



### Key elements of this approach

- The set of labeled objects to be used for evaluating a test object's class
- A distance or similarity metric that can be used to compute the closeness of objects
- The value of k, the number of nearest neighbors
- The method used to determine the class of the target object based on the classes and distances of the k nearest neighbors

© Tan,Steinbach, Kumar	Introduction to Data Mining	4/18/2004	5
© Tan,Steinbach, Kumar	Introduction to Data Mining	4/18/2004	

## KNN is good at

- KNN is particularly well-suited for multimodal classes as well as applications in which an object can have many class labels.
- With multimodal classes, objects of a particular class labels are concentrated in several distinct areas of the data space, not just one. In statistical terms, the probability density function for the class does not have a single "bump" like a Gaussian, but rather, has a number of peaks.

### **Nearest-Neighbor Classifiers**



- Requires three things
  - The set of stored records
  - Distance Metric to compute distance between records
  - The value of k, the number of nearest neighbors to retrieve
  - To classify an unknown record:
    - Compute distance to other training records
    - Identify *k* nearest neighbors
    - Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

© Tan,Steinbach, Kumar

Introduction to Data Mining

4/18/2004

7

#### **Definition of Nearest Neighbor**



(a) 1-nearest neighbor

(b) 2-nearest neighbor

(c) 3-nearest neighbor

K-nearest neighbors of a record x are data points that have the k smallest distance to x

### 1 nearest-neighbor



## **Nearest Neighbor Classification**

• Compute distance between two points:

- Euclidean distance

$$d(p,q) = \sqrt{\sum_{i} (p_i - q_i)^2}$$

• Determine the class from nearest neighbor list

- take the majority vote of class labels among the k-nearest neighbors
- Weigh the vote according to distance

weight factor, w = 1/d<sup>2</sup>

© Tan, Steinbach, Kumar

## Nearest Neighbor Classification...

- Choosing the value of k:
  - If k is too small, sensitive to noise points
  - If k is too large, neighborhood may include points from other classes



© Tan,Steinbach, Kumar

Introduction to Data Mining

4/18/2004

11

## Nearest Neighbor Classification...

- Scaling issues
  - Attributes may have to be scaled to prevent distance measures from being dominated by one of the attributes
  - Example:
    - height of a person may vary from 1.5m to 1.8m
    - weight of a person may vary from 90lb to 300lb
    - income of a person may vary from \$10K to \$1M



# Example: PEBLS

<ul> <li>PEBLS: Para System (Cos</li> <li>Works with features</li> </ul>	allel Examplar-Based at & Salzberg) th both continuous ar	Learning nd nominal	
<ul> <li>For nominal va difference</li> </ul>	inal features, distance be alues is computed using metric (MVDM)	etween two modified value	
<ul> <li>Each reco</li> </ul>	ord is assigned a wei	ght factor	
<ul> <li>Number c</li> </ul>	of nearest neighbor, k	x = 1	
© Tan,Steinbach, Kumar	Introduction to Data Mining	4/18/2004	15