# Statistical Learning and Data Mining
# CS 363D/ SSC 358

# Lecture: Nearest Neighbor Classifiers

Prof. Pradeep Ravikumar
pradeepr@cs.utexas.edu

# Instance-Based Classifiers

## Set of Stored Cases

| Atr1 | ……..… | AtrN | Class |
|------|--------|------|-------|
|      |        |      | A     |
|      |        |      | B     |
|      |        |      | B     |
|      |        |      | C     |
|      |        |      | A     |
|      |        |      | C     |
|      |        |      | B     |

- **Store the training records**

- **Use training records to predict the class label of unseen cases**

## Unseen Case

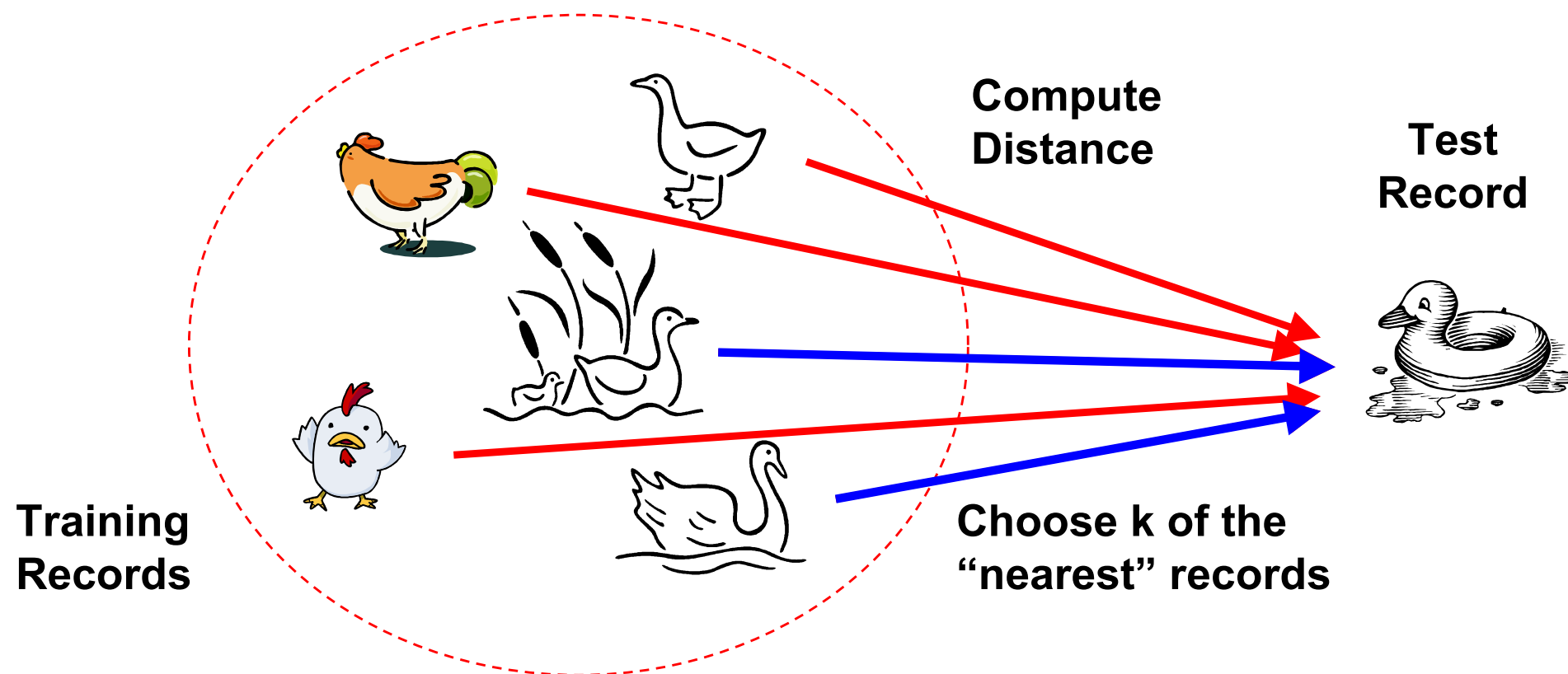| Atr1 | ……… | AtrN |
|------|------|------|
|      |      |      |

