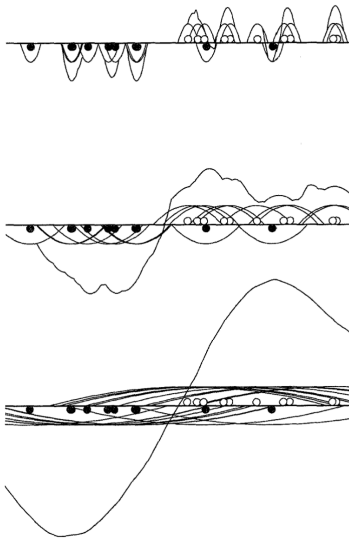


FIGURE 10.2. *Kernel rule on the real line. The figure shows  $\sum_{i=1}^n (2Y_i - 1)K((x - X_i)/h)$  for  $n = 20$ ,  $K(u) = (1 - u^2)I_{\{|u| \leq 1\}}$  (the Epanechnikov kernel), and three smoothing factors  $h$ . One definitely undersmooths and one oversmooths. We took  $p = 1/2$ , and the class-conditional densities are  $f_0(x) = 2(1 - x)$  and  $f_1(x) = 2x$  on  $[0, 1]$ .*

# Histogram, Moving Window, and Kernel Regression

- ▶ Consider regressing an unidimensional continuous random variable on a  $D$ -dimensional continuous random variable.
- ▶ The best regression function under the squared error loss function is  $y^*(\mathbf{x}) = \mathbb{E}_Y[y|\mathbf{x}]$ .
- ▶ Since  $\mathbf{x}$  may not appear in the finite training set  $\{(\mathbf{x}_n, t_n)\}$  available, then we average over the points in  $C(\mathbf{x}, h)$  or  $S(\mathbf{x}, h)$ , or kernel-weighted average over all the points.
- ▶ In other words,

$$y_C(\mathbf{x}) = \frac{\sum_{\mathbf{x}_n \in C(\mathbf{x}, h)} t_n}{|\{\mathbf{x}_n \in C(\mathbf{x}, h)\}|}$$

or

$$y_S(\mathbf{x}) = \frac{\sum_{\mathbf{x}_n \in S(\mathbf{x}, h)} t_n}{|\{\mathbf{x}_n \in S(\mathbf{x}, h)\}|}$$

or

$$y_k(\mathbf{x}) = \frac{\sum_n k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) t_n}{\sum_n k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right)}$$



## Histogram, Moving Window, and Kernel Density Estimation

- ▶ Consider density estimation for a  $D$ -dimensional continuous random variable.
- ▶ Let  $R \subseteq \mathbb{R}^D$  and  $\mathbf{x} \in R$ . Then,

$$P = \int_R p(\mathbf{x}) d\mathbf{x} \simeq p(\mathbf{x}) \text{Volume}(R)$$

and the number of the  $N$  training points  $\{\mathbf{x}_n\}$  that fall inside  $R$  is

$$|\{\mathbf{x}_n \in R\}| \simeq P N$$

and thus

$$p(\mathbf{x}) \simeq \frac{|\{\mathbf{x}_n \in R\}|}{N \text{Volume}(R)}$$

- ▶ Then,

$$p_C(\mathbf{x}) = \frac{|\{\mathbf{x}_n \in C(\mathbf{x}, h)\}|}{N \text{Volume}(C(\mathbf{x}, h))}$$

or

$$p_S(\mathbf{x}) = \frac{|\{\mathbf{x}_n \in S(\mathbf{x}, h)\}|}{N \text{Volume}(S(\mathbf{x}, h))}$$

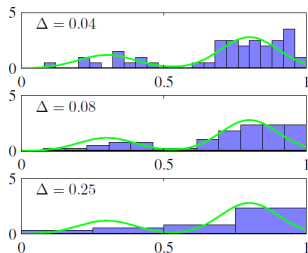
or

$$p_k(\mathbf{x}) = \frac{1}{N} \sum_n k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

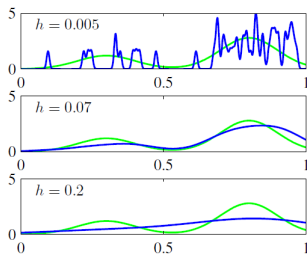
assuming that  $k(u) \geq 0$  for all  $u$  and  $\int k(u) du = 1$ .

