

Lecture 4

Nearest Neighbors
Kernel predictors

Material:
to end of
Lecture I

HW 4 - t.b. posted
today (web)

Q1 - Oct 18

Zoom recordings
t.b. found

Tutorial - today !!

Lecture Notes I – Examples of Predictors

Marina Meilă
`mmp@stat.washington.edu`

Department of Statistics
University of Washington

September 29, 2022

Prediction problems by the type of output ✓

The “learning” paradigm and vocabulary ✓

The Nearest-Neighbor and kernel predictors ←

Linear predictors

Least squares regression

Linear Discriminant Analysis (LDA)

QDA (Quadratic Discriminant Analysis)

Logistic Regression

The PERCEPTRON algorithm

Classification and regression tree(s) (CART)

(The Naive Bayes classifier)

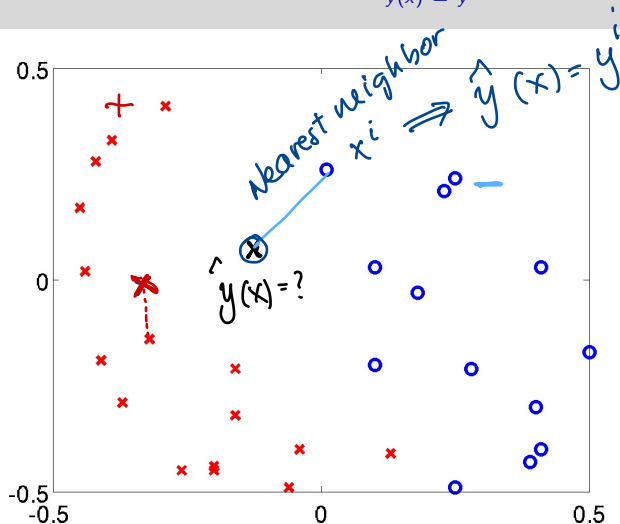
Reading HTF Ch.: 2.3.1 Linear regression, 2.3.2 Nearest neighbor, 4.1–4 Linear classification, 6.1–3. Kernel regression, 6.6.2 kernel classifiers, 6.6.3 Naive Bayes, 9.2 CART, 11.3 Neural networks, Murphy Ch.: 1.4.2 nearest neighbors, 1.4.4 linear regression, 1.4.5 logistic regression, 3.5 and 10.2.1 Naive Bayes, 4.2.1–3 linear and quadratic discriminant, 14.7.3– kernel regression, locally weighted regression, 16.2.1–4 CART, (16.5 neural nets), Bach Ch.:

The Nearest-Neighbor predictor

► **1-Nearest Neighbor** The label of a point x is assigned as follows:

1. find the example x^i that is **nearest** to x in \mathcal{D} (in Euclidean distance)
2. assign x the label y^i , i.e.

$$\hat{y}(x) = y^i$$

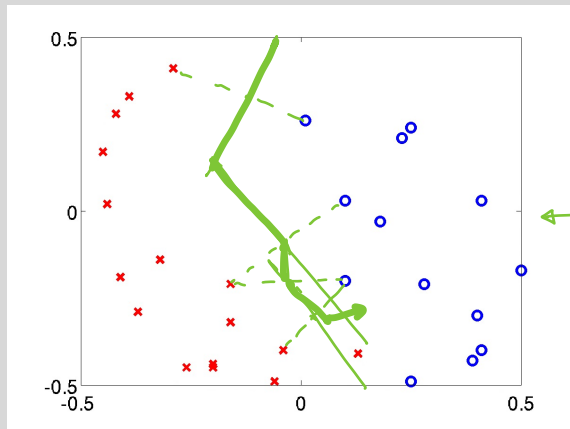


The Nearest-Neighbor predictor

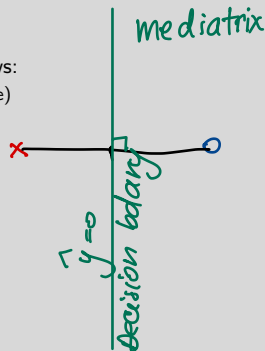
► **1-Nearest Neighbor** The label of a point x is assigned as follows:

1. find the example x^i that is **nearest** to x in \mathcal{D} (in Euclidean distance)
2. assign x the label y^i , i.e.