# Instance-Based Classifiers

- Examples:
  - Rote-learner
    - ◆ Memorizes entire training data and performs classification only if attributes of record match one of the training examples exactly

# Instance-Based Classifiers

- Examples:
  - Rote-learner
    - Memorizes entire training data and performs classification only if attributes of record match one of the training examples exactly

  - Nearest neighbor
    - Uses k "closest" points (nearest neighbors) for performing classification

# Nearest-Neighbor Classifiers

● Basic idea:

  – If it walks like a duck, quacks like a duck, then it's probably a duck



Compute Distance

Test Record

Training Records

Choose k of the "nearest" records
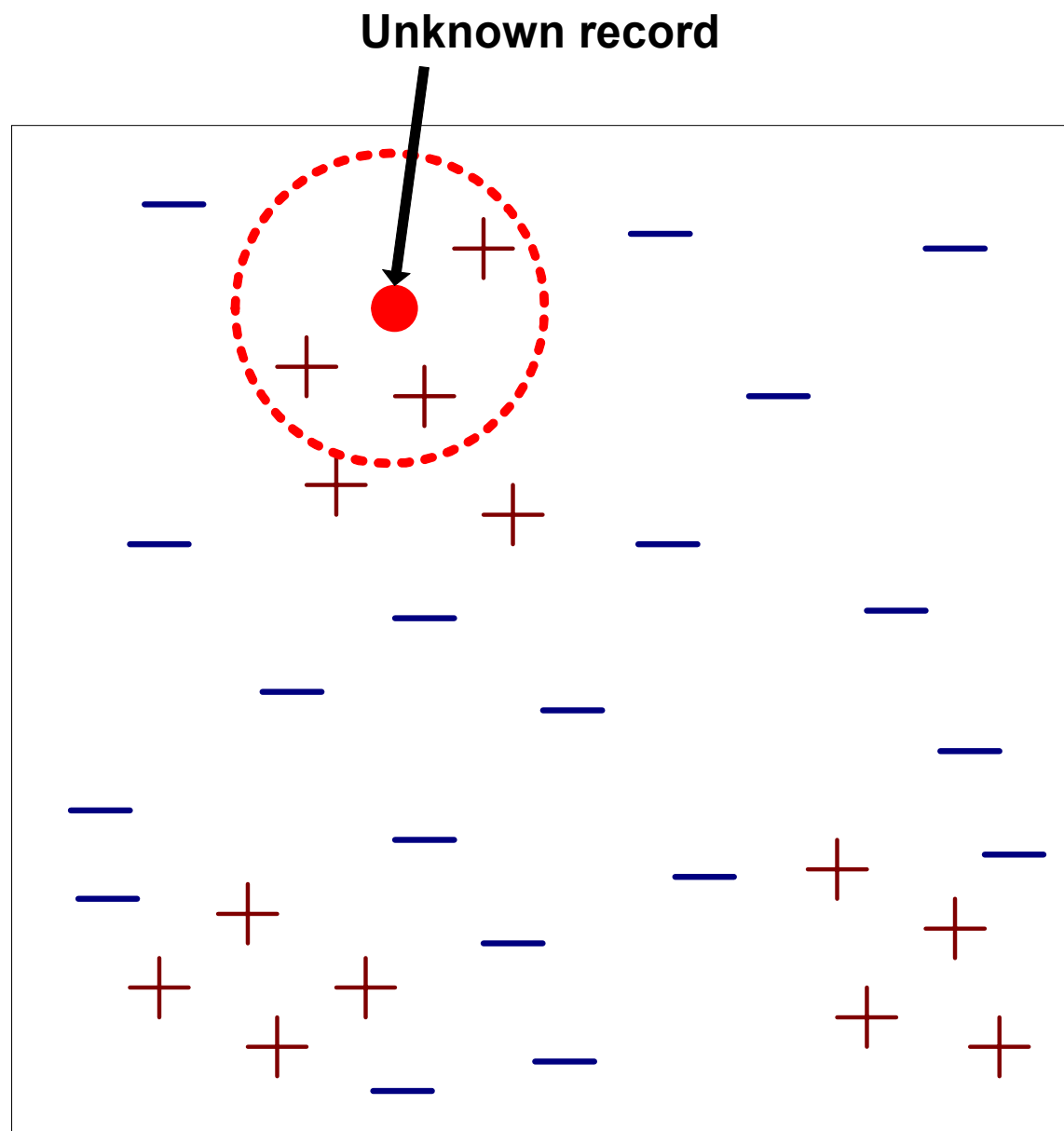
# Nearest-Neighbor Classifiers