# Histogram, Moving Window, and Kernel Density Estimation

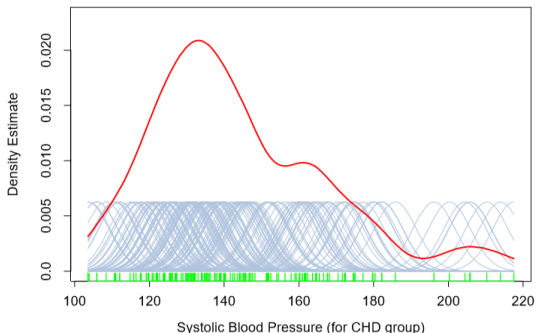
**Figure 2.24** An illustration of the histogram approach to density estimation, in which a data set of 50 data points is generated from the distribution shown by the green curve. Histogram density estimates, based on (2.241), with a common bin width  $\Delta$  are shown for various values of  $\Delta$ .



**Figure 2.25** Illustration of the kernel density model (2.250) applied to the same data set used to demonstrate the histogram approach in Figure 2.24. We see that  $h$  acts as a smoothing parameter and that if it is set too small (top panel), the result is a very noisy density model, whereas if it is set too large (bottom panel), then the bimodal nature of the underlying distribution from which the data is generated (shown by the green curve) is washed out. The best density model is obtained for some intermediate value of  $h$  (middle panel).



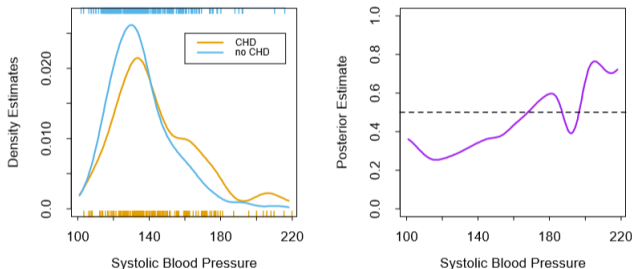
# Histogram, Moving Window, and Kernel Density Estimation



**FIGURE 6.13.** A kernel density estimate for systolic blood pressure (for the CHD group). The density estimate at each point is the average contribution from each of the kernels at that point. We have scaled the kernels down by a factor of 10 to make the graph readable.

- From kernel density estimation to kernel classification:
  1. Estimate  $p(\mathbf{x}|y = 0)$  and  $p(\mathbf{x}|y = 1)$  using the methods just seen.
  2. Estimate  $p(y)$  as class proportions.
  3. Compute  $p(y|\mathbf{x}) \propto p(\mathbf{x}|y)p(y)$  by Bayes theorem.

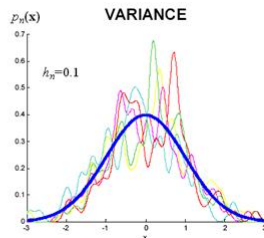
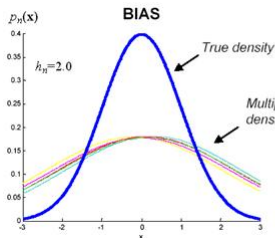
# Histogram, Moving Window, and Kernel Density Estimation



**FIGURE 6.14.** The left panel shows the two separate density estimates for systolic blood pressure in the CHD versus no-CHD groups, using a Gaussian kernel density estimate in each. The right panel shows the estimated posterior probabilities for CHD, using (6.25).

## Kernel Selection

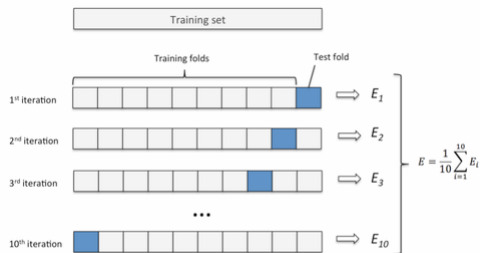
- ▶ How to choose the right kernel and width ? E.g., by cross-validation.
- ▶ What does “right” mean ? E.g., minimize loss function.
- ▶ Note that the width of the kernel corresponds to a bias-variance trade-off.



- ▶ Small width implies considering few points. So, the variance will be large (similar to the variance of a single point). The bias will be small since the points considered are close to  $\mathbf{x}$ .
- ▶ Large width implies considering many points. So, the variance will be small and the bias will be large.

## Kernel Selection

- ▶ Recall the following from previous lectures.
- ▶ Cross-validation is a technique to estimate the prediction error of a model.



- ▶ If the training set contains  $N$  points, note that cross-validation estimates the prediction error when the model is trained on  $N - N/K$  points.
- ▶ Note that the model returned is trained **on  $N$  points**. So, cross-validation overestimates the prediction error of the model returned.
- ▶ This seems to suggest that a large  $K$  should be preferred. However, this typically implies a large variance of the error estimate, since there are only  $N/K$  test points.
- ▶ Typically,  $K = 5, 10$  works well.

