

Instance-based Learning

Machine Learning

Hamid R Rabiee – Zahra Dehghanian
Spring 2025



Sharif University
of Technology

Outline

- Non-parametric approach
 - Unsupervised: Non-parametric density estimation
 - Parzen Windows
 - Kn-Nearest Neighbor Density Estimation
 - Supervised: Instance-based learners
 - Classification
 - kNN classification
 - Weighted (or kernel) kNN
 - Regression
 - kNN regression
 - Locally linear weighted regression

Introduction

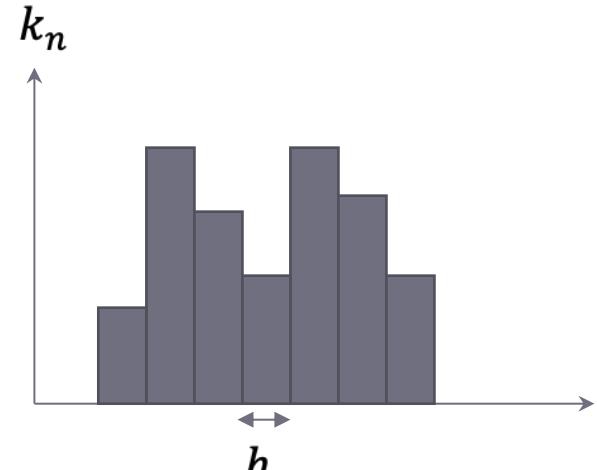
- **Estimation of arbitrary density functions**
 - Parametric density functions cannot usually fit the densities we encounter in practical problems
 - e.g., parametric densities are unimodal
 - Non-parametric methods don't assume that the model (from) of underlying densities is known in advance
- **Non-parametric methods (for classification) can be categorized into**
 - Generative
 - Estimate $p(\mathbf{x}|\mathcal{C}_i)$ from \mathcal{D}_i using non-parametric density estimation
 - Discriminative
 - Estimate $p(\mathcal{C}_i|\mathbf{x})$ from \mathcal{D}

Parametric vs. nonparametric methods

- Parametric methods need to find parameters from data and then use the inferred parameters to decide on new data points
 - Learning: finding parameters from data
- Nonparametric methods
 - Training examples are explicitly used
 - Training phase is not required
- Both supervised and unsupervised learning methods can be categorized into parametric and non-parametric methods.

Histogram approximation idea

- **Histogram** approximation of an unknown pdf
 - $P(b_l) \approx k_n(b_l)/n \quad l = 1, \dots, L$
 - $k_n(b_l)$: number of samples (among n ones) lied in the bin b_l



- The corresponding estimated pdf:

$$\bullet \hat{p}(x) = \frac{P(b_l)}{h} \quad |x - \bar{x}_{b_l}| \leq \frac{h}{2}$$

↓
Mid-point of
the bin b_l

Non-parametric density estimation

- ▶ Probability of falling in a region \mathcal{R} :
 - ▶ $P = \int_{\mathcal{R}} p(x')dx'$ (smoothed version of $p(x)$)
- ▶ $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n$: a set of samples drawn i.i.d. according to $p(x)$
 - ▶ The probability that k of the n samples fall in \mathcal{R} :
 - ▶ $P_k = \binom{n}{k} P^k (1 - P)^{n-k}$
 - ▶ $E[k] = nP$

Non-parametric density estimation

- ▶ Probability of falling in a region \mathcal{R} :
 - ▶ $P = \int_{\mathcal{R}} p(x')dx'$ (smoothed version of $p(x)$)
 - ▶ $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n$: a set of samples drawn i.i.d. according to $p(x)$
 - ▶ The probability that k of the n samples fall in \mathcal{R} :
 - ▶ $P_k = \binom{n}{k} P^k (1 - P)^{n-k}$
 - ▶ $E[k] = nP$
 - ▶ This binomial distribution peaks sharply around the mean:
 - ▶ $k \approx nP \Rightarrow \frac{k}{n}$ as an estimate for P
 - ▶ More accurate for larger n

Non-parametric density estimation

- We can estimate smoothed $p(\mathbf{x})$ by estimating P :
- Assumptions: $p(\mathbf{x})$ is continuous and the region \mathcal{R} enclosing \mathbf{x} is so small that p is near constant in it:

$$P = \int_{\mathcal{R}} p(\mathbf{x}') d\mathbf{x}' = p(\mathbf{x}) \times V$$
$$V = Vol(\mathcal{R})$$

$$\mathbf{x} \in \mathcal{R} \Rightarrow p(\mathbf{x}) = \frac{P}{V} \approx \frac{k/n}{V}$$

- Let V approach zero if we want to find $p(\mathbf{x})$ instead of the averaged version.

Necessary conditions for converge

- $p_n(\mathbf{x})$ is the estimate of $p(\mathbf{x})$ using n samples:

- V_n : the volume of region around \mathbf{x}
- k_n : the number of samples falling in the region

$$p_n(\mathbf{x}) = \frac{k_n/n}{V_n}$$

- Necessary conditions for converge of $p_n(\mathbf{x})$ to $p(\mathbf{x})$:
- $\lim_{n \rightarrow \infty} V_n = 0$
- $\lim_{n \rightarrow \infty} k_n = \infty$
- $\lim_{n \rightarrow \infty} k_n/n = 0$

Non-parametric density estimation: Main approaches

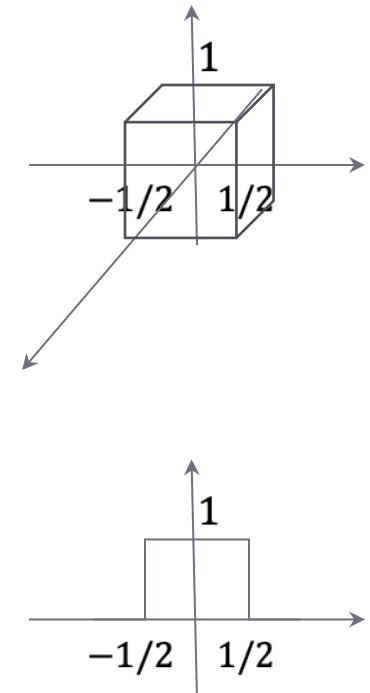
- Two approaches of satisfying conditions:
 - Kernel density estimator (Parzen window): fix V and determine K from the data
 - Number of points falling inside the volume can vary from point to point
 - converges to the true probability density in the limit $n \rightarrow \infty$ when V shrinks suitably with n (e.g., $V_n = V_1/\sqrt{n}$)
 - k-nearest neighbor density estimator: fix k and determine the value of V from the data
 - Volume grows until it contains k neighbors of x
 - converges to the true probability density in the limit $n \rightarrow \infty$ when k grows with n (e.g., $k_n = k_1\sqrt{n}$)

Parzen window

- Extension of histogram idea:
 - Hyper-cubes with length of side h (i.e., volume h^d) are located on the samples
- Hypercube as a simple window function:

$$\varphi(\mathbf{u}) = \begin{cases} 1 & (|u_1| \leq \frac{1}{2} \wedge \dots \wedge |u_d| \leq \frac{1}{2}) \\ 0 & \text{o.w.} \end{cases}$$

$$p_n(x) = \frac{1}{n} \times \frac{1}{h_n^d} \sum_{i=1}^n \varphi\left(\frac{\mathbf{x} - \mathbf{x}^{(i)}}{h_n}\right)$$



Parzen window

- Extension of histogram idea:

- Hyper-cubes with length of side h (i.e., volume h^d) are located on the samples
- Hypercube as a simple window function:

$$\varphi(\mathbf{u}) = \begin{cases} 1 & (|u_1| \leq \frac{1}{2} \wedge \dots \wedge |u_d| \leq \frac{1}{2}) \\ 0 & \text{o.w.} \end{cases}$$

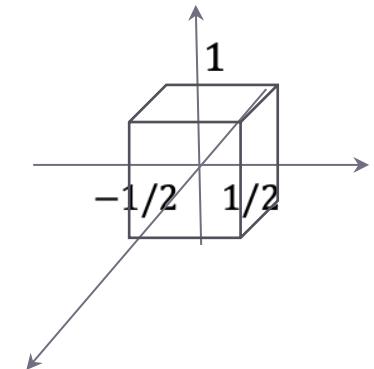
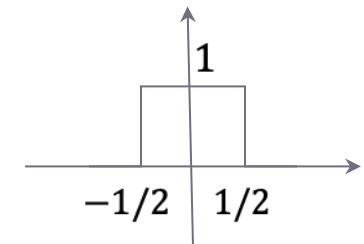
$$p_n(\mathbf{x}) = \frac{1}{n} \times \frac{1}{h_n^d} \sum_{i=1}^n \varphi\left(\frac{\mathbf{x} - \mathbf{x}^{(i)}}{h_n}\right)$$