- Basic idea:
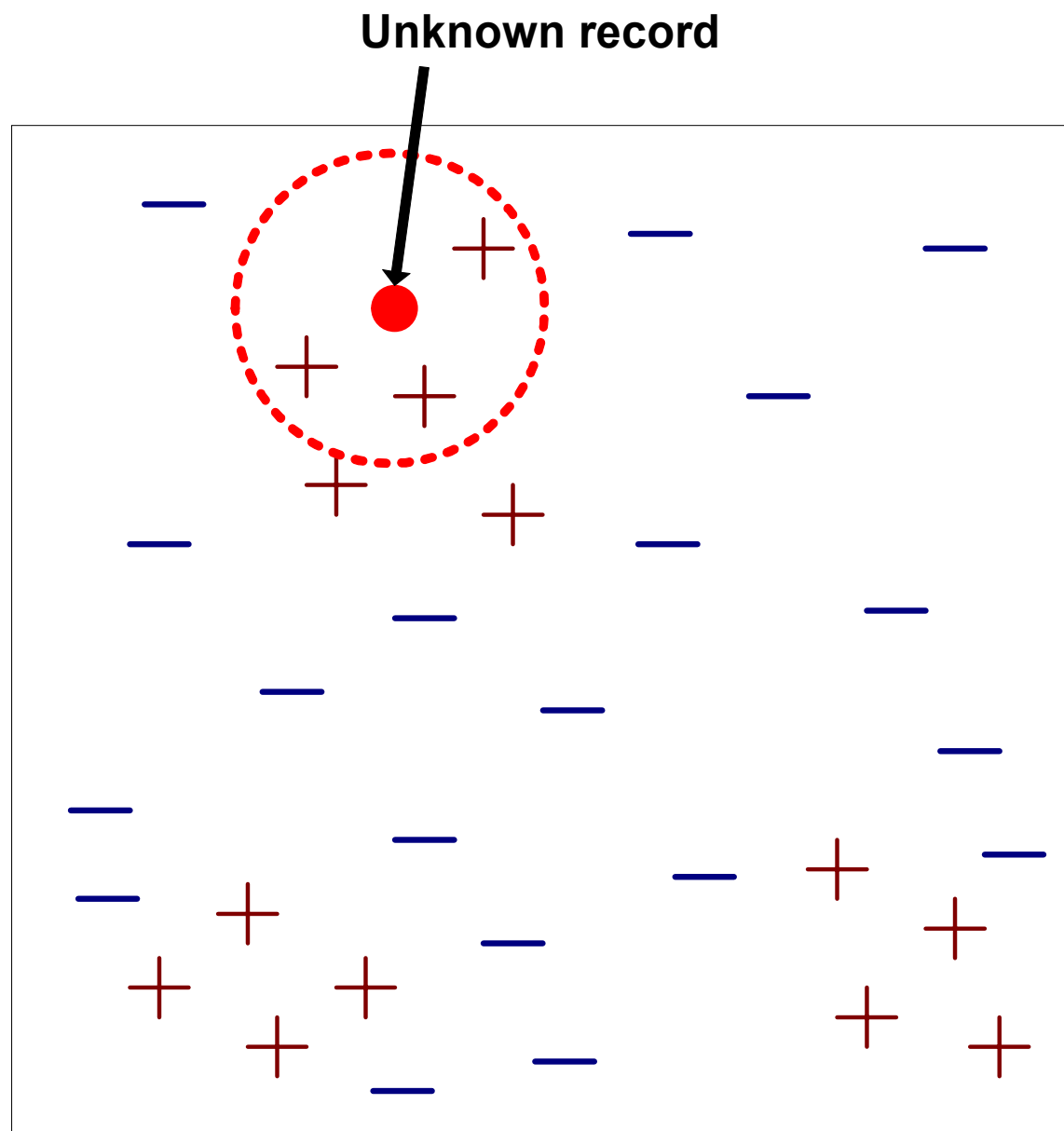
# Nearest-Neighbor Classifiers