$$\hat{y}(x) = y^i$$



- decision bldary can take arbitrary shape



← Decision bldary:

- piecewise linear

Union of subsets of hyperplanes

mediatrices btw $(x^i, x^j), i \neq j$

The Nearest-Neighbor predictor

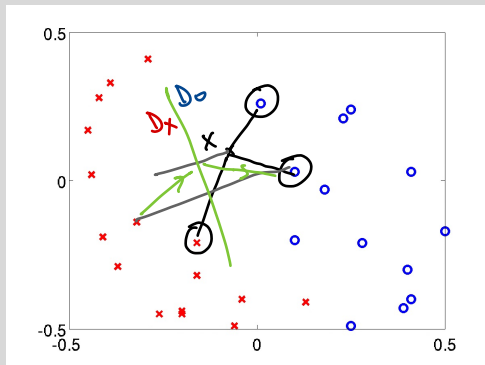
► **1-Nearest Neighbor** The label of a point x is assigned as follows:

1. find the example x^i that is nearest to x in \mathcal{D} (in Euclidean distance)
2. assign x the label y^i , i.e.

$$\hat{y}(x) = y^i$$

► **K-Nearest Neighbor** (with $K = 3, 5$ or larger)

1. find the K nearest neighbors of x in \mathcal{D} : x^1, \dots, x^K
2. ► for classification $f(x)$ = the most frequent label among the K neighbors (well suited for multiclass)
- for regression $f(x) = \frac{1}{K} \sum_{i \text{ neighbor of } x} y^i$ = mean of neighbors' labels



$K=3$

$x: 2 \circ, 1 \times \Rightarrow \hat{y}(x) = \circ$

Decision boundary:
- piecewise linear
Exercise

The Nearest-Neighbor predictor

► 1-Nearest Neighbor The label of a point x is assigned as follows:

1. find the example x^i that is nearest to x in \mathcal{D} (in Euclidean distance)
2. assign x the label y^i , i.e.



$$\hat{y}(x) = y^i$$

► K-Nearest Neighbor (with $K = 3, 5$ or larger)

1. find the K nearest neighbors of x in \mathcal{D} : x^{i_1}, \dots, x^{i_K}
2.
 - for classification $f(x)$ = the most frequent label among the K neighbors (well suited for multiclass)
 - for regression $f(x) = \frac{1}{K} \sum_{i \text{ neighbor of } x} y^i$ = mean of neighbors' labels

2-class

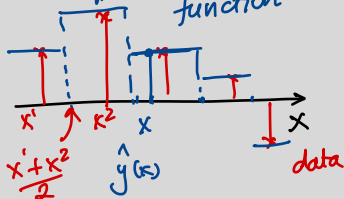
OR
multiclass

↓
ties possible
⇒ on decision
bdary.

1-NN regression

piece wise
constant
function

pcwise constant



The Nearest-Neighbor predictor

- **1-Nearest Neighbor** The label of a point x is assigned as follows:

1. find the example x^i that is nearest to x in \mathcal{D} (in Euclidean distance)
2. assign x the label y^i , i.e.

$$\hat{y}(x) = y^i$$

- **K-Nearest Neighbor** (with $K = 3, 5$ or larger)

1. find the K nearest neighbors of x in \mathcal{D} : x^{i_1}, \dots, x^{i_K}
2. ► for classification $f(x)$ = the most frequent label among the K neighbors (well suited for multiclass)
- for regression $f(x) = \frac{1}{K} \sum_{i \text{ neighbor of } x} y^i$ = mean of neighbors' labels

• smoothness parameter ^{↑ with K}
 (hyperparameter)

- No parameters to estimate!

- No training!

- But all data must be stored (also called memory-based learning)

— computing $\hat{y} \rightarrow$ requires search of \mathcal{D} !

