

# L8: Nearest neighbors

**Nearest neighbors density estimation**

**The k nearest neighbors classification rule**

**kNN as a lazy learner**

**Characteristics of the kNN classifier**

**Optimizing the kNN classifier**

# Non-parametric density estimation: review

## Recap

- As we saw in the previous lecture that the general expression for non-parametric density estimation is

$$p(x) \cong \frac{k}{NV} \text{ where } \begin{cases} V & \text{volume surrounding } x \\ N & \text{total \#examples} \\ k & \text{\#examples inside } V \end{cases}$$

- At that time, we mentioned that this estimate could be computed by
  - Fixing  $V$  and determining the number  $k$  of data points inside  $V$ 
    - This is the approach used in **kernel density estimation**
  - Fixing  $k$  and determining the minimum volume  $V$  that encompasses  $k$  points in the dataset
    - This gives rise to the **k-nearest-neighbors** (kNN) approach, which is the subject of this lecture

# kNN density estimation

## Approach

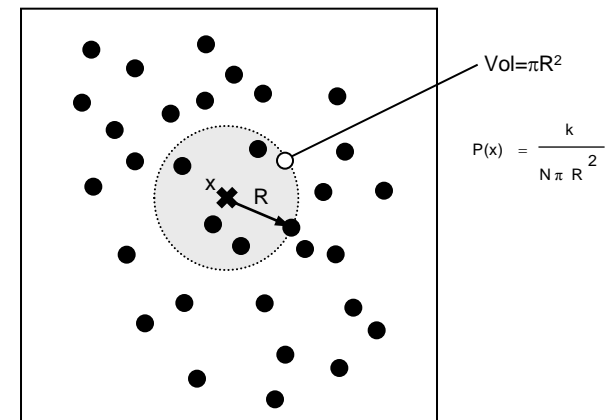
- In kNN we grow the volume surrounding the estimation point  $x$  until it encloses a total of  $k$  data points
- The density estimate then becomes

$$p(x) \cong \frac{k}{NV} = \frac{k}{N c_D R_k^D(x)}$$

- where  $R_k^D(x)$  is the distance between the estimation point  $x$  and its  $k$ -th closest neighbor
- and  $c_D$  is the volume of the unit sphere in  $D$  dimensions, which is equal to (Bishop, 1995)

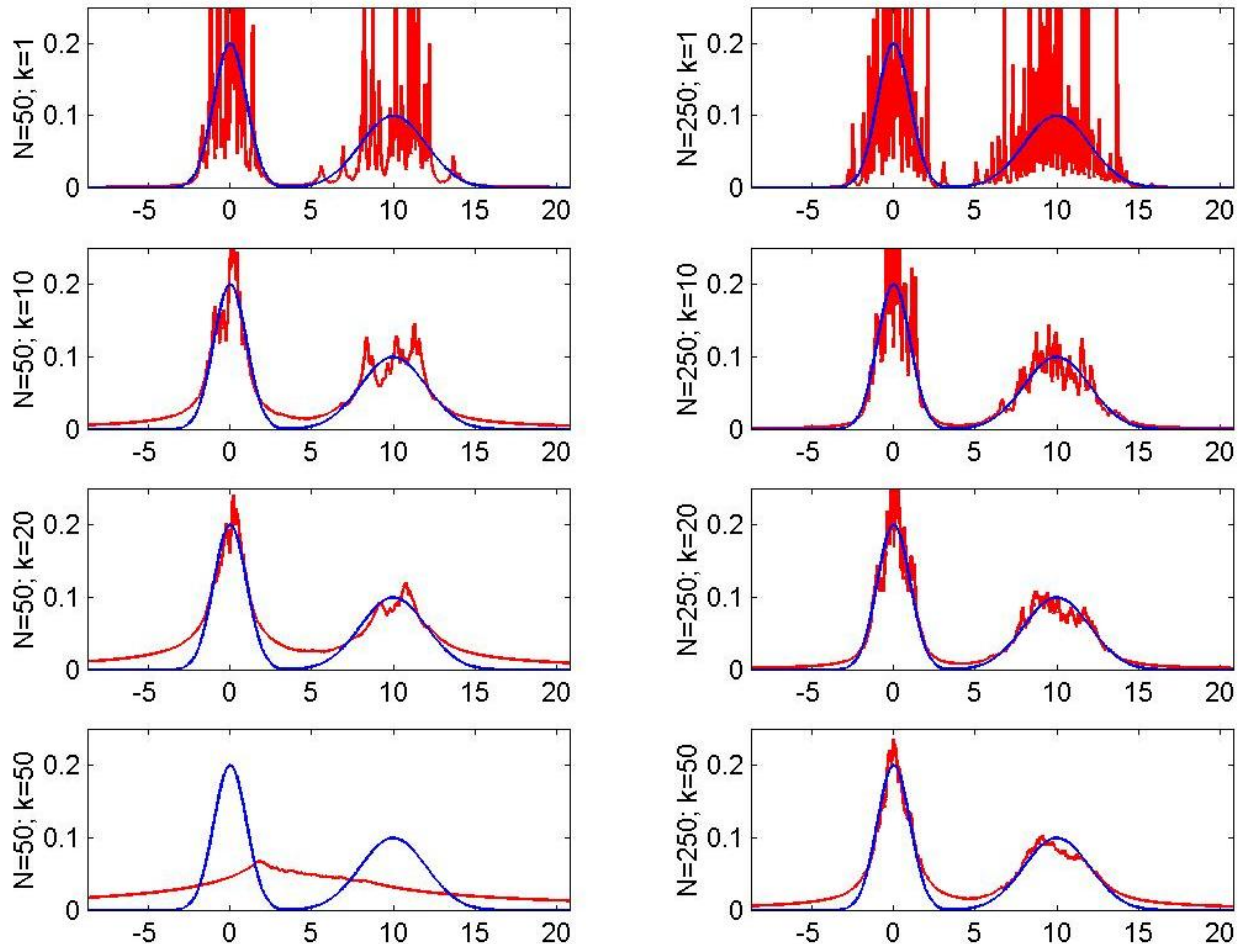
$$c_D = \frac{2\pi^{D/2}}{D \cdot \Gamma(D/2)}$$