- $p_n(\mathbf{x}) = \frac{k_n}{nV_n}$

- $k_n = \sum_{i=1}^n \varphi\left(\frac{\mathbf{x} - \mathbf{x}^{(i)}}{h_n}\right)$

- $V_n = (h_n)^d$

number of samples in the hypercube around \mathbf{x}



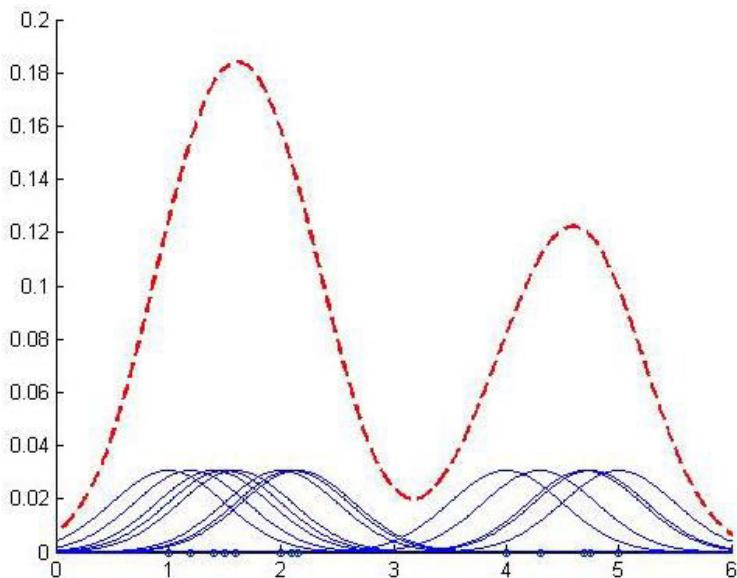
Window function

- Necessary conditions for window function to find legitimate density function:
 - $\varphi(x) \geq 0$
 - $\int \varphi(x)dx = 1$
- Windows are also called **kernels** or potential functions.

Density estimation: non-parametric

$$\hat{p}_n(x) = \frac{1}{n} \sum_{i=1}^n N(x|x^{(i)}, h^2) \rightarrow \frac{1}{\sqrt{2\pi}h} e^{-\frac{(x-x^{(i)})^2}{2h^2}}$$

1 1.2 1.4 1.5 1.6 2 2.1 2.15 4 4.3 4.7 4.75 5

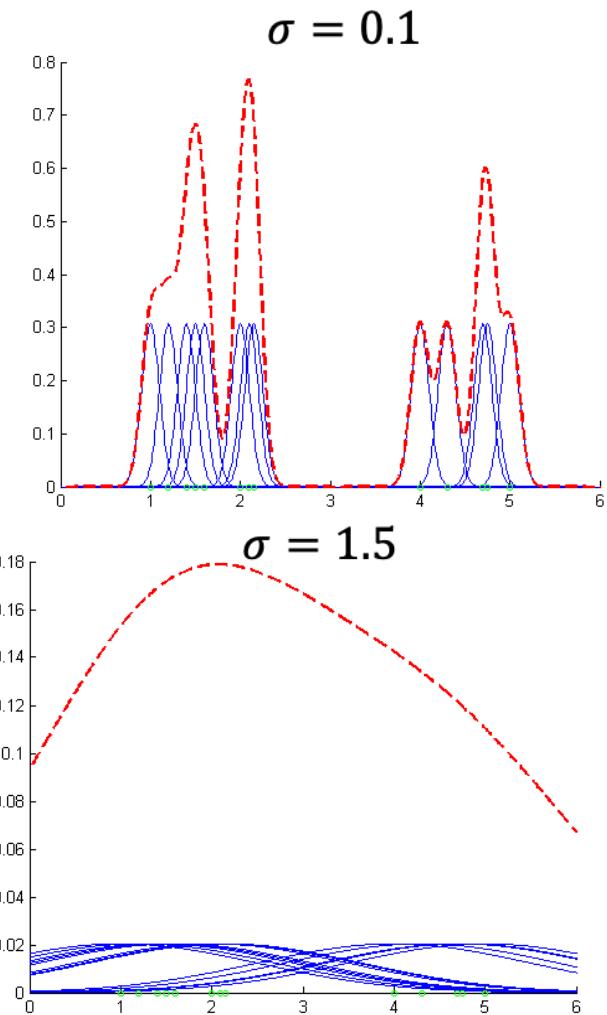
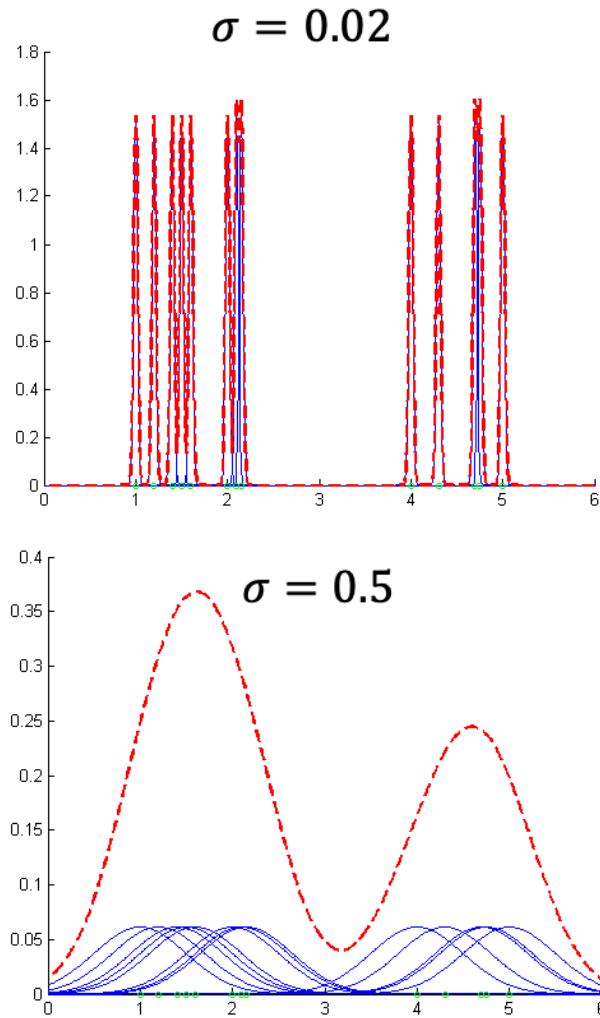


$$\sigma = h$$

$$\begin{aligned}\hat{p}(x) &= \frac{1}{n} \sum_{i=1}^n N(x|x^{(i)}, \sigma^2) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-x^{(i)})^2}{2\sigma^2}}\end{aligned}$$

Choice of σ is crucial.

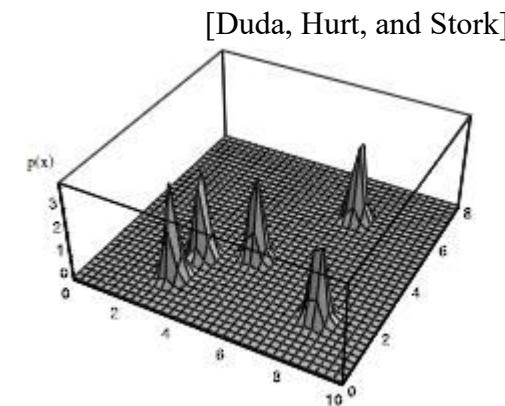
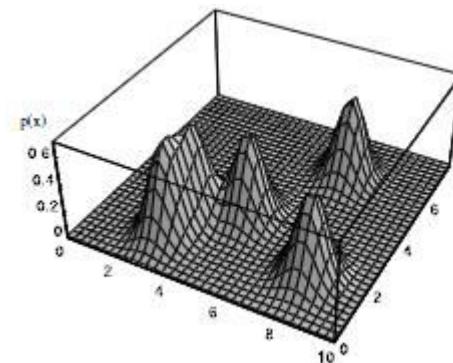
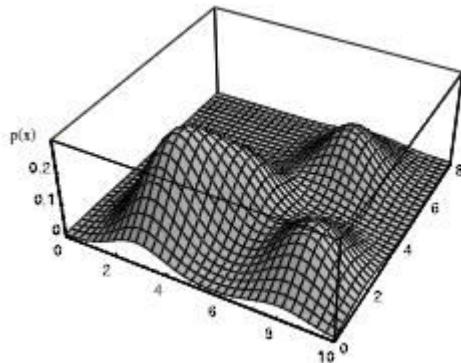
Density estimation: non-parametric



Window (or kernel) function: Width parameter

- $$p_n(x) = \frac{1}{n} \times \frac{1}{h_n^d} \sum_{i=1}^n \varphi\left(\frac{x - x^{(i)}}{h_n}\right)$$