15-Nearest Neighbor Classifier

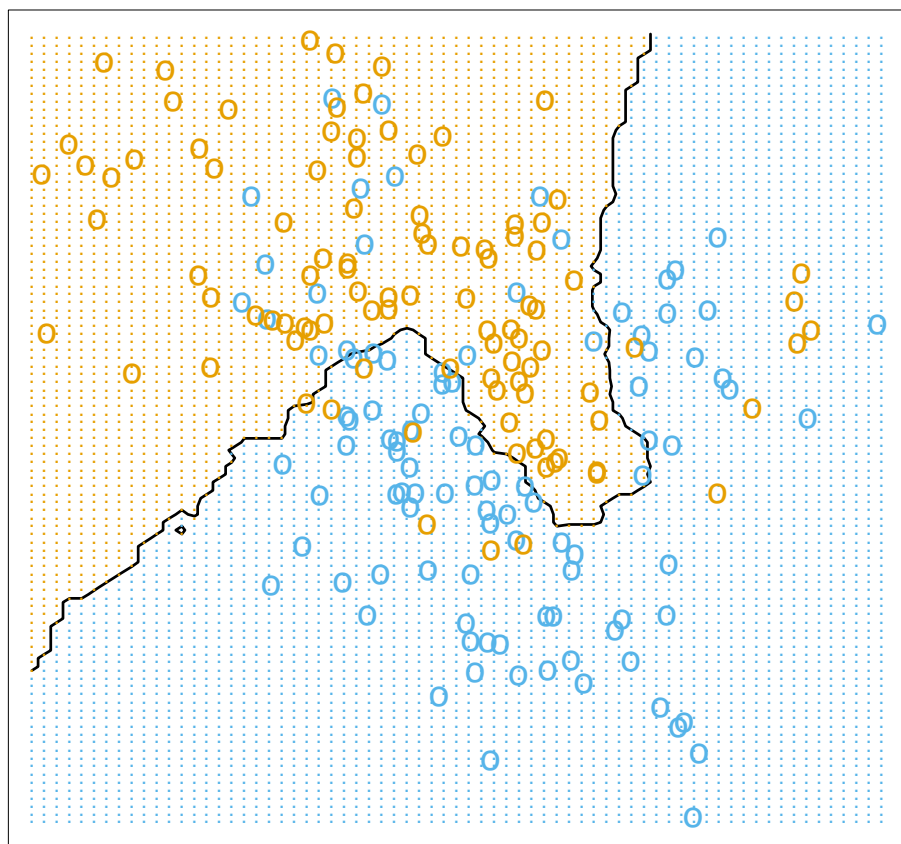


FIGURE 2.2. *The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1) and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors.*

