

# Nearest neighbor methods

## Lecture 11

David Sontag  
New York University

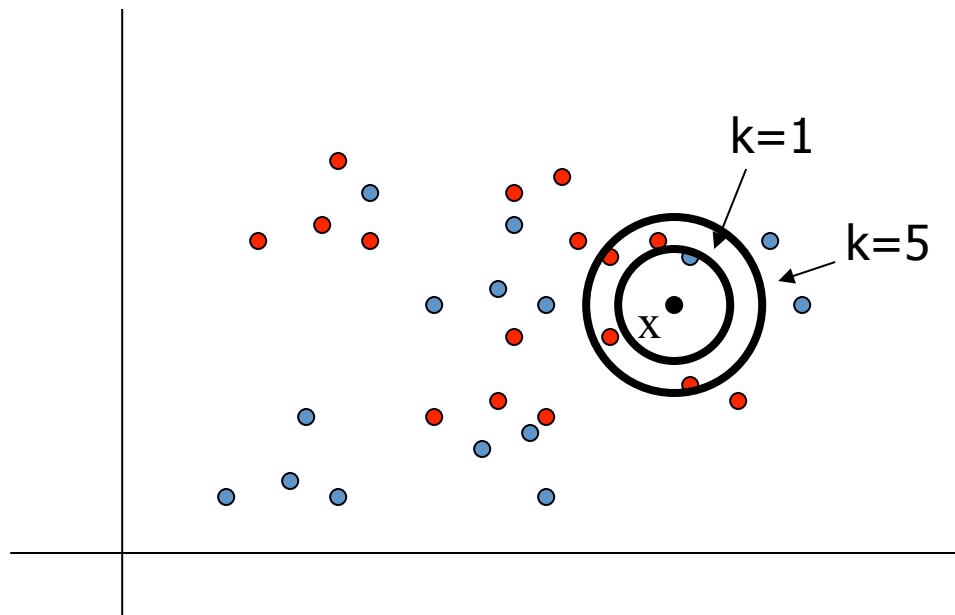
Slides adapted from Vibhav Gogate, Carlos Guestrin,  
Mehryar Mohri, & Luke Zettlemoyer

# Nearest Neighbor Algorithm

- Learning Algorithm:
  - Store training examples
- Prediction Algorithm:
  - To classify a new example  $\mathbf{x}$  by finding the training example  $(\mathbf{x}^i, y^i)$  that is *nearest* to  $\mathbf{x}$
  - Guess the class  $y = y^i$

# K-Nearest Neighbor Methods

- To classify a new input vector  $x$ , examine the  $k$ -closest training data points to  $x$  and assign the object to the most frequently occurring class