- Choosing h_n :
 - Too large: low resolution
 - Too small: much variability



- For unlimited n , by letting V_n slowly approach zero as n increases $p_n(x)$ converges to $p(x)$

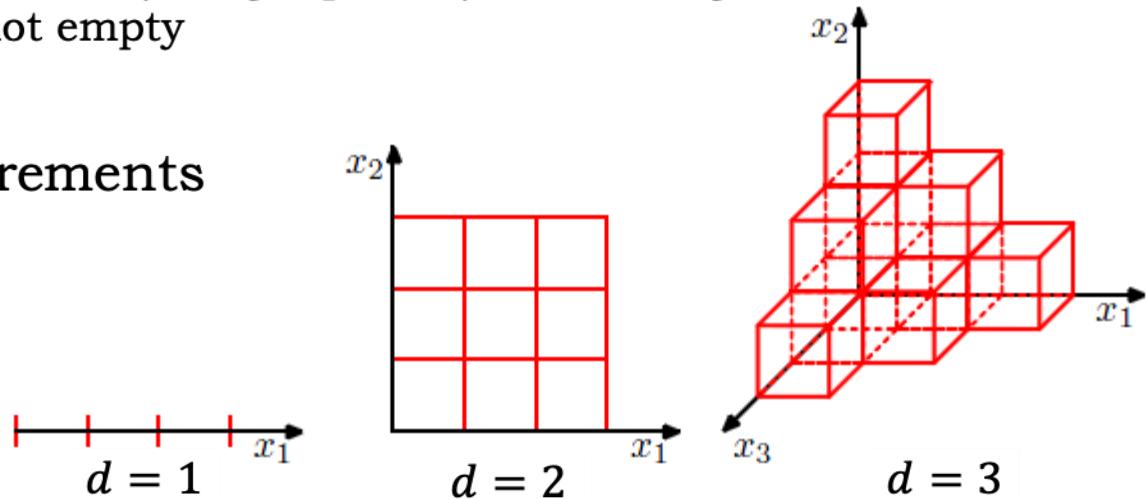
Width parameter

- ▶ For fixed n , a smaller h results in higher variance while a larger h leads to higher bias.
- ▶ For a fixed h , the variance decreases as the number of sample points n tends to infinity
 - ▶ for a large enough number of samples, the smaller h the better the accuracy of the resulting estimate
- ▶ In practice, where only a finite number of samples is possible, a compromise between h and n must be made.
 - ▶ h can be set using techniques like cross-validation where the density estimation used for learning tasks such as classification

Practical issues: Curse of dimensionality

- ▶ Large n is necessary to find an acceptable density estimation in high dimensional feature spaces
- ▶ n must grow exponentially with the dimensionality d .
 - ▶ If n equidistant points are required to densely fill a one-dim interval, n^d points are needed to fill the corresponding d -dim hypercube.
 - We need an exponentially large quantity of training data to ensure that the cells are not empty

- ▶ Also complexity requirements



Non-parametric density estimation: Main approaches

- ▶ Two approaches of satisfying conditions:
 - ▶ Kernel density estimator (Parzen window): fix V and determine K from the data
 - ▶ Number of points falling inside the volume can vary from point to point
 - ▶ converges to the true probability density in the limit $n \rightarrow \infty$ when V shrinks suitably with n (e.g., $V_n = V_1/\sqrt{n}$)
 - ▶ k-nearest neighbor density estimator: fix k and determine the value of V from the data
 - ▶ Volume grows until it contains k neighbors of x
 - ▶ converges to the true probability density in the limit $n \rightarrow \infty$ when k grows with n (e.g., $k_n = k_1\sqrt{n}$)

k_n -nearest neighbor estimation

- ▶ Cell volume is a function of the point location
 - ▶ To estimate $p(x)$, let the cell around x grow until it captures k_n samples called k_n nearest neighbors of x .
- ▶ Two possibilities can occur:
 - ▶ high density near $x \Rightarrow$ cell will be small which provides a good resolution
 - ▶ low density near $x \Rightarrow$ cell will grow large and stop until higher density regions are reached

k_n -nearest neighbor estimation

-

$$p_n(\mathbf{x}) = \frac{k_n/n}{V_n} \Rightarrow V_n \approx \frac{1}{p(\mathbf{x})} \times k_n/n$$

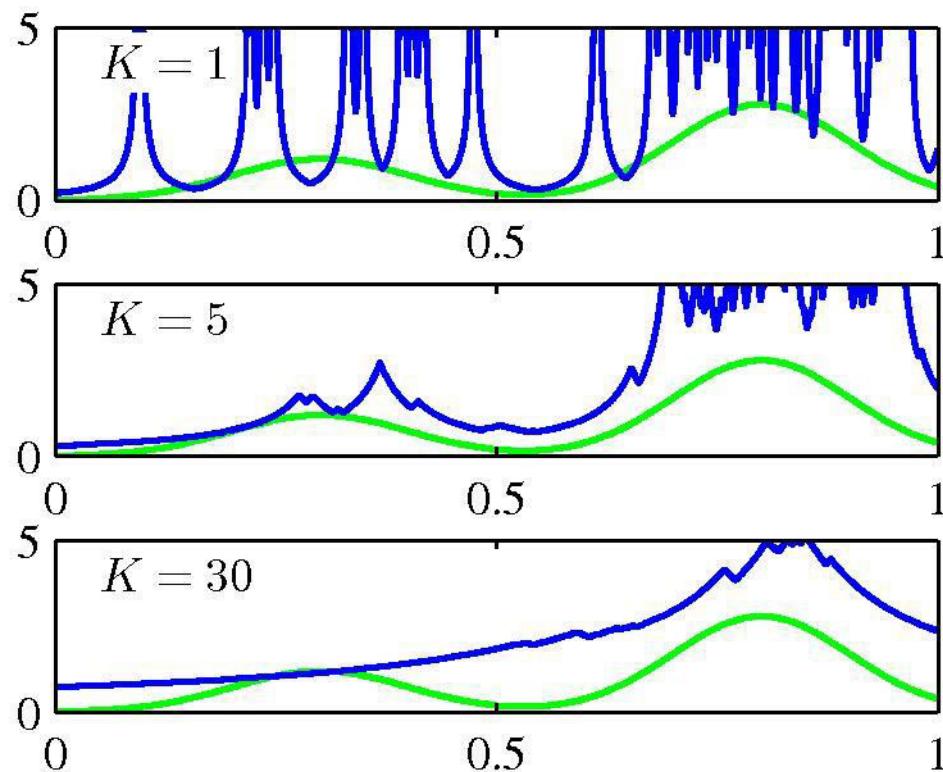
- A family of estimates by setting $k_n = k_1\sqrt{n}$ and choosing different values for k_1 :

$$p_n(\mathbf{x}) = \frac{k_n/n}{V_n} \Rightarrow V_n \approx \frac{1/p(\mathbf{x})}{\sqrt{n}} \quad k_1 = 1$$

↓
 V_n is a function of \mathbf{x}

k_n -Nearest Neighbor Estimation: Example

- Discontinuities in the slopes

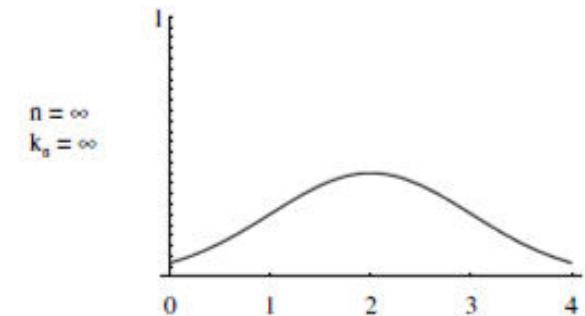
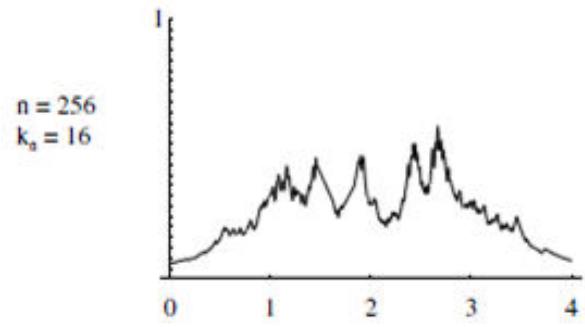
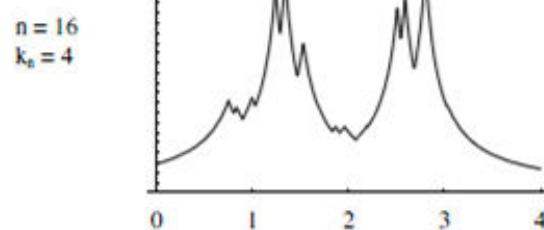
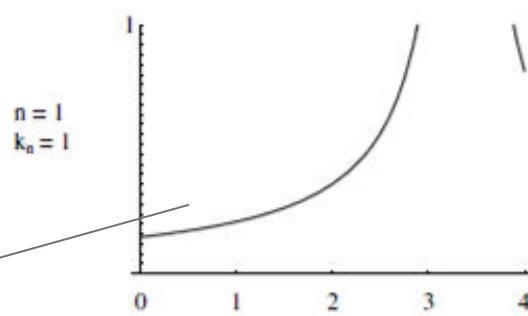