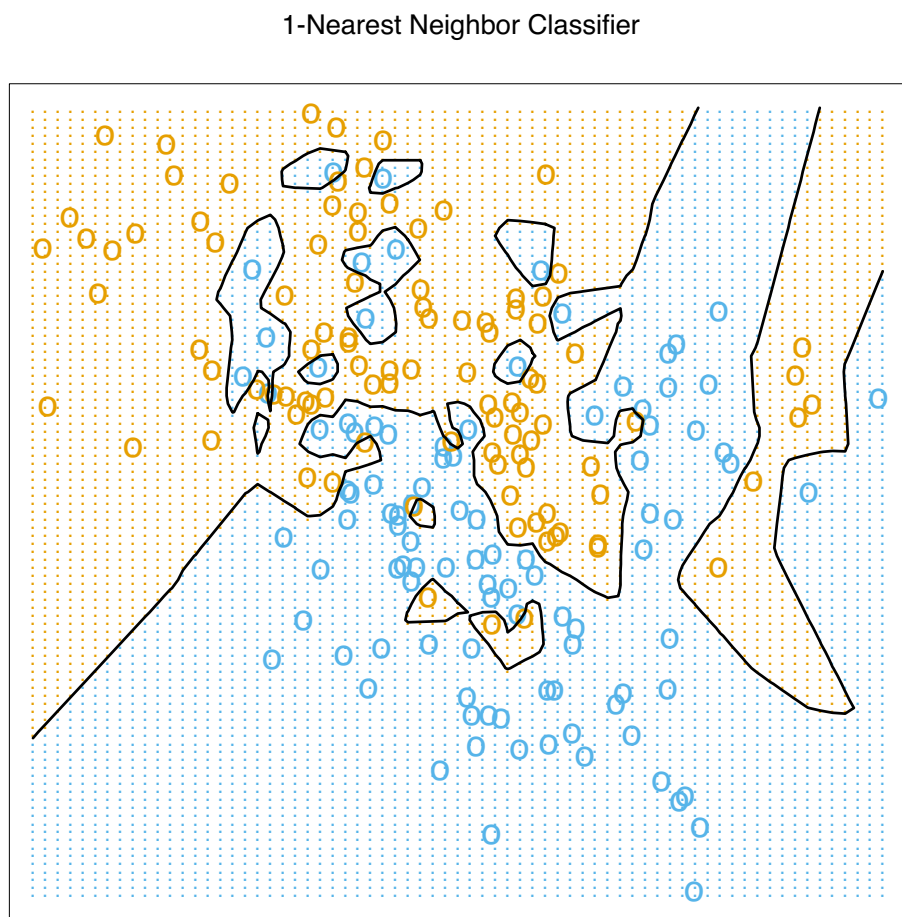


FIGURE 2.3. *The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1), and then predicted by 1-nearest-neighbor classification.*

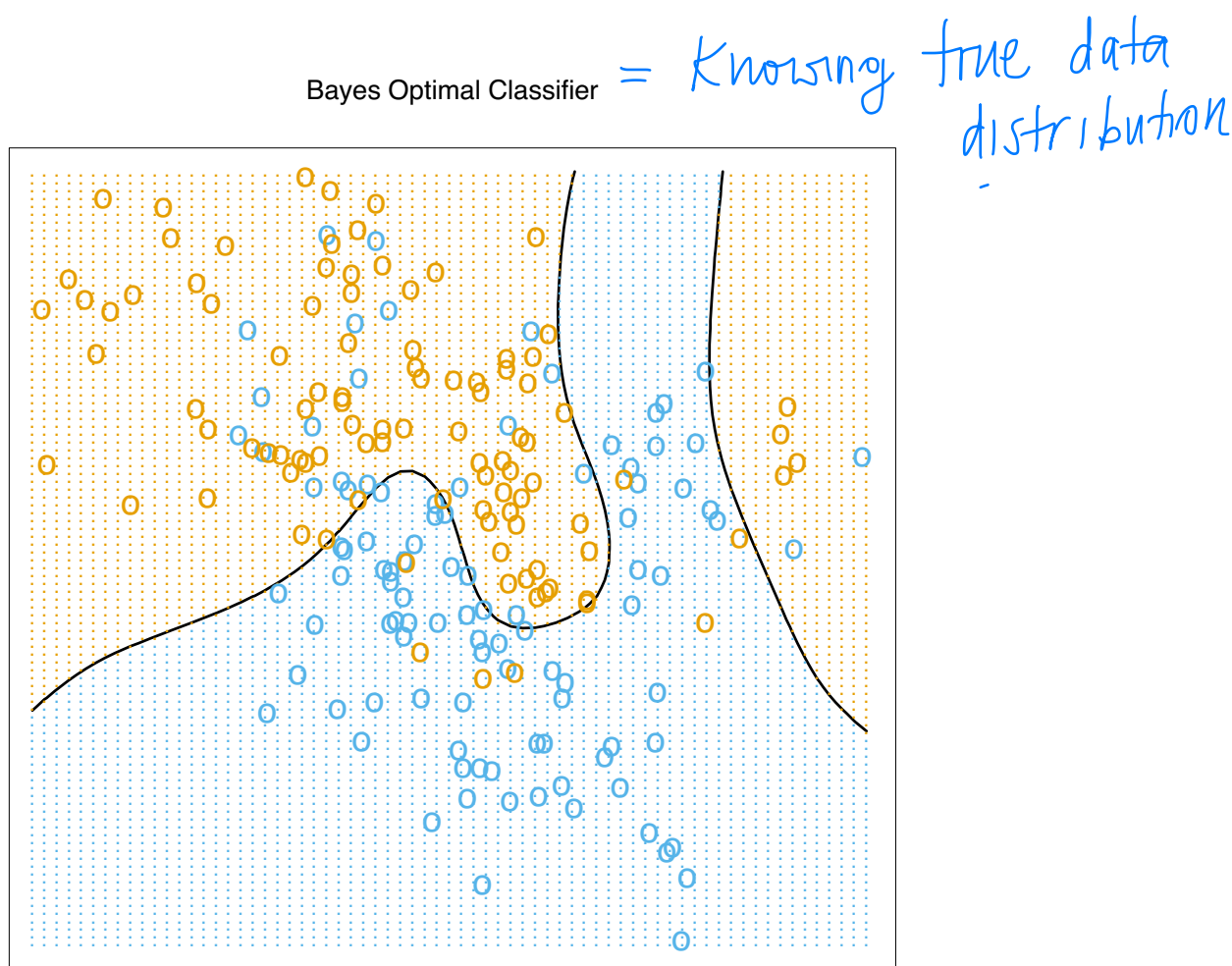


FIGURE 2.5. *The optimal Bayes decision boundary for the simulation example of Figures 2.1, 2.2 and 2.3. Since the generating density is known for each class, this boundary can be calculated exactly (Exercise 2.2).*