**Unknown record**

- 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

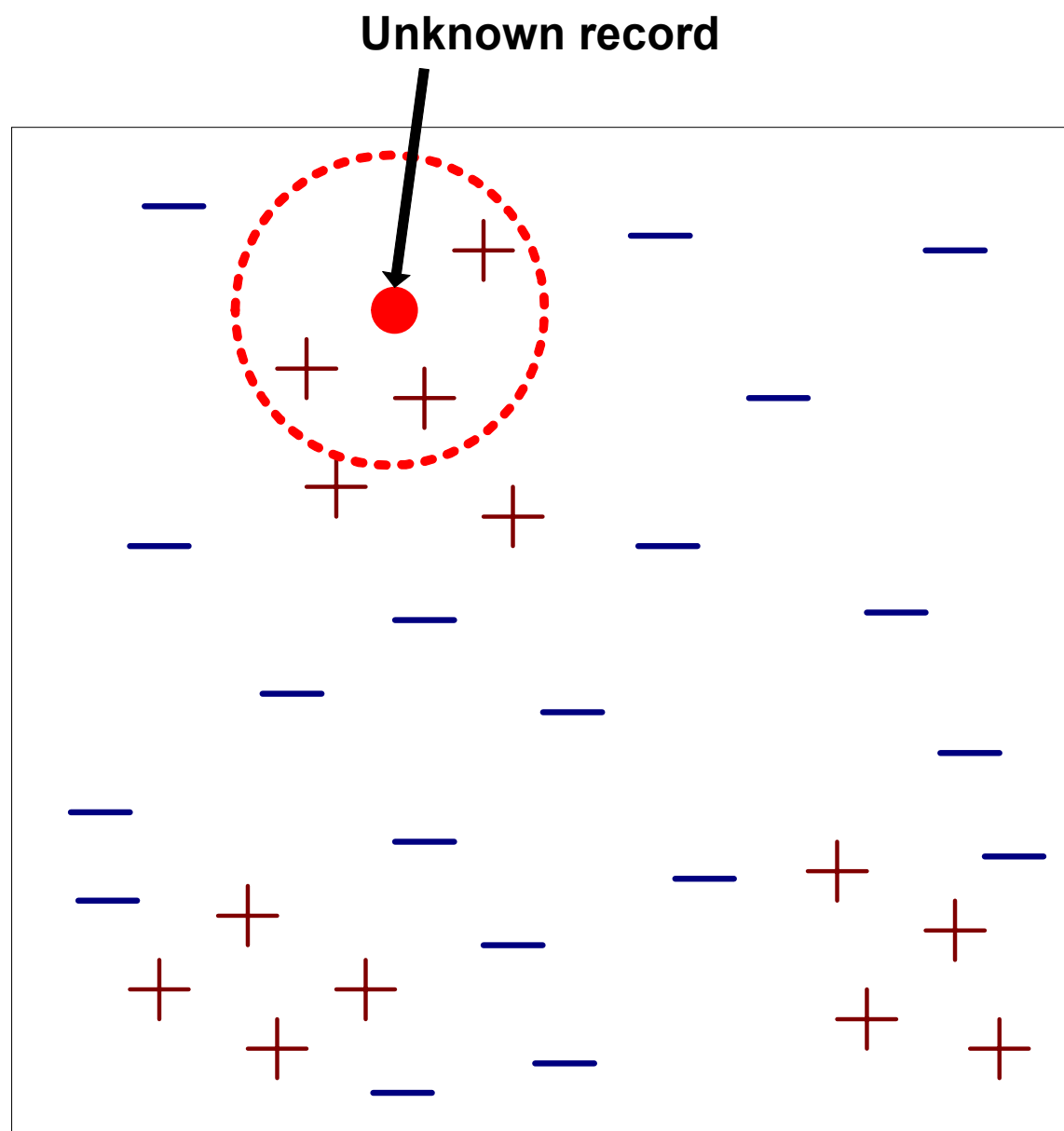# Nearest-Neighbor Classifiers