[Bishop]

k_n -Nearest Neighbor Estimation: Example

- $k_n = \sqrt{n}$

$$p_1(x) = \frac{1}{2|x - x^{(1)}|}$$



[Duda, Hurt, and Stork]

Non-parametric density estimation: Summary

- Generality of distributions
 - ▶ With enough samples, convergence to an arbitrarily complicated target density can be obtained.
- ▶ The number of required samples must be very large to assure convergence
 - ▶ grows exponentially with the dimensionality of the feature space
- ▶ These methods are very sensitive to the choice of window width or number of nearest neighbors
- ▶ There may be severe requirements for computation time and storage (needed to save all training samples).
 - ▶ ‘training’ phase simply requires storage of the training set
 - ▶ computational cost of evaluating $p(x)$ grows linearly with the size of the data set

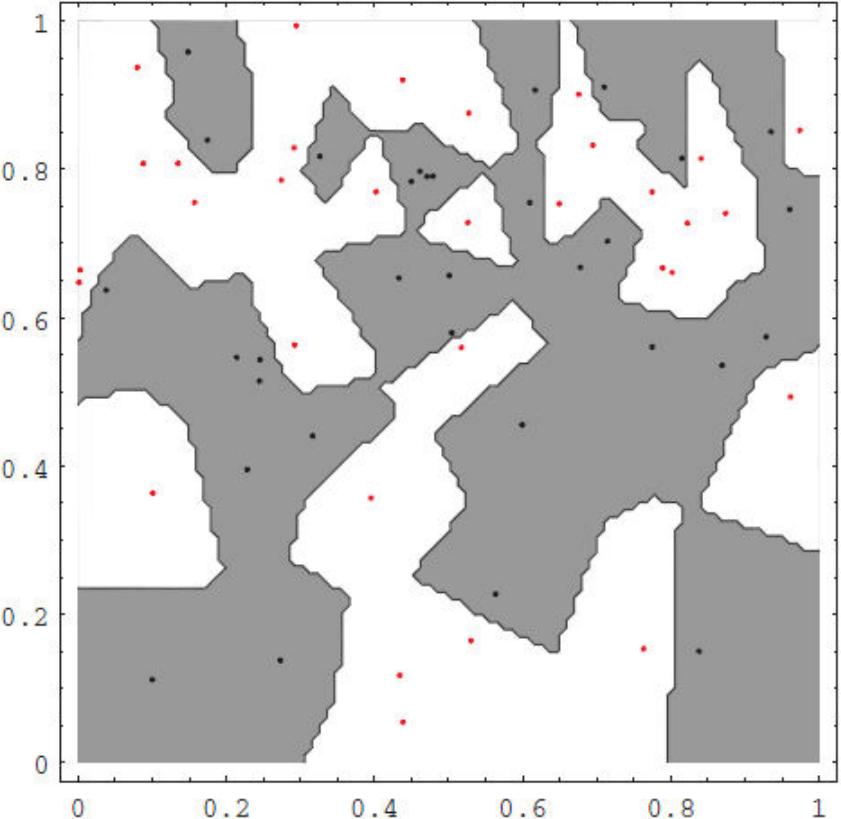
Nonparametric learners

- **Memory-based or instance-based** learners
 - lazy learning: (almost) all the work is done at the test time.
- Generic description:
 - Memorize training $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})$.
 - Given test \mathbf{x} predict: $\hat{y} = f(\mathbf{x}; \mathbf{x}^{(1)}, y^{(1)}, \dots, \mathbf{x}^{(n)}, y^{(n)})$.
- f is typically expressed in terms of the similarity of the test sample \mathbf{x} to the training samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$

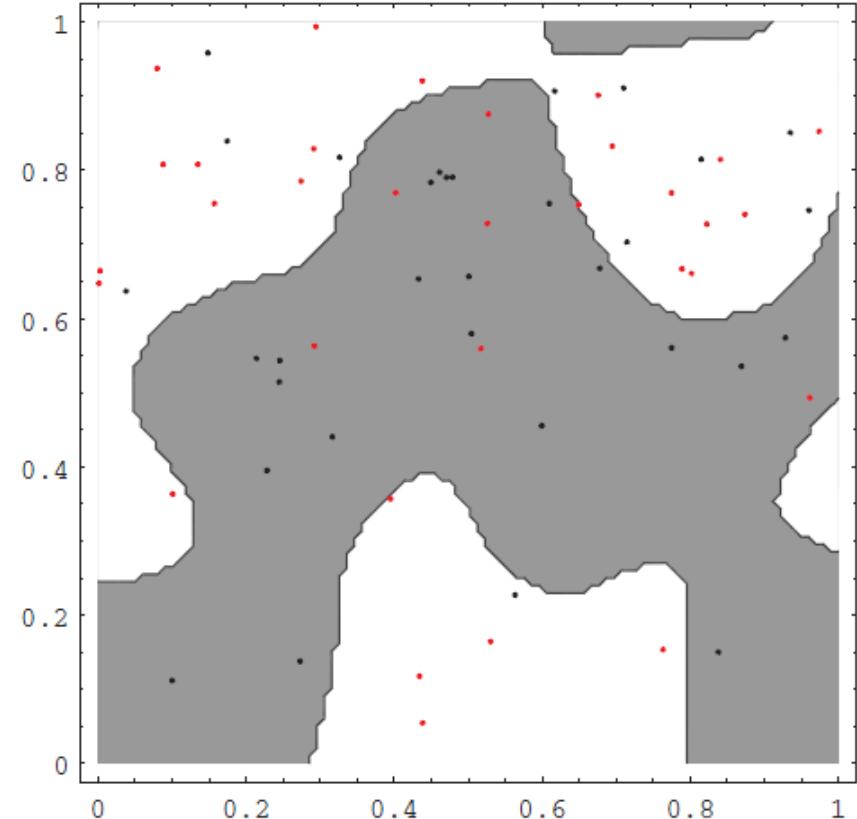
Parzen window & generative classification

- If $\frac{\frac{1}{n_1} \times \frac{1}{h^d} \sum_{x^{(i)} \in \mathcal{D}_1} \varphi\left(\frac{x-x^{(i)}}{h}\right)}{\frac{1}{n_2} \times \frac{1}{h^d} \sum_{x^{(i)} \in \mathcal{D}_2} \varphi\left(\frac{x-x^{(i)}}{h}\right)} \times \frac{P(\mathcal{C}_1)}{P(\mathcal{C}_2)} > 1$ decide \mathcal{C}_1
otherwise decide \mathcal{C}_2
- ▶ $n_j = |\mathcal{D}_j|$ ($j = 1, 2$): number of training samples in class \mathcal{C}_j
 - ▶ \mathcal{D}_j : set of training samples labels as \mathcal{C}_j
- ▶ For large n , it needs both high time and memory requirements

Parzen window & generative classification: Example



Smaller h



larger h

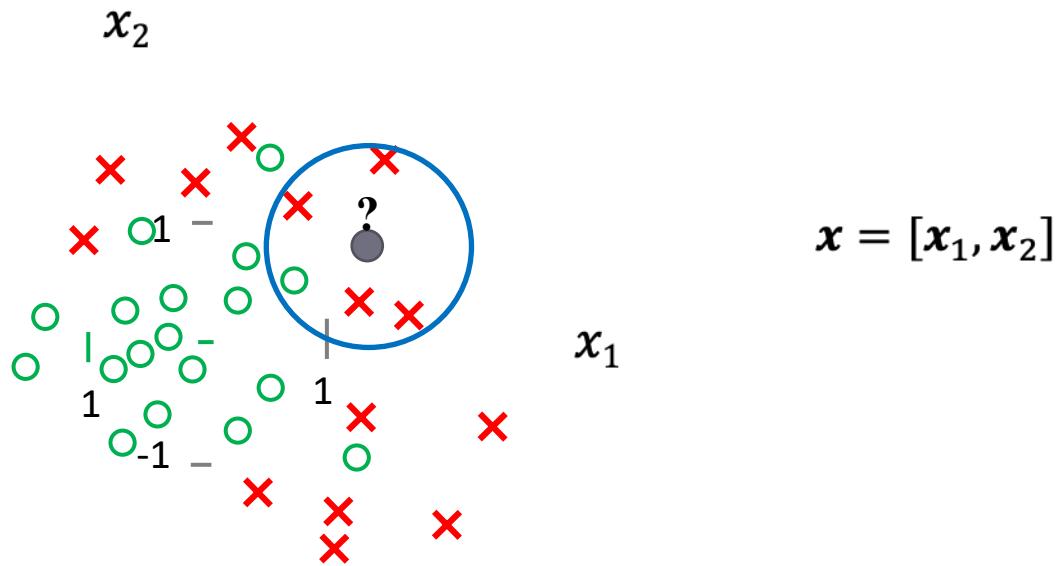
[Duda, Hurt, and Stork]

Instance-based Learning

Classification: k-Nearest-Neighbor (kNN)

- **k-NN classifier: $k > 1$ nearest neighbors**
 - Label for x predicted by majority voting among its k-NN.