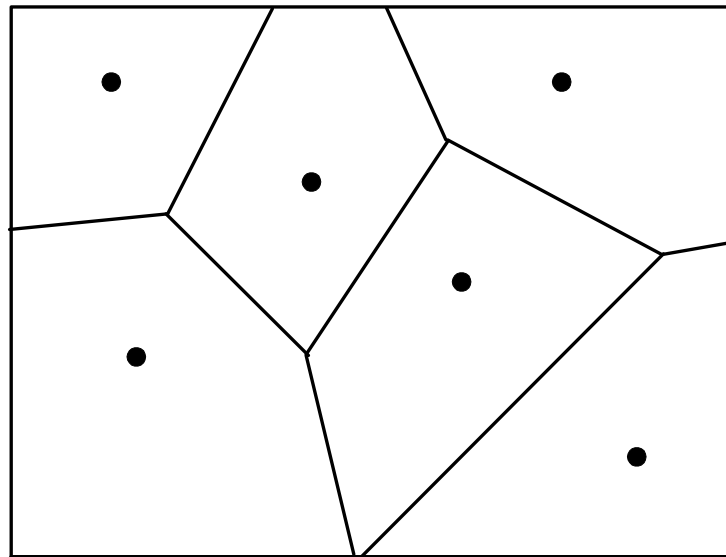


common values for  $k$ : 3, 5

# Decision Boundaries

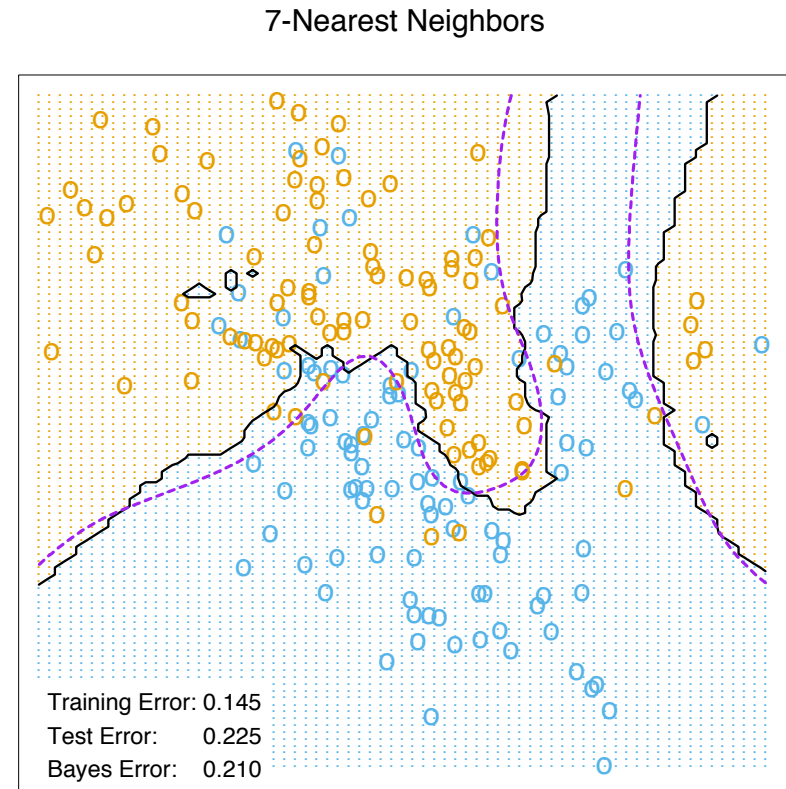
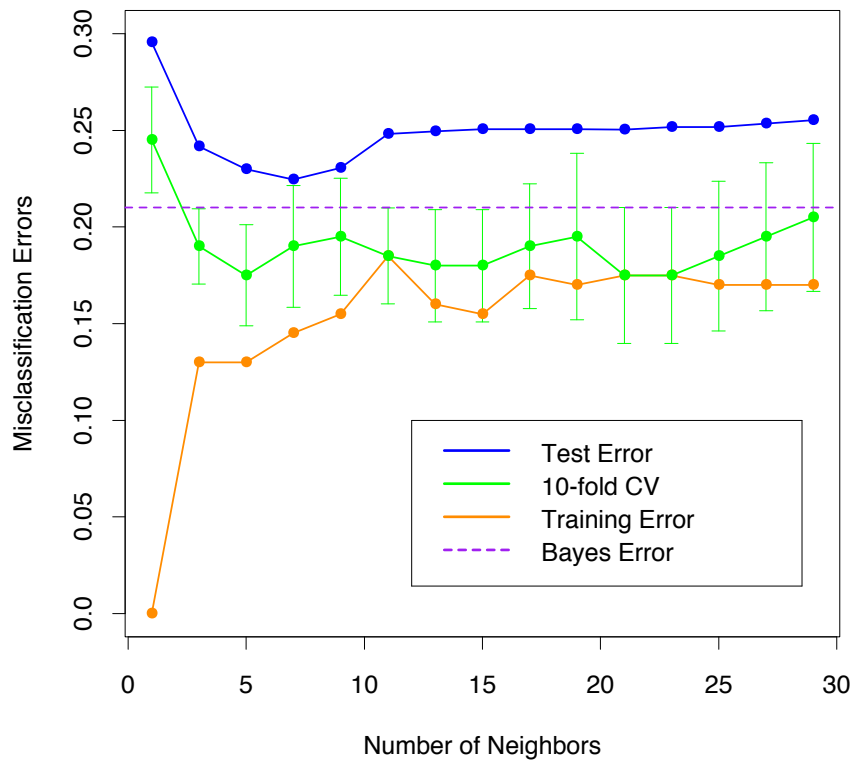
- The nearest neighbor algorithm does not explicitly compute decision boundaries. However, the decision boundaries form a subset of the Voronoi diagram for the training data.

*1-NN Decision Surface*



- The more examples that are stored, the more complex the decision boundaries can become

# Example results for k-NN



[Figures from Hastie and Tibshirani, Chapter 13]

# Nearest Neighbor

## When to Consider

- Instance map to points in  $R^n$
- Less than 20 attributes per instance
- Lots of training data

## Advantages

- Training is very fast
- Learn complex target functions
- Do not lose information

## Disadvantages

- Slow at query time
- Easily fooled by irrelevant attributes

# Issues

- Distance measure
  - Most common: Euclidean
- Choosing  $k$ 
  - Increasing  $k$  reduces variance, increases bias
- For high-dimensional space, problem that the nearest neighbor may not be very close at all!
- Memory-based technique. Must make a pass through the data for each classification. This can be prohibitive for large data sets.

## Distance

- Notation: object with  $p$  features

$$\mathbf{x}^i = (\mathbf{x}_1^i, \mathbf{x}_2^i, \dots, \mathbf{x}_p^i)$$

- Most common distance metric is *Euclidean* distance:

$$d_E(\mathbf{x}^i, \mathbf{x}^j) = \left( \sum_{k=1}^p (\mathbf{x}_k^i - \mathbf{x}_k^j)^2 \right)^{\frac{1}{2}}$$

- ED makes sense when different features are commensurate; each is variable measured in the same units.
- If the features are different, say length and weight, it is not clear.



## Normalization of features

Can divide features by them by the standard deviation, making them all equally important

The estimate for the standard deviation of feature **k**:

$$\hat{\sigma}_k = \left( \frac{1}{n} \sum_{i=1}^n \left( x_k^i - \bar{x}_k \right)^2 \right)^{\frac{1}{2}}$$

where  $\bar{x}_k$  is the sample mean:

$$\bar{x}_k = \frac{1}{n} \sum_{i=1}^n x_k^i$$

## Weighted Euclidean distance

Finally, if we have some idea of the relative importance of each variable, we can weight them:

$$d_{WE}(i, j) = \left( \sum_{k=1}^p w_k (x_k^i - x_k^j)^2 \right)^{\frac{1}{2}}$$