**Unknown record**



- 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

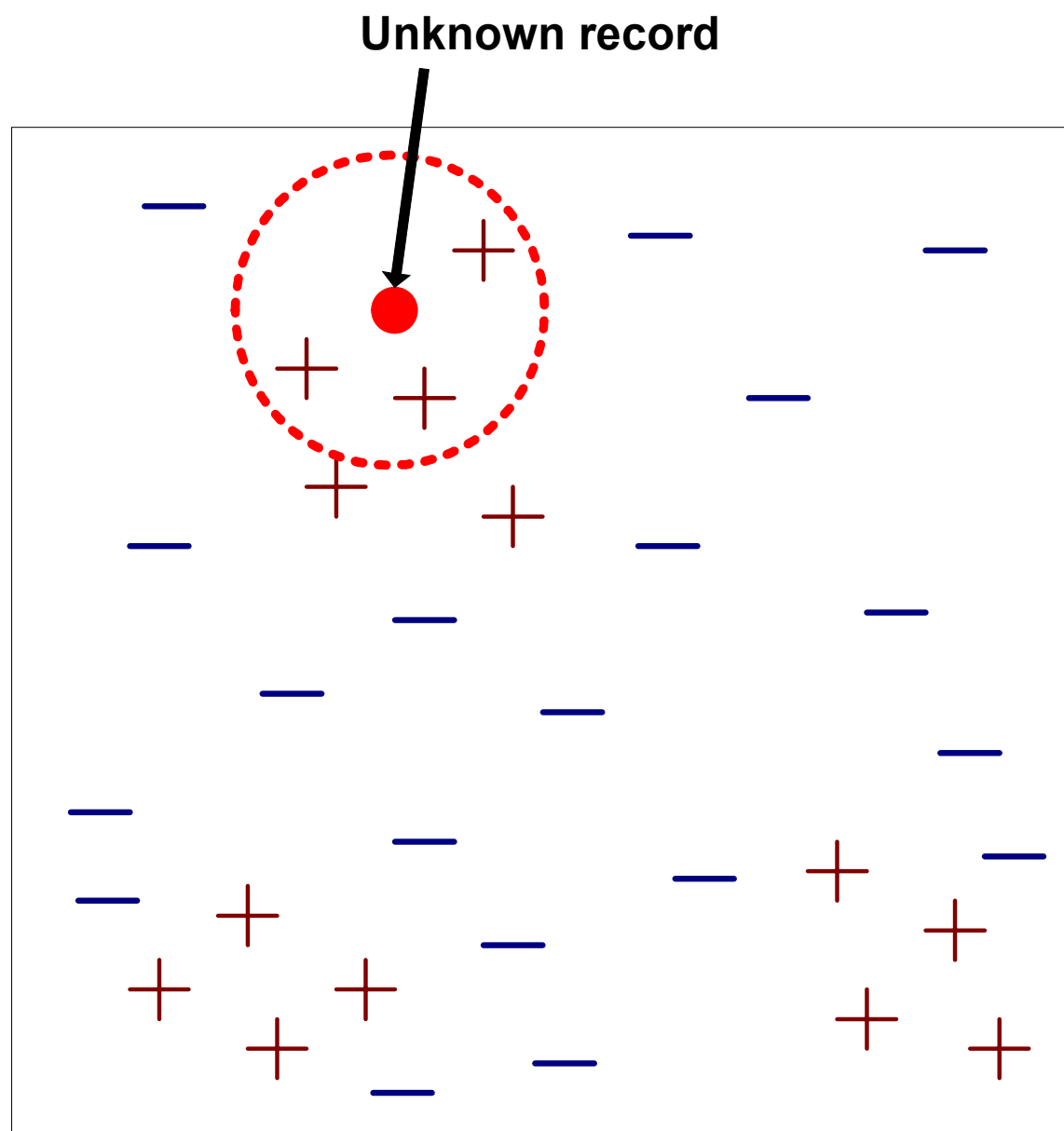- To classify an unknown record:

# Nearest-Neighbor Classifiers

**Unknown record**



- 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
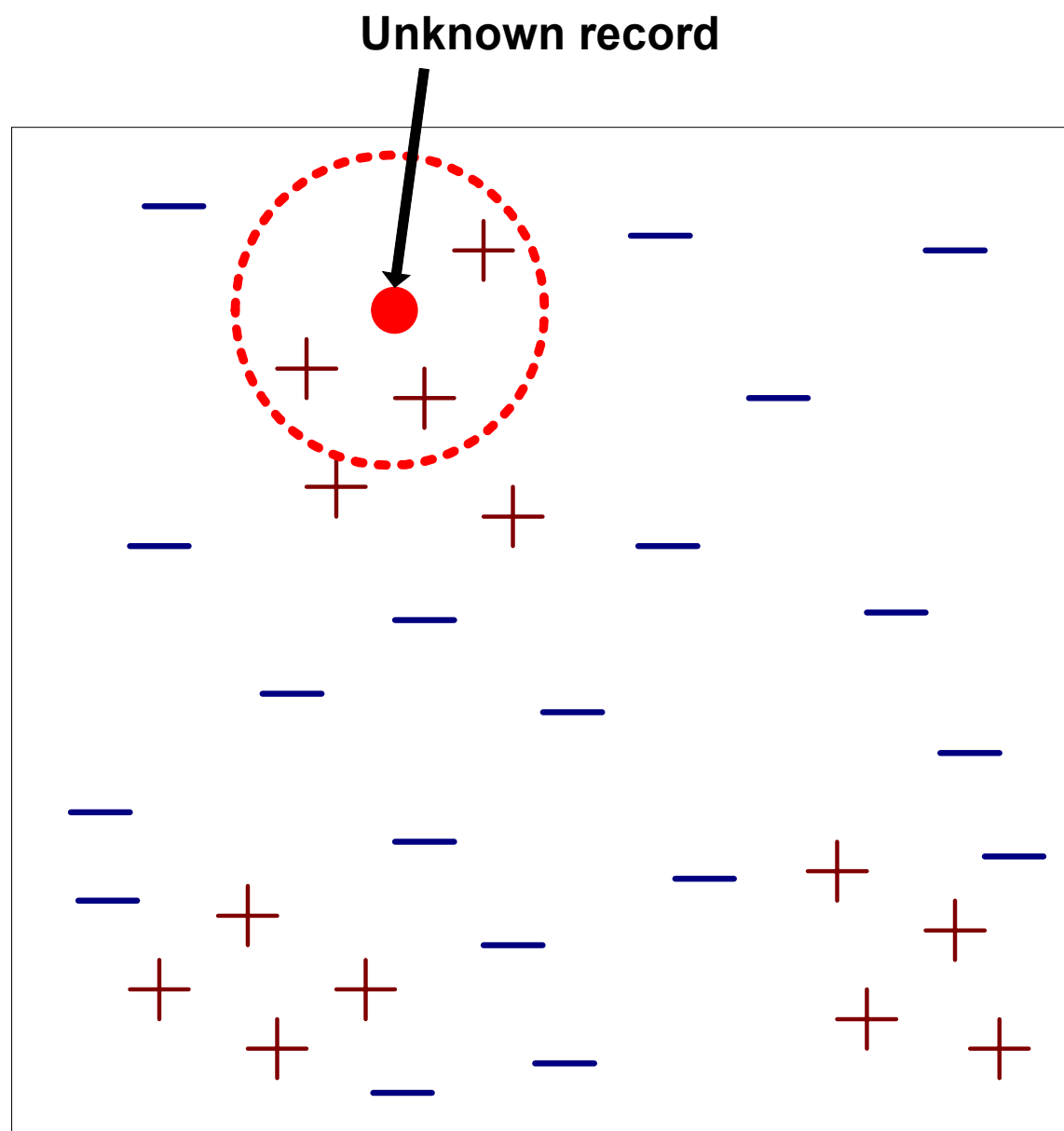
# Nearest-Neighbor Classifiers