Kernel regression and classification

Kernel Predictors

- ▶ Like the K -nearest neighbor but with "smoothed" neighborhoods
- ▶ The predictor

$$f(\underline{x}) = \sum_{i=1}^n \beta_i \underbrace{b(\underline{x}, \underline{x}^i)}_{\text{kernel}} \underbrace{y^i}_{\text{data of } y^i \text{ of neighbors}} \quad (1)$$

← \approx weighted avg of y^i of neighbors

weights

where β_i are coefficients

Kernel $b(z)$

- $b(z) \geq 0$ for all z

- $\int_{-\infty}^{\infty} b(z) dz < \infty$

$$\left[\int_{-\infty}^{\infty} b(z) dz = 1 \text{ in KDE} \right]$$

w.l.o.g (without loss of generality)

Typically:

- $b(z) = b(-z)$
 - smooth $\rightarrow C^\infty$ b', b'', \dots exist a.e. (almost everywhere)
 - compact support $\Leftrightarrow b(z) = 0$ for $|z| > R \Rightarrow \text{supp } b = \overline{B_R(0)} = \{z \mid \|z\| \leq R\}$
- $\text{supp } f = \{z \mid f(z) \neq 0\}$

Kernel regression and classification

$$x \in \mathbb{R}^d$$

- ▶ Like the K -nearest neighbor but with "smoothed" neighborhoods
- ▶ The predictor

$$f(x) = \sum_{i=1}^n \beta_i b(x, x^i) y^i \quad (1)$$

where β_i are coefficients

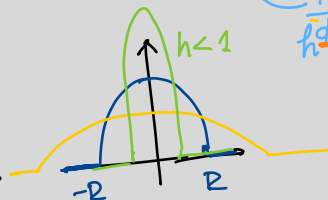
Example:



Gaussian

$$\underline{b(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}}$$

$$b: \mathbb{R} \rightarrow [0, \infty)$$



Quadratic
(compact support)

$\frac{1}{h^d} b\left(\frac{\|x - x^i\|}{h}\right)$
kernel bandwidth

• h is smoothing parameter

Kernel regression and classification

- ▶ Like the K -nearest neighbor but with “smoothed” neighborhoods
- ▶ The predictor

$$f(x) = \sum_{i=1}^n \beta_i b(x, x^i) y^i \quad (1)$$

where β_i are coefficients

- ▶ Intuition: center a “bell-shaped” *kernel* function b on each data point, and obtain the prediction $f(x)$ as a weighted sum of the values y^i , where the weights are $\beta_i b(x, x^i)$
- ▶ Requirements for a kernel function $b(x, x')$
 1. non-negativity
 2. symmetry in the arguments x, x'
 3. optional: radial symmetry, bounded support, smoothness
- ▶ A typical kernel function is the **Gaussian kernel** (or **Radial Basis Function (RBF)**)

$$b(z) \propto e^{-z^2/2} \quad (2)$$

$$b_h(x, x') \propto e^{-\frac{\|x-x'\|^2}{2h^2}} \quad \text{with } h = \text{the kernel width} \quad (3)$$

Regression example

A special case in wide use is the Nadaraya-Watson regressor

• $x = \text{query point}$

$$f(x) = \frac{\sum_{i=1}^n b\left(\frac{\|x - x^i\|}{h}\right) y^i}{\sum_{i=1}^n b\left(\frac{\|x - x^i\|}{h}\right)} = \sum_{i=1}^n \beta \frac{b\left(\frac{\|x - x^i\|}{h}\right)}{\sum_{i=1}^n b\left(\frac{\|x - x^i\|}{h}\right)} y^i \quad (4)$$

normalizing constant (x)

In this regressor, $f(x)$ is always a convex combination of the y^i 's, and the weights are proportional to $b_h(x, x^i)$.