- $k = 5$



- What is the effect of k ?

kNN classifier

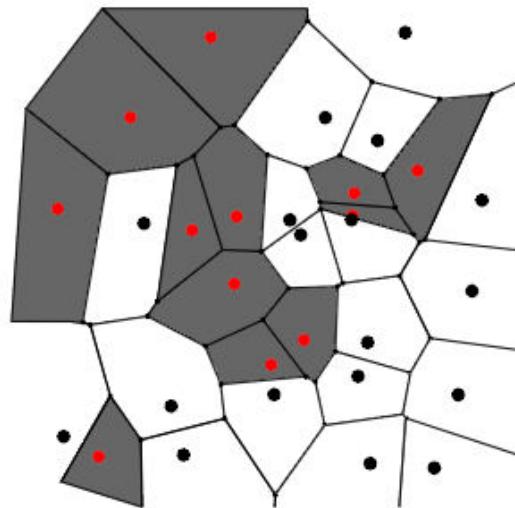
- Given
 - Training data $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$ are simply stored.
 - Test sample: \mathbf{x}
- To classify \mathbf{x} :
 - Find k nearest training samples to \mathbf{x}
 - Out of these k samples, identify the number of samples k_j belonging to class \mathcal{C}_j ($j = 1, \dots, C$).
 - Assign \mathbf{x} to the class \mathcal{C}_{j^*} where $j^* = \operatorname{argmax}_{j=1,\dots,C} k_j$
- It can be considered as a discriminative method.

Probabilistic perspective of kNN

- ▶ kNN as a discriminative nonparametric classifier
 - ▶ Non-parametric density estimation for $P(\mathcal{C}_j|x)$
 - ▶ $P(\mathcal{C}_j|x) \approx \frac{k_j}{k}$ where k_j shows the number of training samples among k nearest neighbors of x that are labeled \mathcal{C}_j
 - ▶ Bayes decision rule for assigning labels

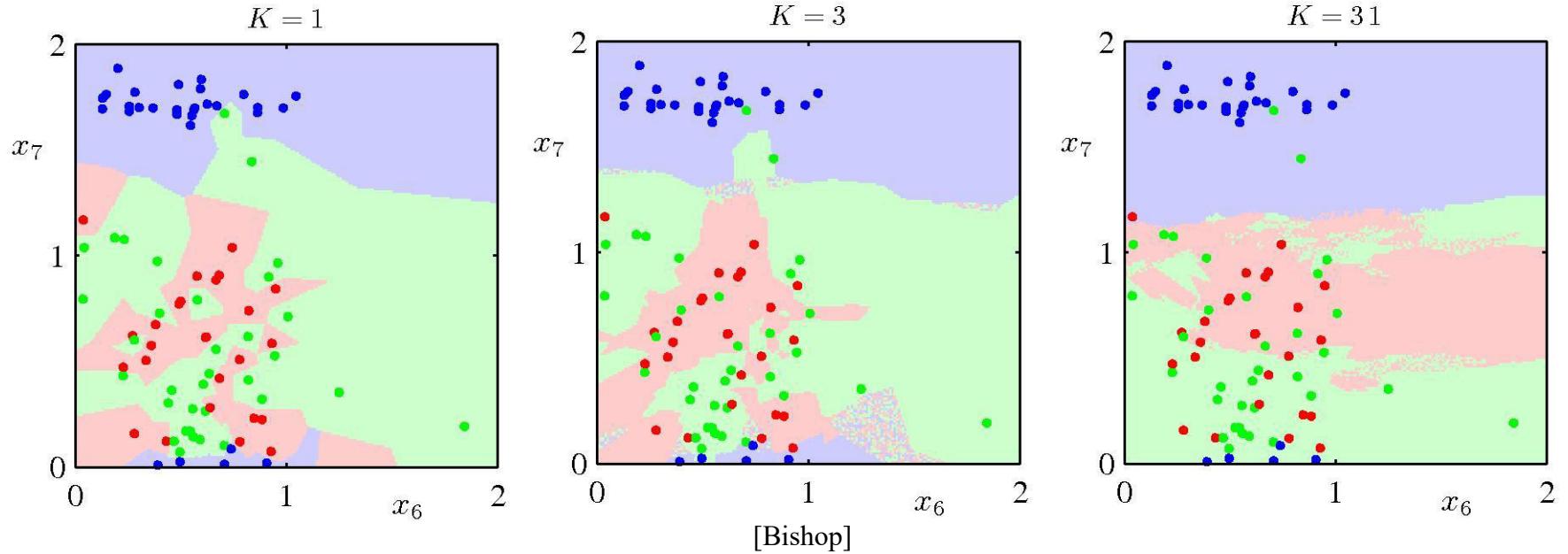
Nearest-neighbor classifier: Example

- Voronoi tessellation:
 - Each cell consists of all points closer to a given training point than to any other training points
 - All points in a cell are labeled by the category of the corresponding training point.



[Duda, Hurt, and Strok's Book]

kNN classifier: Effect of k



Need to determine an appropriate value for k (e.g., by crossvalidation)

Nearest neighbor classifier: error bound

- ▶ Nearest-Neighbor: kNN with $k = 1$
 - ▶ Decision rule: $\hat{y} = y^{NN(\mathbf{x})}$ where $NN(\mathbf{x}) = \operatorname{argmin}_{i=1,\dots,N} \|\mathbf{x} - \mathbf{x}^{(i)}\|$
 - ▶ Cover & Hart 67: asymptotic risk of NN classifier satisfies:

$$R^* \leq R_{\infty}^{NN} \leq 2R^*(1 - R^*) \leq 2R^*$$

R_n : expected risk of NN classifier with n training examples drawn from $p(\mathbf{x}, y)$

$$R_{\infty}^{NN} = \lim_{n \rightarrow \infty} R_n^{NN}$$

R^* : the optimal Bayes risk

k-NN classifier: error bound

- Devr 96: the asymptotic risk of the kNN classifier $R_\infty = \lim_{n \rightarrow \infty} R_n$ satisfies

$$R^* \leq R_\infty^{kNN} \leq R^* + \sqrt{\frac{2R_\infty^{NN}}{k}}$$

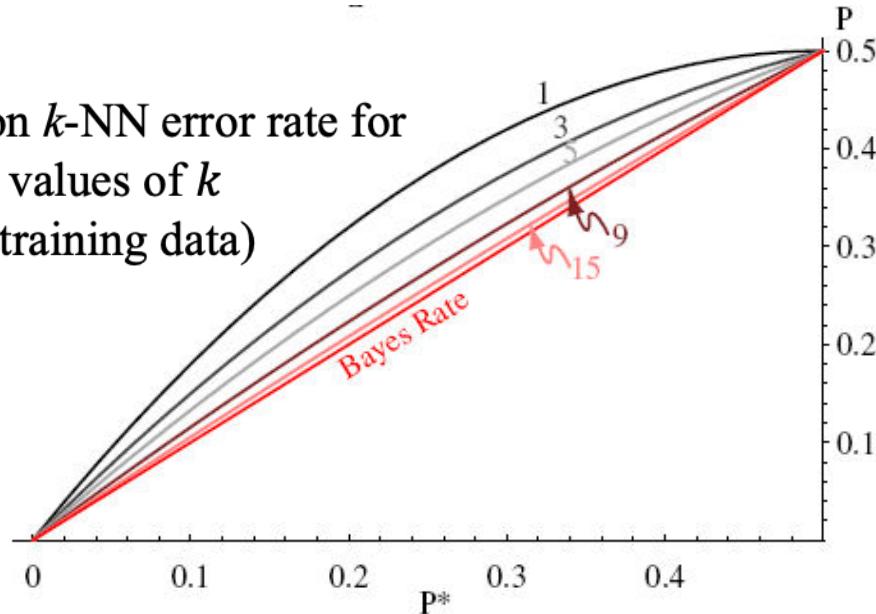
where R^* is the optimal Bayes risk.