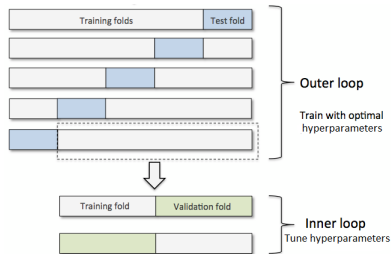
## Kernel Selection

- ▶ **Nested** cross-validation is a technique to estimate the prediction error of model **selection**, e.g. selection of the kernel and its width for classification, regression or density estimation.
- ▶ Nested cross-validation consists of an inner and an outer loop:
  - ▶ Inner loop: Model selection via cross-validation.
  - ▶ Outer loop: Error estimation of the model selected.

### Cross-validation for estimating model prediction error



### Nested cross-validation for estimating model **selection** prediction error



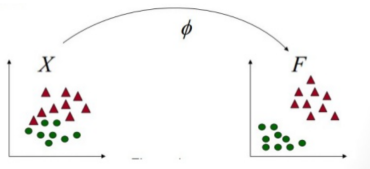
- ▶ Error overestimation may not be a concern for model selection. So,  $K = 2$  may suffice in the inner loop.
- ▶ Note that the model returned is the result of model selection via cross-validation **on  $N$  points**.

## Kernel Trick

- ▶ The kernel function  $k\left(\frac{\mathbf{x}-\mathbf{x}'}{h}\right)$  is invariant to translations, and it can be generalized as  $k(\mathbf{x}, \mathbf{x}')$ . For instance,
  - ▶ Polynomial kernel:  $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M$
  - ▶ Gaussian kernel:  $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$
- ▶ If the matrix

$$\begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \dots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is symmetric and positive semi-definite for all choices of  $\{\mathbf{x}_n\}$ , then  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$  where  $\phi(\cdot)$  is a mapping from the input space to the feature space.



- ▶ The feature space may be non-linear and even infinite dimensional. For instance,

$$\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}cx_1, \sqrt{2}cx_2, c)$$

for the polynomial kernel with  $M = D = 2$ .



## Kernel Trick

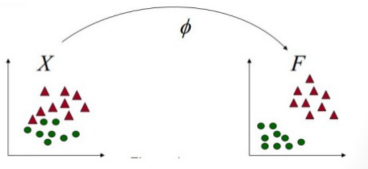
- ▶ Consider again moving window classification, regression, and density estimation.
- ▶ Note that  $\mathbf{x}_n \in S(\mathbf{x}, h)$  if and only if  $\|\mathbf{x} - \mathbf{x}_n\| \leq h$ .
- ▶ Note that

$$\|\mathbf{x} - \mathbf{x}_n\| = \sqrt{(\mathbf{x} - \mathbf{x}_n)^T (\mathbf{x} - \mathbf{x}_n)} = \sqrt{\mathbf{x}^T \mathbf{x} + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}^T \mathbf{x}_n}$$

- ▶ Then,

$$\begin{aligned}\|\phi(\mathbf{x}) - \phi(\mathbf{x}_n)\| &= \sqrt{\phi(\mathbf{x}^T)\phi(\mathbf{x}) + \phi(\mathbf{x}_n^T)\phi(\mathbf{x}_n) - 2\phi(\mathbf{x}^T)\phi(\mathbf{x}_n)} \\ &= \sqrt{k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}, \mathbf{x}_n)}\end{aligned}$$

- ▶ So, the distance is now computed in a (hopefully) more convenient space.



- ▶ Note that we do not need to compute  $\phi(\mathbf{x})$  and  $\phi(\mathbf{x}_n)$ .

# Kernel Trick

- ▶ Two alternatives for building  $k(\mathbf{x}, \mathbf{x}')$ :
  - ▶ Choose a convenient  $\phi(\mathbf{x})$  and let  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$ .
  - ▶ Build it from existing kernel functions as follows.

## Techniques for Constructing New Kernels.

Given valid kernels  $k_1(\mathbf{x}, \mathbf{x}')$  and  $k_2(\mathbf{x}, \mathbf{x}')$ , the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

where  $c > 0$  is a constant,  $f(\cdot)$  is any function,  $q(\cdot)$  is a polynomial with nonnegative coefficients,  $\phi(\mathbf{x})$  is a function from  $\mathbf{x}$  to  $\mathbb{R}^M$ ,  $k_3(\cdot, \cdot)$  is a valid kernel in  $\mathbb{R}^M$ ,  $\mathbf{A}$  is a symmetric positive semidefinite matrix,  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are variables (not necessarily disjoint) with  $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ , and  $k_a$  and  $k_b$  are valid kernel functions over their respective spaces.

# Summary

- ▶ Kernel methods: Smoothing models.
- ▶ Model selection: Nested cross-validation.
- ▶ Kernel trick: It allows to work in the feature space without constructing it.