- Thus  $c_1 = 2$ ,  $c_2 = \pi$ ,  $c_3 = 4\pi/3$ , and so on



- In general, the estimates that can be obtained with the kNN method are not very satisfactory
  - The estimates are prone to local noise
  - The method produces estimates with very heavy tails
  - Since the function  $R_k(x)$  is not differentiable, the density estimate will have discontinuities
  - The resulting density is not a true probability density since its integral over all the sample space diverges
- These properties are illustrated in the next few slides

- To illustrate the behavior of kNN we generated several density estimates for a bimodal Gaussian:  $p(x) = \frac{1}{2}N(0,1) + \frac{1}{2}N(10,4)$



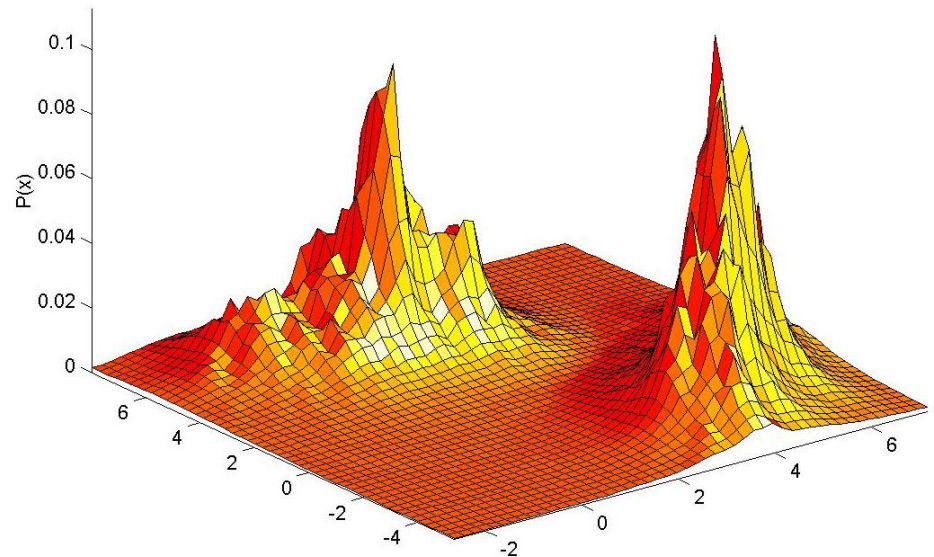
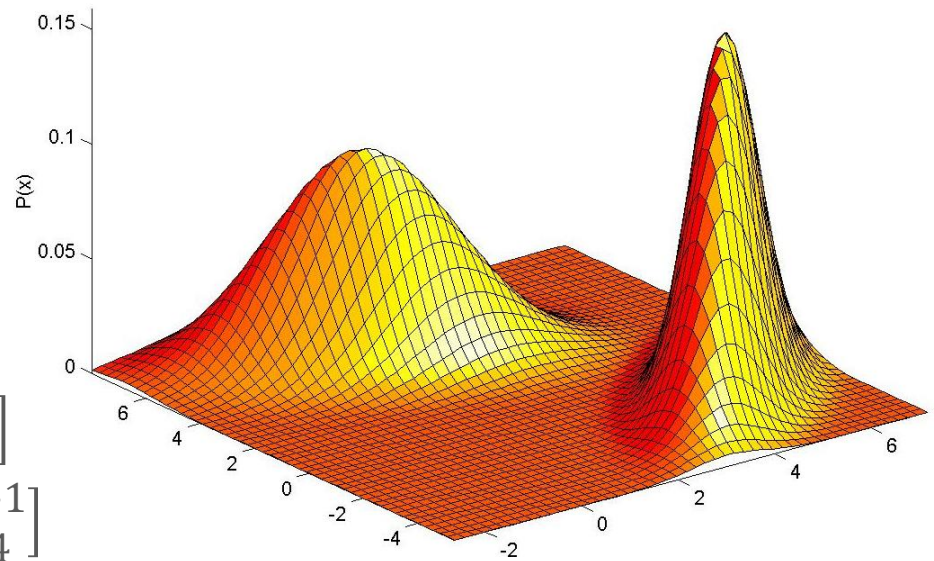
– The performance of kNN on 2D is illustrated in these figures