**Unknown record**



- 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
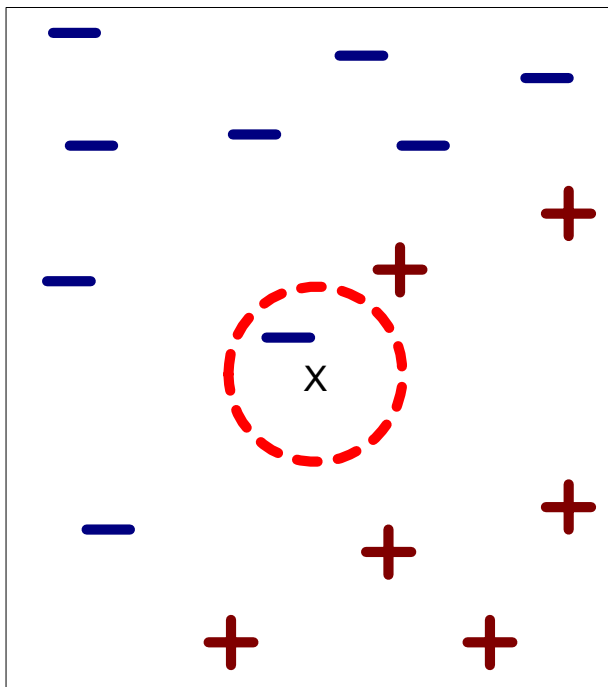
# Nearest-Neighbor Classifiers

**Unknown record**

- 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)
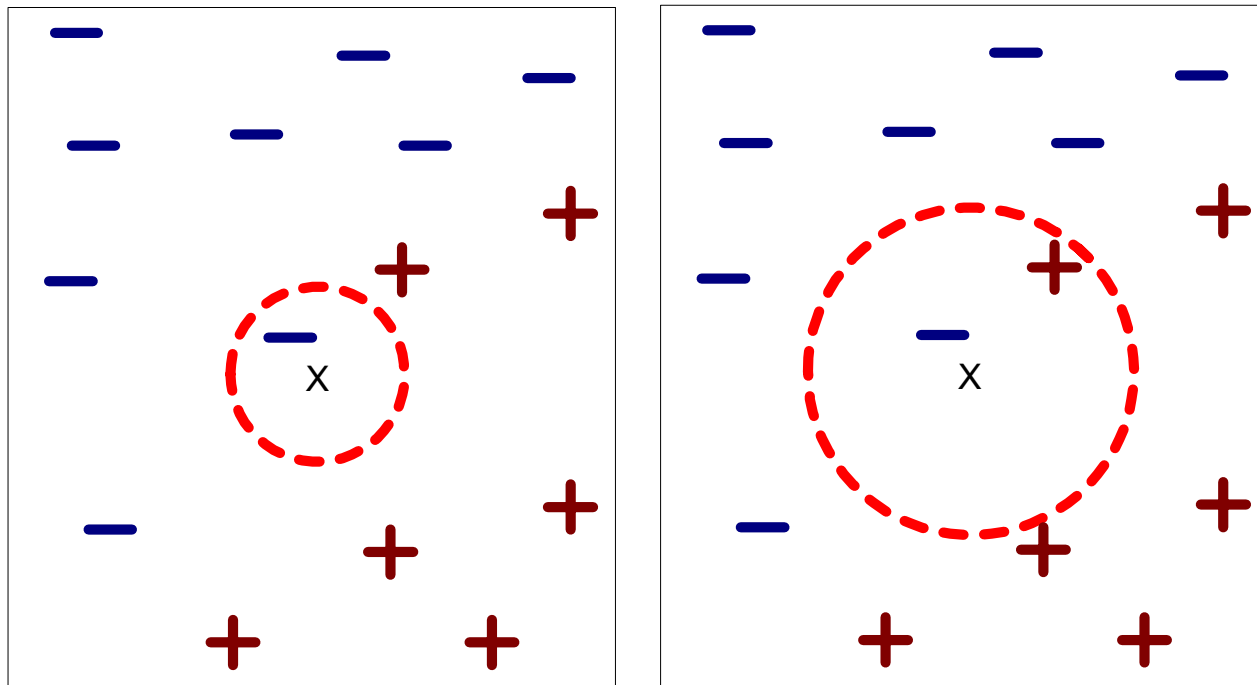
# Definition of Nearest Neighbor