- As k increases, the upper bounds get closer to the lower bound (the Bayes Error rate).
 - When $k \rightarrow \infty$, the two bounds meet and k-NN rule becomes optimal.

Bound on the k-Nearest Neighbor Error Rate

-

bounds on k -NN error rate for different values of k
(infinite training data)



[Duda, Hurt, and Strok]

- As k increases, the upper bounds get closer to the lower bound (the Bayes Error rate).
 - When $k \rightarrow \infty$, the two bounds meet and k -NN rule becomes optimal.

Instance-based learner

- Main things to construct an instance-based learner:
 - A distance metric
 - Number of nearest neighbors of the test data that we look at
 - A weighting function (optional)
 - How to find the output based on neighbors?

Distance measures

- Euclidean distance

$$d(\mathbf{x}, \mathbf{x}') = \sqrt{\|\mathbf{x} - \mathbf{x}'\|_2^2} = \sqrt{(x_1 - x'_1)^2 + \cdots + (x_d - x'_d)^2}$$

Sensitive to irrelevant features

- Distance learning methods for this purpose

- Weighted Euclidean distance

- $d_w(\mathbf{x}, \mathbf{x}') = \sqrt{w_1(x_1 - x'_1)^2 + \cdots + w_d(x_d - x'_d)^2}$

- Mahalanobis distance

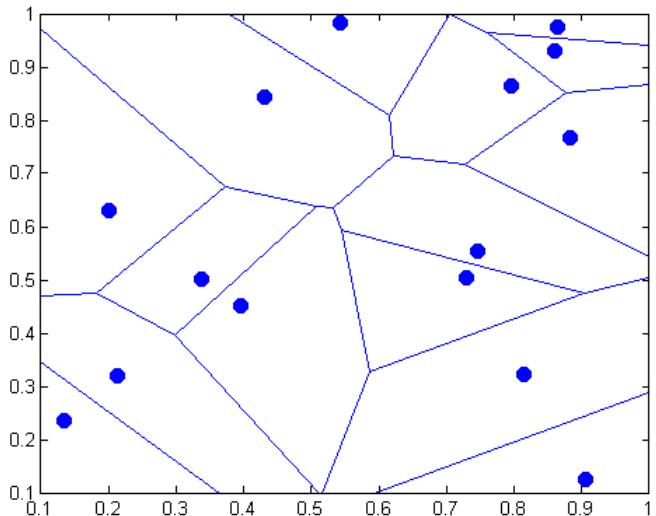
- $d_A(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \mathbf{A} (\mathbf{x} - \mathbf{x}')}}$

- Other distances:

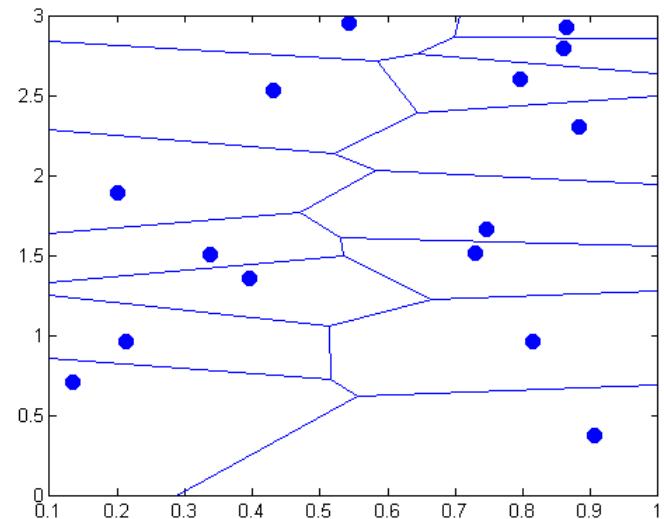
- Hamming, angle, ...

- $L_p(\mathbf{x}, \mathbf{x}') = \sqrt[p]{\sum_{i=1}^d (x_i - x'_i)^p}$

Distance measure: example

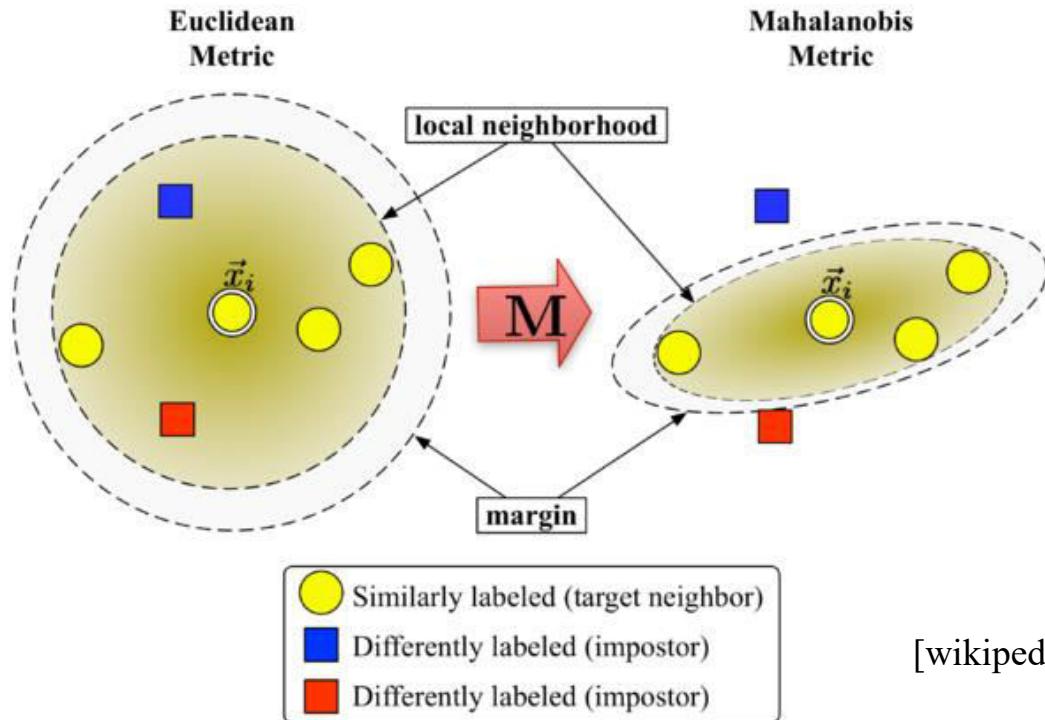


$$d(\mathbf{x}, \mathbf{x}') = \sqrt{(x_1 - x'_1)^2 + (x_2 - x'_2)^2}$$



$$d(\mathbf{x}, \mathbf{x}') = \sqrt{(x_1 - x'_1)^2 + 3(x_2 - x'_2)^2}$$

Metric Learning Example: LMNN



$$\begin{aligned} & \min_{\mathbf{M}} \sum_{i,j \in N_i} d(\vec{x}_i, \vec{x}_j) + \sum_{i,j,l} \xi_{ijl} \\ & \forall i,j \in N_i, l, y_l \neq y_i \\ & d(\vec{x}_i, \vec{x}_j) + 1 \leq d(\vec{x}_i, \vec{x}_l) + \xi_{ijl} \\ & \xi_{ijl} \geq 0 \\ & \mathbf{M} \succeq 0 \end{aligned}$$

[wikipedia]

$$d^2(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}')$$

Weighted kNN classification

- Weight nearer neighbors more heavily:

$$\hat{y} = f(\mathbf{x}) = \operatorname{argmax}_{c=1,\dots,C} \sum_{j \in N_k(\mathbf{x})} w_j(\mathbf{x}) \times I(c = y^{(j)})$$

$$w_j(\mathbf{x}) = \frac{1}{\|\mathbf{x} - \mathbf{x}^{(j)}\|^2}$$

An example of
weighting function

- In the weighted kNN, we can use all training examples instead of just k (Stepard's method):