- The top figure shows the true density, a mixture of two bivariate Gaussians

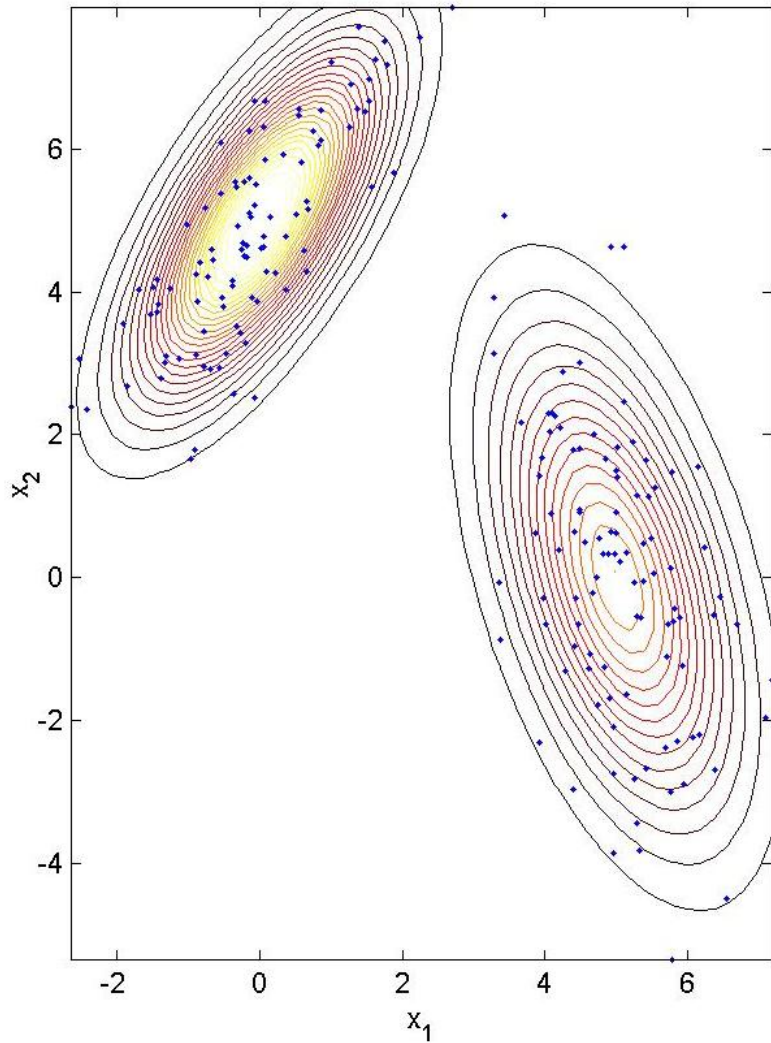
$$p(x) = \frac{1}{2}N(\mu_1, \Sigma_1) + \frac{1}{2}N(\mu_2, \Sigma_2)$$

$$\text{with } \begin{cases} \mu_1 = [0 \ 5]^T & \Sigma_1 = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix} \\ \mu_2 = [5 \ 0]^T & \Sigma_2 = \begin{bmatrix} 1 & -1 \\ -1 & 4 \end{bmatrix} \end{cases}$$

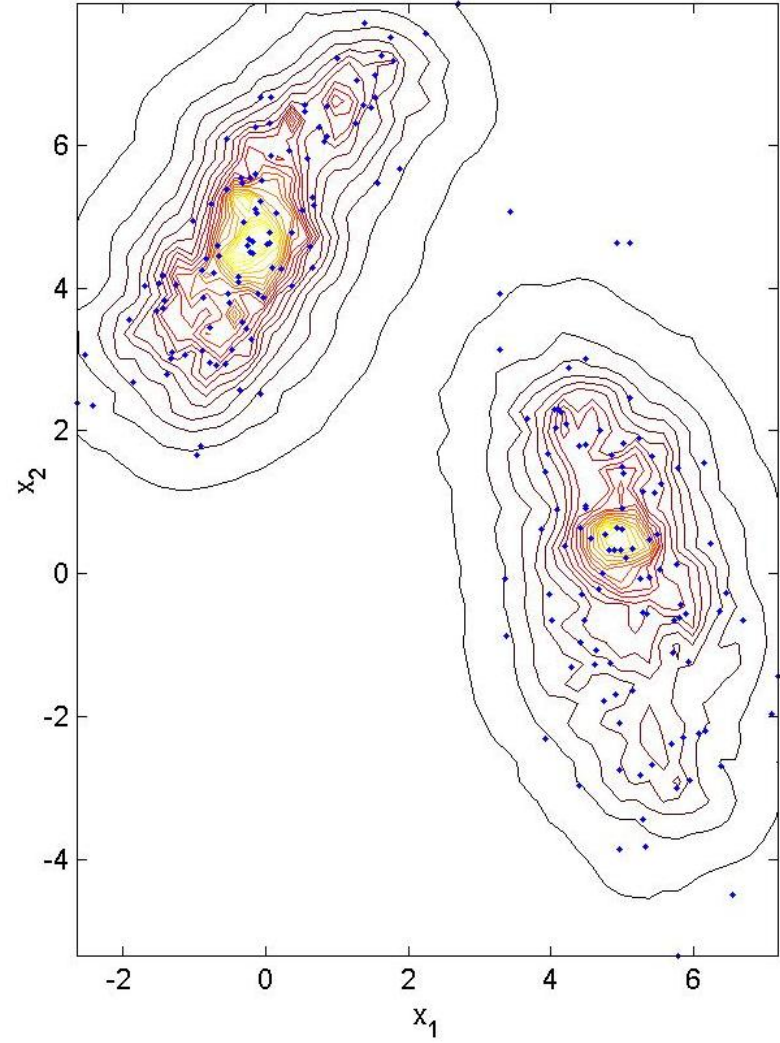
- The bottom figure shows the density estimate for  $k = 10$  neighbors and  $N = 200$  examples
- In the next slide we show the contours of the two distributions overlapped with the training data used to generate the estimate