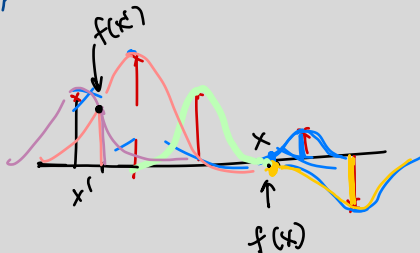
The Nadaraya-Watson regressor is biased if the density of P_X varies around x .

$$\beta = \frac{1}{\sum_{i=1}^n b\left(\frac{\|x - x^i\|}{h}\right)} \cdot \frac{1}{h^d}$$

$$f(x) = \sum w_i y^i \quad \leftarrow \text{convex combination of } y^i\text{'s}$$

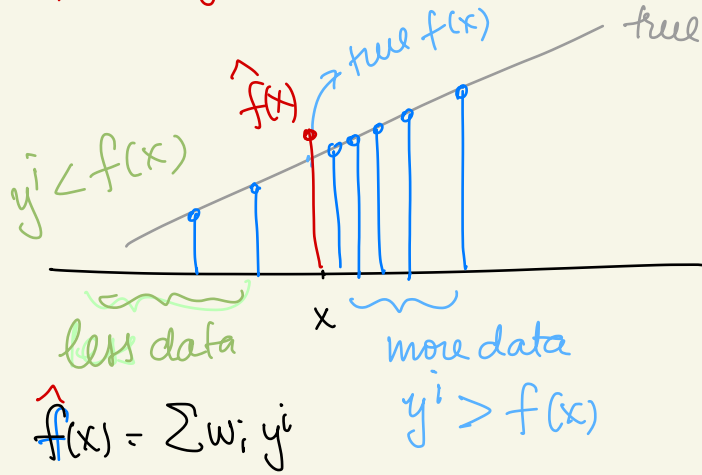
$$\sum_{i=1}^n w_i(x, h) = 1$$

- $\sum \alpha_i y_i = \text{convex combination}$ if $\alpha_i \geq 0$ for all i and $\sum \alpha_i = 1$



Nadaraya Watson bias → defined later

(FYI bias $\propto \nabla f(x) \cdot \nabla p(x)$)
[Without proof]
function density



Local Linear Regression

To correct for the bias (to first order) one can estimate a regression line around x .

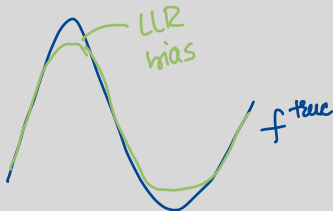
1. Given **query point** x
2. Compute kernel $b_h(x, x^i) = w_i$ for all $i = 1, \dots, N$
3. Solve **weighted regression** $\min_{\beta, \beta_0} \sum_{i=1}^d w_i (y^i - \beta^T x^i - \beta_0)^2$ to obtain β, β_0
(β, β_0 depend on x through w_i)
4. Calculate $f(x) = \beta^T x + \beta_0$ ← *evaluated at 1 point!*

Exercise Show that Nadaraya-Watson solves a local linear regression with fixed $\beta = 0$

Idea

N-Watson : $f^{NW}(x) = \sum w_i y_i$ ↗ best constant fcn fit to weighted data

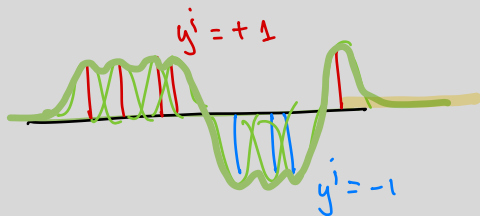
LLR : $f^{LLR}(x) =$ best linear fcn fit to weighted data



Kernel binary classifiers

- Obtained from Nadaraya-Watson by setting y^i to ± 1 .
- Note that the classifier can be written as the difference of two non-negative functions

$$f(x) \propto \sum_{i:y^i=1} b\left(\frac{\|x-x^i\|}{h}\right) - \sum_{i:y^i=-1} b\left(\frac{\|x-x^i\|}{h}\right). \quad (5)$$



$$\hat{y} = \text{sign } f(x)$$