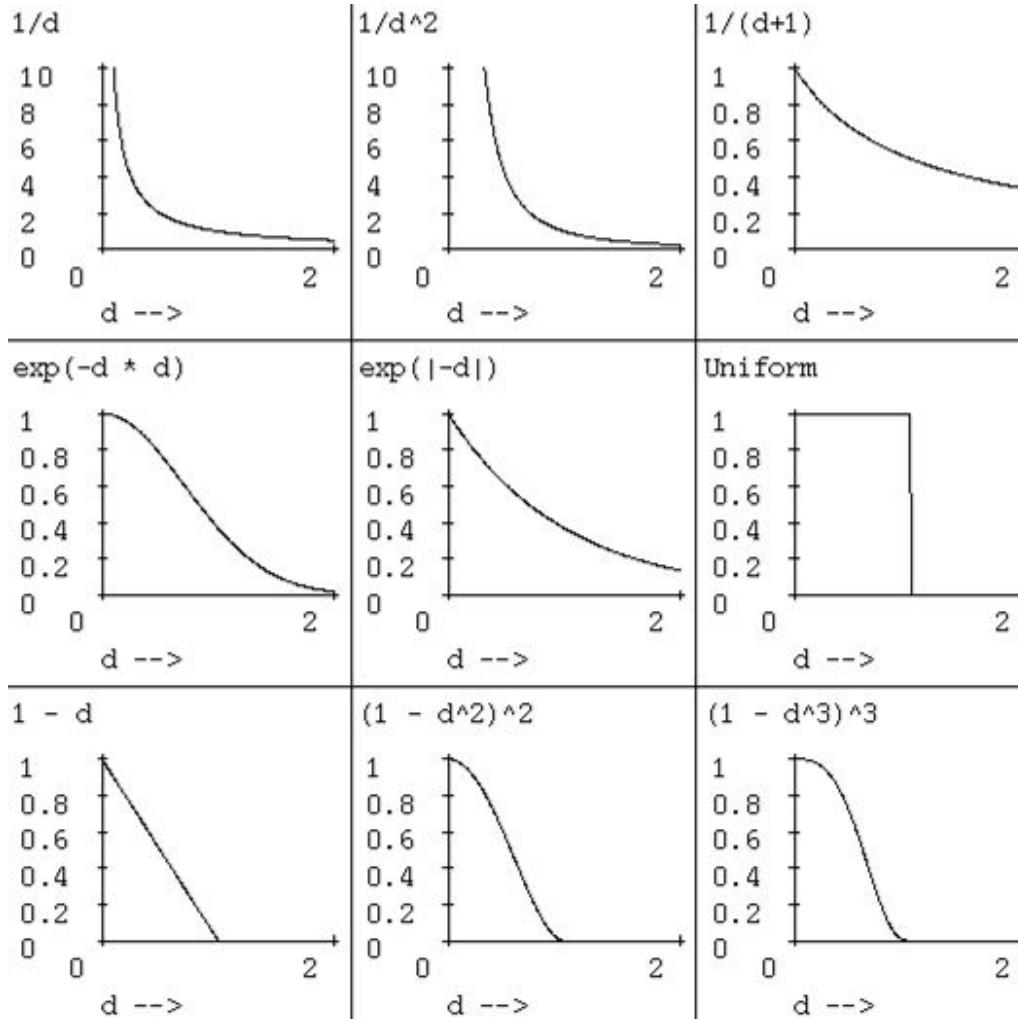
$$\hat{y} = f(\mathbf{x}) = \operatorname{argmax}_{c=1,\dots,C} \sum_{j=1}^n w_j(\mathbf{x}) \times I(c = y^{(j)})$$

- Weights can be found using a kernel function $w_j(\mathbf{x}) = K(\mathbf{x}, \mathbf{x}^{(j)})$:

- e.g., $K(\mathbf{x}, \mathbf{x}^{(j)}) = e^{-\frac{d(\mathbf{x}, \mathbf{x}^{(j)})}{\sigma^2}}$

Weighting functions

$$d = d(\mathbf{x}, \mathbf{x}')$$



[Fig. has been adopted from Andrew Moore's tutorial on "Instance-based learning"]

kNN regression

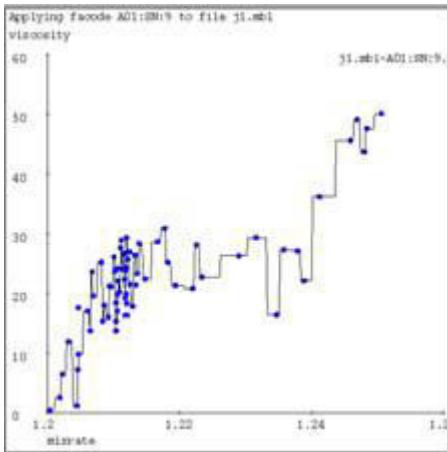
- Simplest k-NN regression:
 - Let $x'^{(1)}, \dots, x'^{(k)}$ be the k nearest neighbors of x and $y'^{(1)}, \dots, y'^{(k)}$ be their labels.

$$\hat{y} = \frac{1}{k} \sum_{j \in N_k(x)} y'^{(j)}$$

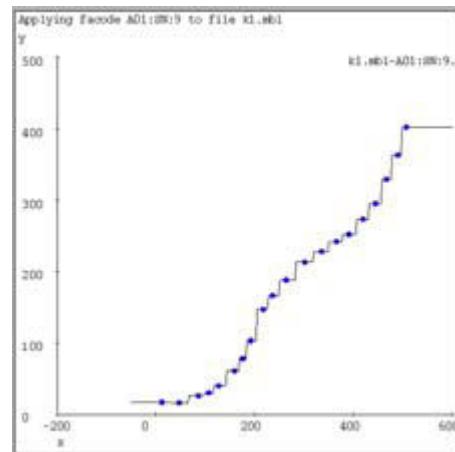
- Problems of kNN regression for fitting functions:
 - Problem 1: Discontinuities in the estimated function
 - Solution: Weighted (or kernel) regression
 - INN: noise-fitting problem
 - kNN ($k > 1$) smoothes away noise, but there are other deficiencies.
 - flats the ends

kNN regression: examples

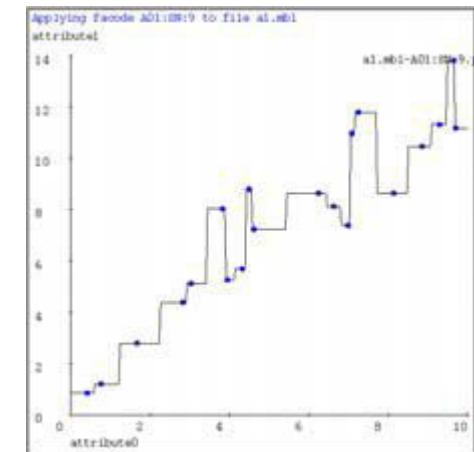
Dataset 1



Dataset 2



Dataset 3



$k = 1$

$k = 9$

Instance-based Learning

Weighted (or kernel) kNN regression

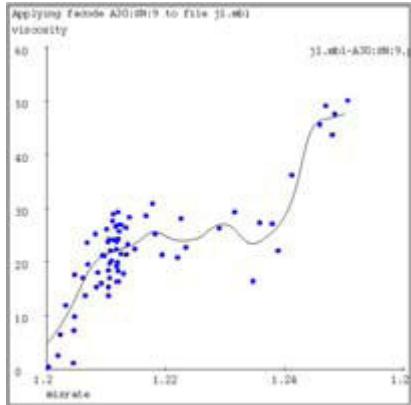
- Higher weights to nearer neighbors:

$$\hat{y} = f(\mathbf{x}) = \frac{\sum_{j \in N_k(\mathbf{x})} w_j(\mathbf{x}) \times y^{(j)}}{\sum_{j \in N_k(\mathbf{x})} w_j(\mathbf{x})}$$

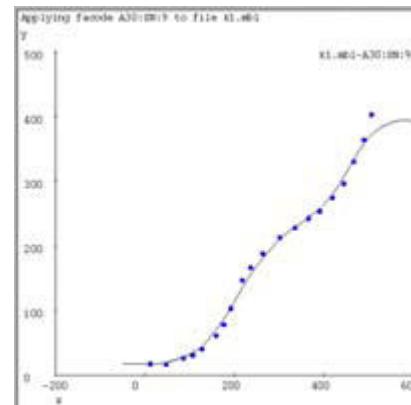
- In the weighted kNN regression, we can use all training examples instead of just k in the weighted form:

$$\hat{y} = f(\mathbf{x}) = \frac{\sum_{j=1}^n w_j(\mathbf{x}) \times y^{(j)}}{\sum_{j=1}^n w_j(\mathbf{x})}$$