**True density contours**



**kNN density estimate contours**



# kNN Density Estimation as a Bayesian classifier

**The main advantage of kNN is that it leads to a very simple approximation of the (optimal) Bayes classifier**

- Assume that we have a dataset with  $N$  examples,  $N_i$  from class  $\omega_i$ , and that we are interested in classifying an unknown sample  $x_u$ 
  - We draw a hyper-sphere of volume  $V$  around  $x_u$ . Assume this volume contains a total of  $k$  examples,  $k_i$  from class  $\omega_i$
- We can then approximate the likelihood functions as

$$p(x|\omega_i) = \frac{k_i}{N_i V}$$

- Similarly, the unconditional density can be estimated as

$$p(x) = \frac{k}{NV}$$

- And the priors are approximated by

$$P(\omega_i) = \frac{N_i}{N}$$

- Putting everything together, the Bayes classifier becomes

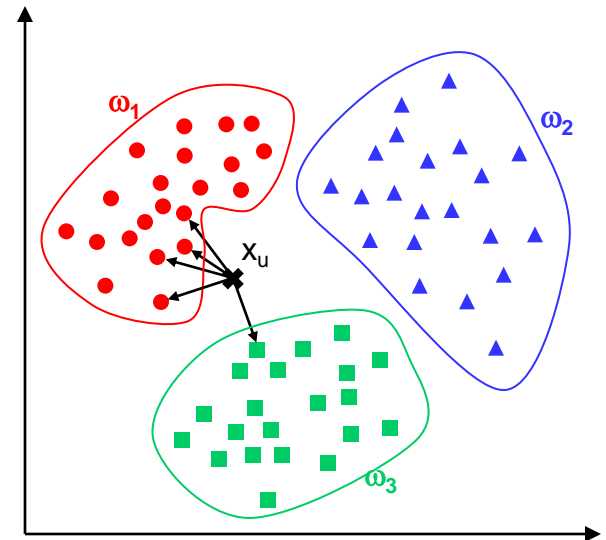
$$P(\omega_i|x) = \frac{p(x|\omega_i)P(\omega_i)}{p(x)} = \frac{\frac{k_i}{N_i V} \cdot \frac{N_i}{N}}{\frac{k}{NV}} = \frac{k_i}{k}$$



# The kNN classifier

## Definition

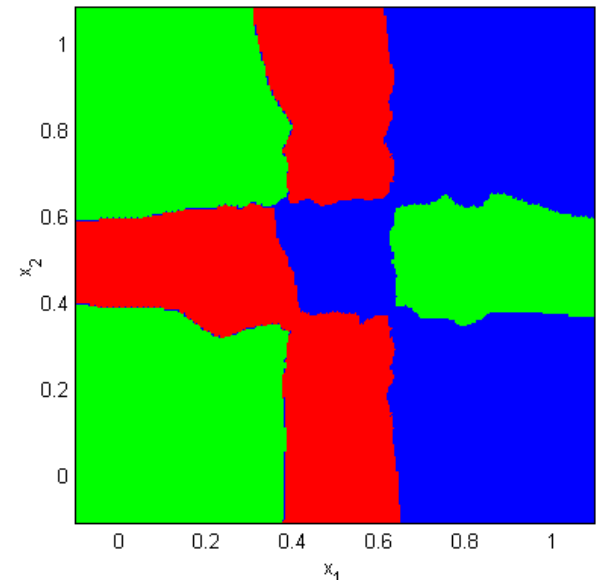
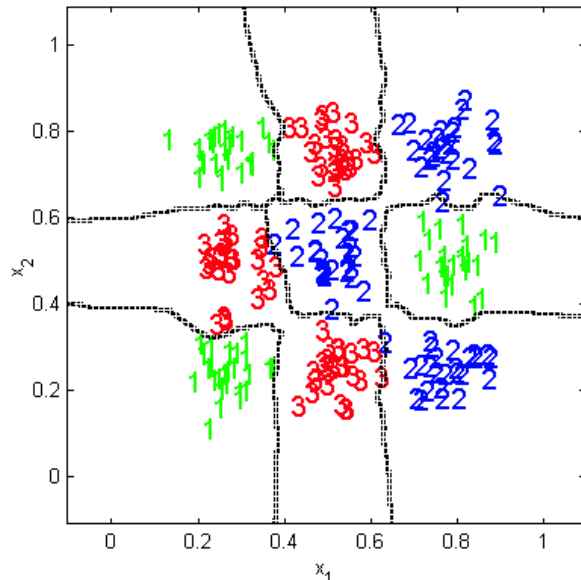
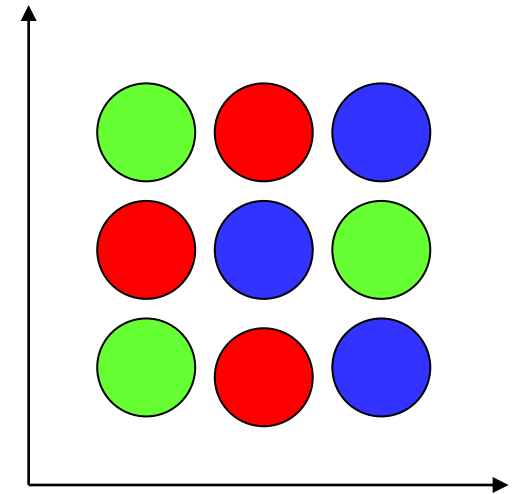
- The kNN rule is a very intuitive method that classifies unlabeled examples based on their similarity to examples in the training set
- For a given unlabeled example  $x_u \in \mathcal{R}^D$ , find the  $k$  “closest” labeled examples in the training data set and assign  $x_u$  to the class that appears most frequently within the k-subset
- The kNN only requires
  - An integer  $k$
  - A set of labeled examples (training data)
  - A metric to measure “closeness”
- Example
  - In the example here we have three classes and the goal is to find a class label for the unknown example  $x_u$
  - In this case we use the Euclidean distance and a value of  $k = 5$  neighbors
  - Of the 5 closest neighbors, 4 belong to  $\omega_1$  and 1 belongs to  $\omega_3$ , so  $x_u$  is assigned to  $\omega_1$ , the predominant class



# kNN in action

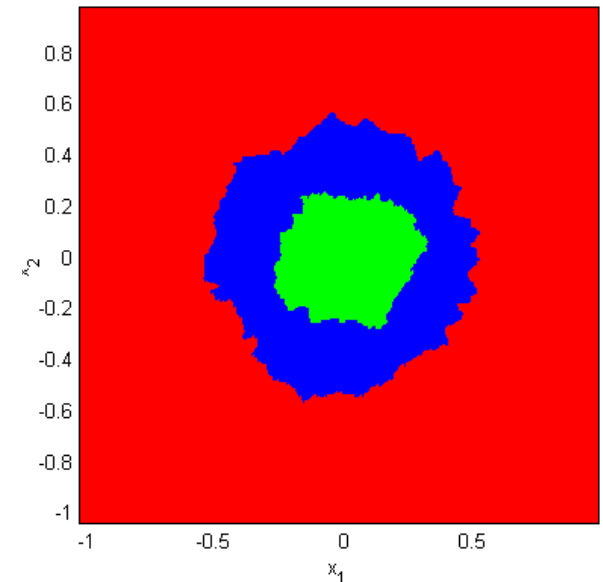
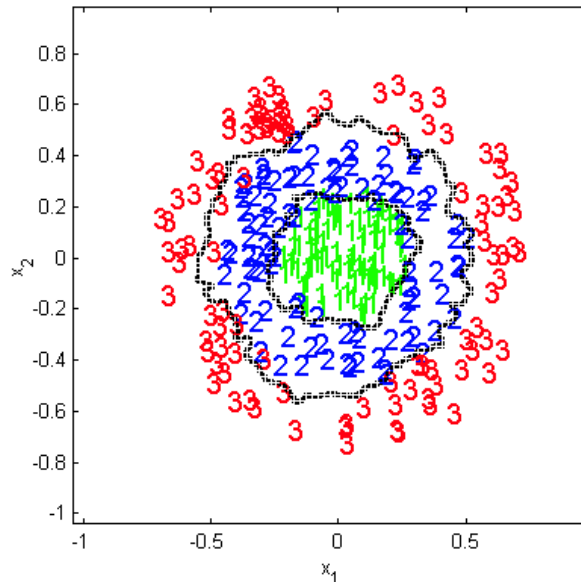
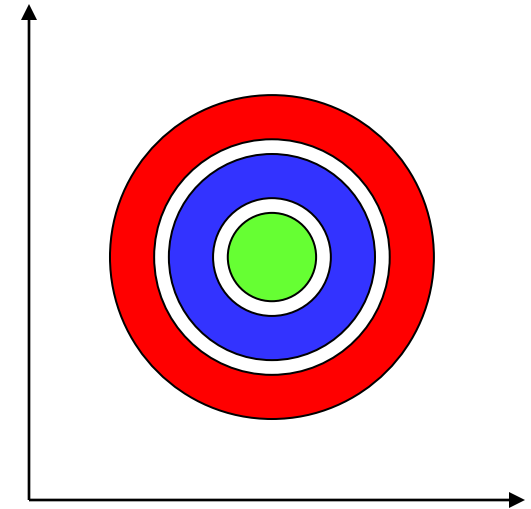
## Example I

- Three-class 2D problem with non-linearly separable, multimodal likelihoods
- We use the kNN rule ( $k = 5$ ) and the Euclidean distance
- The resulting decision boundaries and decision regions are shown below



## Example II

- Two-dim 3-class problem with unimodal likelihoods with a common mean; these classes are also not linearly separable
- We used the kNN rule ( $k = 5$ ), and the Euclidean distance as a metric



# kNN as a machine learning algorithm

## kNN is considered a lazy learning algorithm

- Defers data processing until it receives a request to classify unlabeled data
- Replies to a request for information by combining its stored training data
- Discards the constructed answer and any intermediate results

## This strategy is opposed to an eager learning algorithm which

- Compiles its data into a compressed description or model
  - A density estimate or density parameters (statistical PR)
  - A graph structure and associated weights (neural PR)
- Discards the training data after compilation of the model
- Classifies incoming patterns using the induced model, which is retained for future requests

## Tradeoffs

- Lazy algorithms have fewer computational costs than eager algorithms during training
- Lazy algorithms have greater storage requirements and higher computational costs on recall