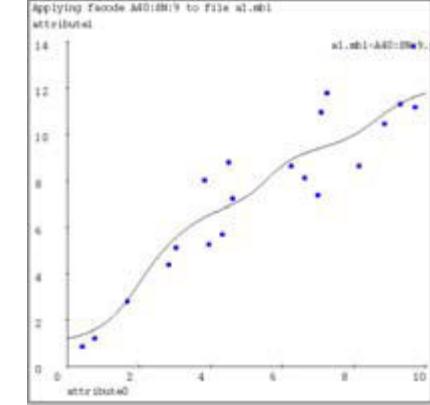
Kernel kNN regression



$$\sigma = \frac{1}{32} \text{ of x-axis width}$$



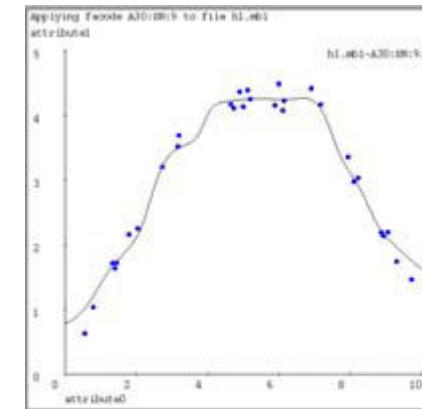
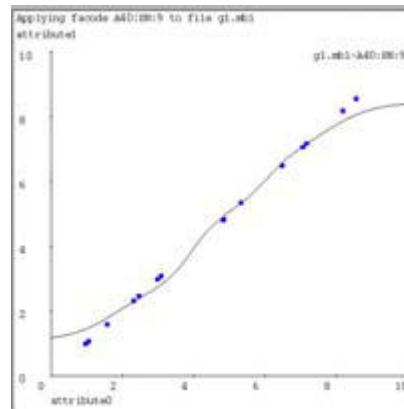
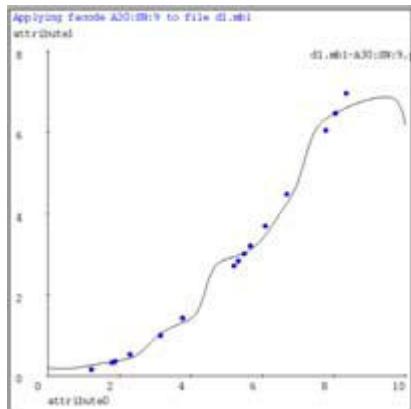
$$\sigma = \frac{1}{32} \text{ of x-axis width}$$



$$\sigma = \frac{1}{16} \text{ of x-axis width}$$

- Choosing a good parameter (kernel width) is important.

Kernel kNN regression



In these datasets, some regions are without samples
Best kernel widths have been used

- Disadvantages:
 - not capturing the simple structure of the data
 - failure to extrapolate at edges

[Figs. have been adopted from Andrew Moore's tutorial on “Instance-based learning”]

Locally weighted linear regression

- For each test sample, it produces linear approximation to the target function in a local region
- Instead of finding the output using weighted averaging (as in the kernel regression), we fit a parametric function locally:

$$\hat{y} = f(\mathbf{x}, \mathbf{x}^{(1)}, y^{(1)}, \dots, \mathbf{x}^{(n)}, y^{(n)})$$

$$\hat{y} = f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_d x_d$$

$$J(\mathbf{w}) = \sum_{i \in N_k(\mathbf{x})} (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

Unweighted linear



\mathbf{w} is found for each test sample

Locally weighted linear regression

-

$$\hat{y} = f(\mathbf{x}, \mathbf{x}^{(1)}, y^{(1)}, \dots, \mathbf{x}^{(n)}, y^{(n)})$$

$$J(\mathbf{w}(\mathbf{x})) = \sum_{i \in N_k(\mathbf{x})} (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

unweighted

$$J(\mathbf{w}(\mathbf{x})) = \sum_{i \in N_k(\mathbf{x})} K(\mathbf{x}, \mathbf{x}^{(i)}) (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

e.g., $K(\mathbf{x}, \mathbf{x}^{(i)}) = e^{-\frac{\|\mathbf{x}-\mathbf{x}^{(i)}\|^2}{2\sigma^2}}$

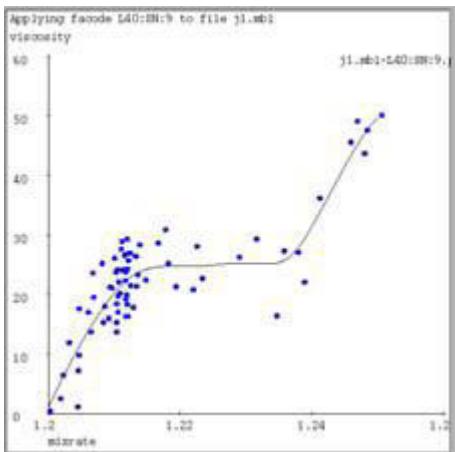
weighted

$$J(\mathbf{w}(\mathbf{x})) = \sum_{i=1}^n K(\mathbf{x}, \mathbf{x}^{(i)}) (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

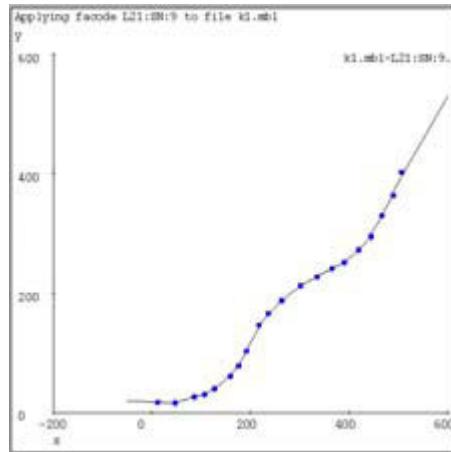


Weighted on all training examples

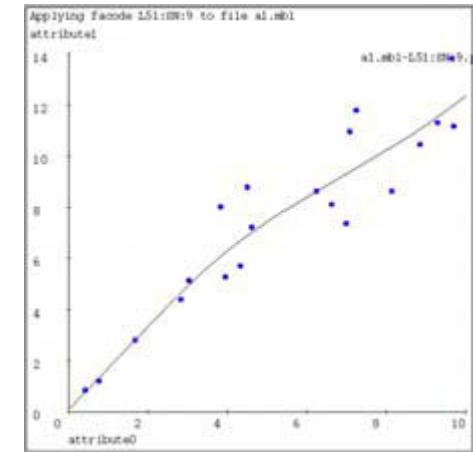
Locally weighted linear regression: example



$$\sigma = \frac{1}{16} \text{ of x-axis width}$$



$$\sigma = \frac{1}{32} \text{ of x-axis width}$$



$$\sigma = \frac{1}{8} \text{ of x-axis width}$$

- More proper result than weighted kNN regression

Locally weighted regression: summary

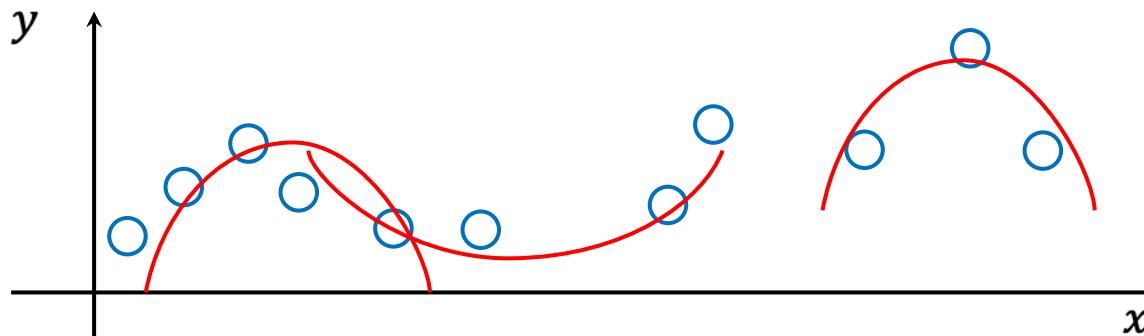
- Idea 1: weighted kNN regression

- using the weighted average on the output of x 's neighbors (or on the outputs of all training data):

$$\hat{y} = \frac{\sum_{i=1}^k y^{(i)} K(x, x^{(i)})}{\sum_{j=1}^k K(x, x^{(j)})}$$

- Idea 2: Locally weighted parametric regression

- Fit a parametric model (e.g. linear function) to the neighbors of x (or on all training data).
- Implicit assumption: the target function is reasonably smooth.



Parametric vs. nonparametric methods

- Is SVM classifier parametric?

$$\hat{y} = \text{sign}(w_0 + \sum_{\alpha_i > 0} \alpha_i y^{(i)} K(\mathbf{x}, \mathbf{x}^{(i)}))$$

- In general, we can not summarize it in a simple parametric form.
 - Need to keep around support vectors (possibly all of the training data).
- However, α_i are kind of parameters that are found in the training phase

Instance-based learning: summary

- Learning is just storing the training data
 - prediction on a new data based on the training data themselves
- An instance-based learner does not rely on assumption concerning the structure of the underlying density function.
- With large datasets, instance-based methods are slow for prediction on the test data
 - kd-tree, Locally Sensitive Hashing (LSH), and other kNN approximations can help.

Reference

- T. Mitchell, “Machine Learning”, 1998. [Chapter 8]
- C.M. Bishop, “Pattern Recognition and Machine Intelligence”, Section 2.5.