[Aha, 1997]

# Characteristics of the kNN classifier

## Advantages

- Analytically tractable
- Simple implementation
- Nearly optimal in the large sample limit ( $N \rightarrow \infty$ )  
$$P_{Bayes}[error] < P_{1NN}[error] < 2P_{Bayes}[error]$$
- Uses local information, which can yield highly adaptive behavior
- Lends itself very easily to parallel implementations

## Disadvantages

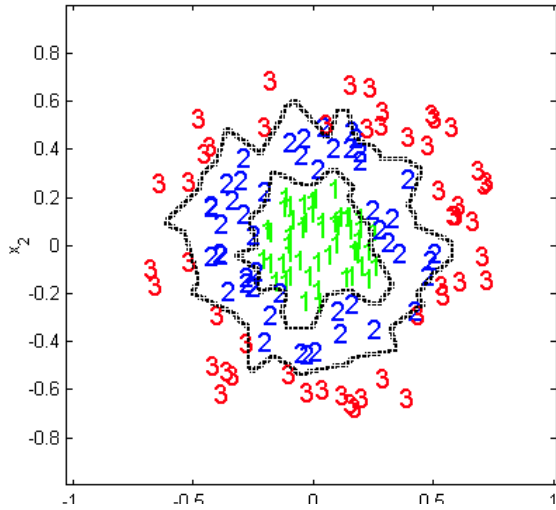
- Large storage requirements
- Computationally intensive recall
- Highly susceptible to the curse of dimensionality

## 1NN versus kNN

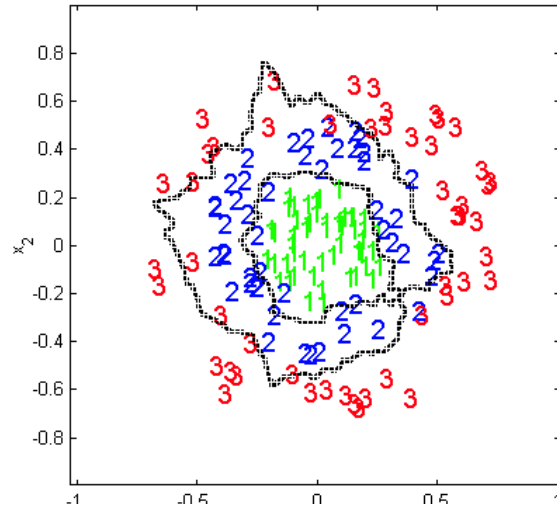
- The use of large values of  $k$  has two main advantages
  - Yields smoother decision regions
  - Provides probabilistic information, i.e., the ratio of examples for each class gives information about the ambiguity of the decision
- However, too large a value of  $k$  is detrimental
  - It destroys the locality of the estimation since farther examples are taken into account
  - In addition, it increases the computational burden

# kNN versus 1NN

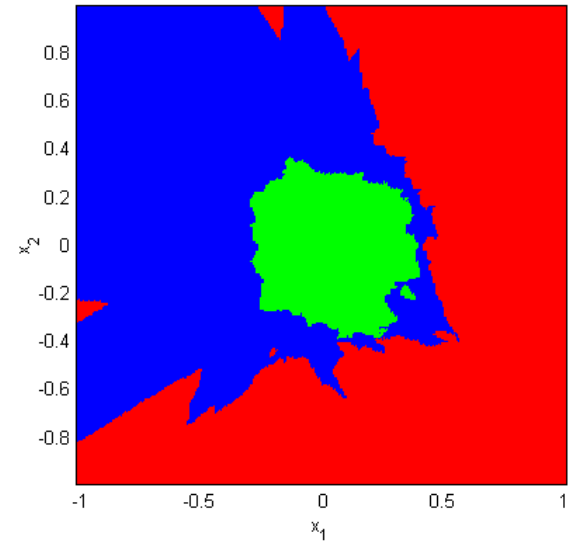
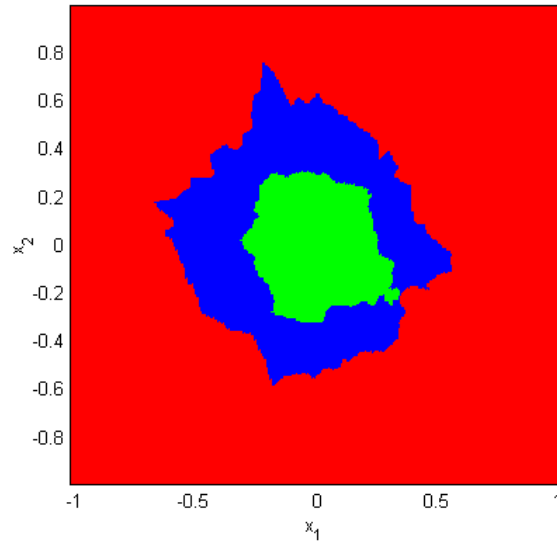
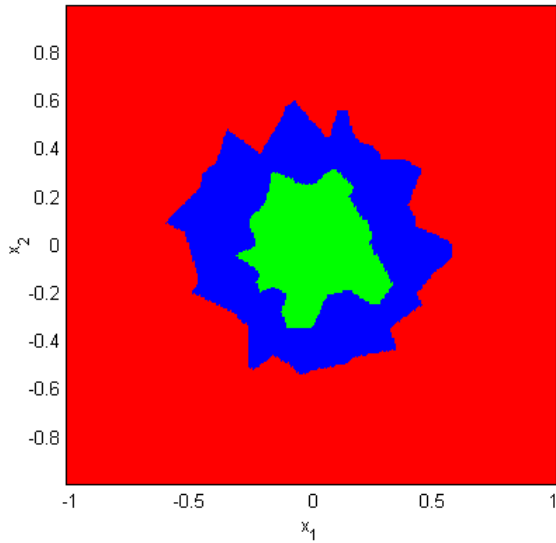
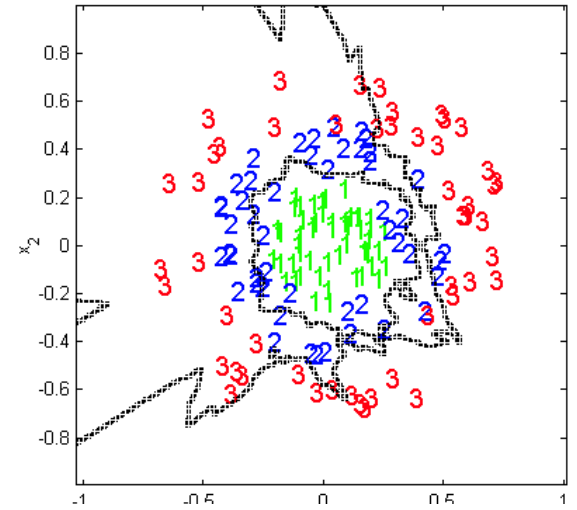
## 1-NN



## 5-NN



## 20-NN



# Optimizing storage requirements

**The basic kNN algorithm stores all the examples in the training set, creating high storage requirements (and computational cost)**

- However, the entire training set need not be stored since the examples may contain information that is highly redundant
  - A degenerate case is the earlier example with the multimodal classes, where each of the clusters could be replaced by its mean vector, and the decision boundaries would be practically identical
- In addition, almost all of the information that is relevant for classification purposes is located around the decision boundaries

**A number of methods, called edited kNN, have been derived to take advantage of this information redundancy**

- One alternative [Wilson 72] is to classify all the examples in the training set and remove those examples that are misclassified, in an attempt to separate classification regions by removing ambiguous points
- The opposite alternative [Ritter 75], is to remove training examples that are classified correctly, in an attempt to define the boundaries between classes by eliminating points in the interior of the regions

**A different alternative is to reduce the training examples to a set of prototypes that are representative of the underlying data**

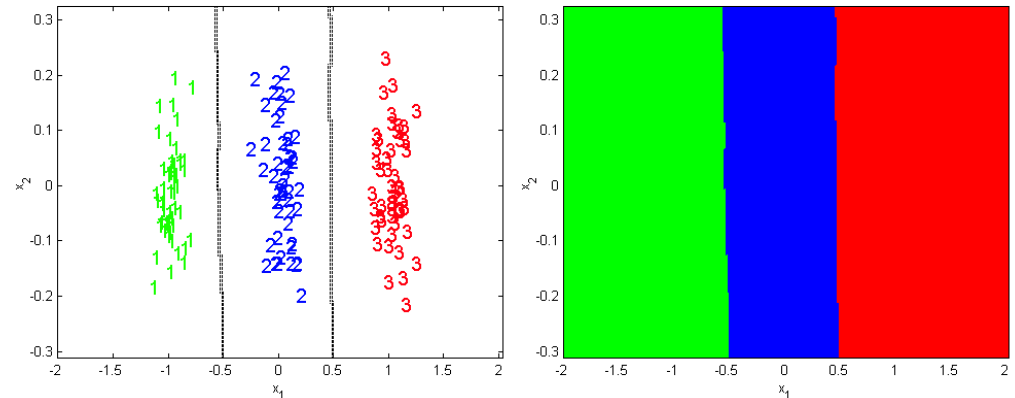
- The issue of selecting prototypes will be the subject of the lectures on clustering

# kNN and feature weighting

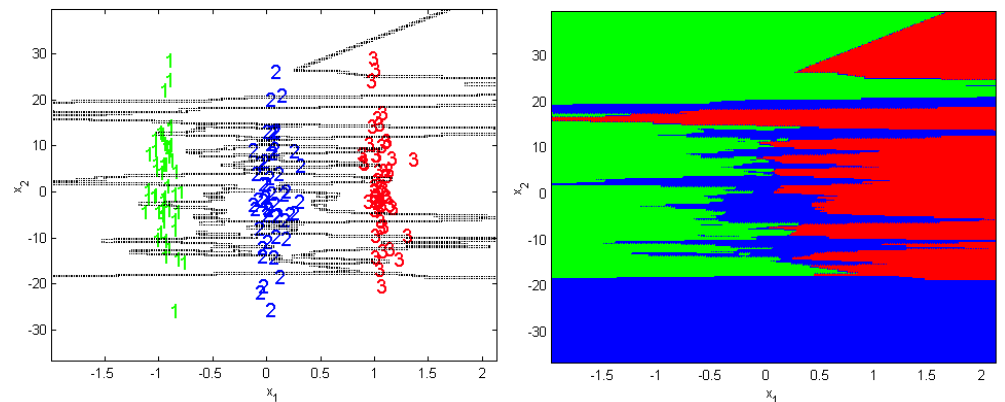
**kNN is sensitive to noise since it is based on the Euclidean distance**

- To illustrate this point, consider the example below
  - The first axis contains all the discriminatory information
  - The second axis is white noise, and does not contain classification information

- In a first case, both axes are scaled properly
  - kNN ( $k = 5$ ) finds decision boundaries fairly close to the optimal



- In a second case, the scale of the second axis has been increased 100 times
  - kNN is biased by the large values of the second axis and its performance is very poor





# Feature weighting

**The previous example illustrated the Achilles' heel of kNN: its sensitivity to noisy features**

- As a potential solution, one may attempt to normalize each feature to  $N(0,1)$
- Unfortunately, the Euclidean distance (see below) becomes very noisy for high dimensional problems if only a few of the features carry the classification information

$$d(x_u, x) = \sqrt{\sum_{k=1}^D (x_{u,k} - x_k)^2}$$

## Feature weighting

- The solution is to modify the Euclidean metric by a set of weights that capture the information content or “goodness” of each feature

$$d_w(x_u, x) = \sqrt{\sum_{k=1}^D [w_k (x_{u,k} - x_k)]^2}$$

- Note this is equivalent to performing a linear transformation with a diagonal matrix
  - Hence, feature weighting is a special case of feature extraction where the features are not allowed to interact
  - In turn, feature subset selection can be viewed as a special case of feature weighting where the weights can only take binary  $[0,1]$  values
- Do not confuse feature-weighting with distance-weighting, a kNN variant that weights the contribution of each neighbor according to its distance to the unlabeled example
  - Distance-weighting distorts the kNN estimate of  $P(\omega_i|x)$  and is NOT recommended
  - Studies have shown that distance-weighting does not improve kNN classification performance

# Feature weighting methods

## Feature weighting methods are divided in two groups

- Performance bias methods
  - These methods find a set of weights through an iterative procedure that uses the classifier's performance to select the next set of weights
  - These methods generally give good solutions since they can incorporate the classifier's feedback into the selection of weights
- Preset bias methods
  - These methods use a pre-determined function that measures the information content of each feature, e.g., mutual information and correlation between each feature and the class label
  - These methods have the advantage of executing very fast
- The issue of performance bias versus preset bias will be revisited when we cover feature subset selection (FSS)
  - In FSS the performance bias methods are called wrappers and preset bias methods are called filters

# Improving the NN search procedure

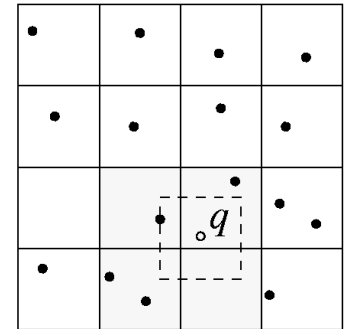
## The NN search procedure can be stated as follows

- Given a set of  $N$  points in  $D$ -dimensional space and an unlabeled example  $x_u \in \mathcal{X}^D$ , find the point that minimizes the distance to  $x_u$
- The naïve approach of computing a set of  $N$  distances, and finding the ( $k$ ) smallest becomes impractical for large values of  $N$  and  $D$

## Two classical algorithms can be used to speed up the NN search

### – Bucketing (a.k.a Elias's algorithm) [Welch 1971]

- The space is divided into identical cells; for each cell, the data points inside it are stored in a list
- Cells are examined in order of increasing distance from the query point; for each cell, the distance is computed between its internal data points and the query point
- The search terminates when the distance from the query point to the cell exceeds the distance to the closest point already visited

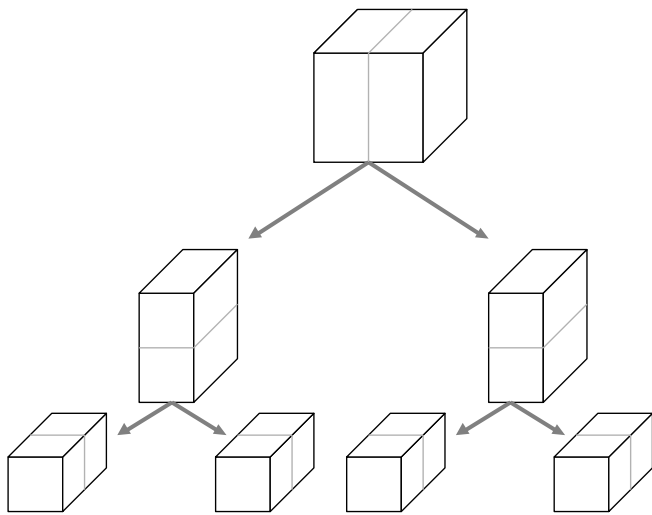


### – k-d trees [Bentley, 1975; Friedman et al, 1977]

- A k-d tree is a generalization of a binary search tree in high dimensions
  - Each internal node in a k-d tree is associated with a hyper-rectangle and a hyper-plane orthogonal to one of the coordinate axis
  - The hyper-plane splits the hyper-rectangle into two parts, which are associated with the child nodes
  - The partitioning process goes on until the # data points in the hyper-rectangle falls below some given threshold
- k-d trees partition the sample space according to the underlying distribution of the data: the partitioning being finer in regions where the density of data points is higher
  - For a given query point, the algorithm works by first descending the tree to find the data points lying in the cell that contains the query point
  - Then it examines surrounding cells if they overlap the ball centered at the query point and the closest data point so far

# k-d tree example

## Data structure (3D case)



## Partitioning (2D case)