(a) 1-nearest neighbor

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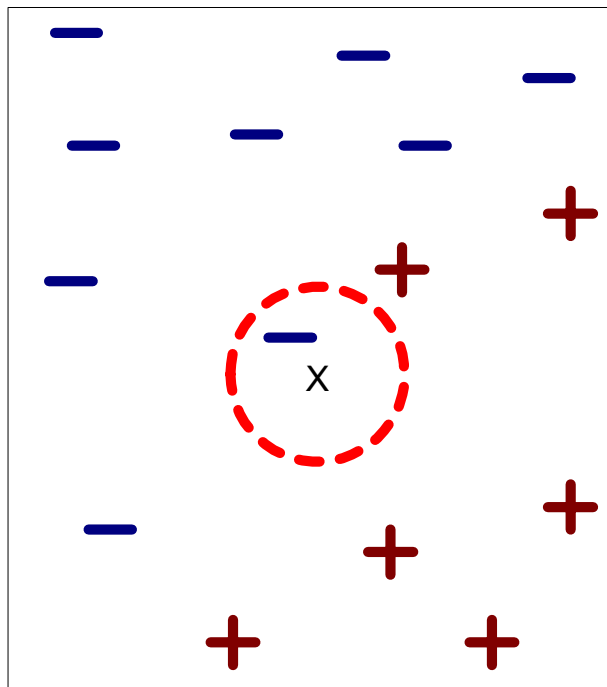that have the k smallest distance to x

# Definition of Nearest Neighbor



(a) 1-nearest neighbor     (b) 2-nearest neighbor

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

# Definition of Nearest Neighbor
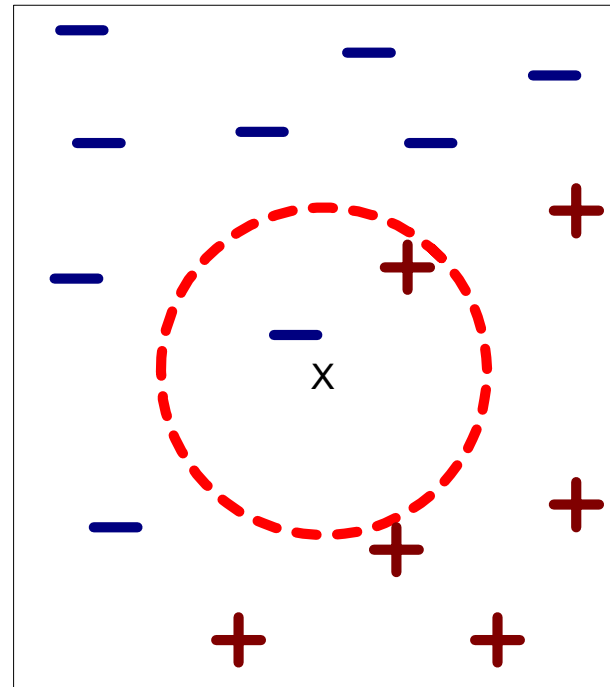


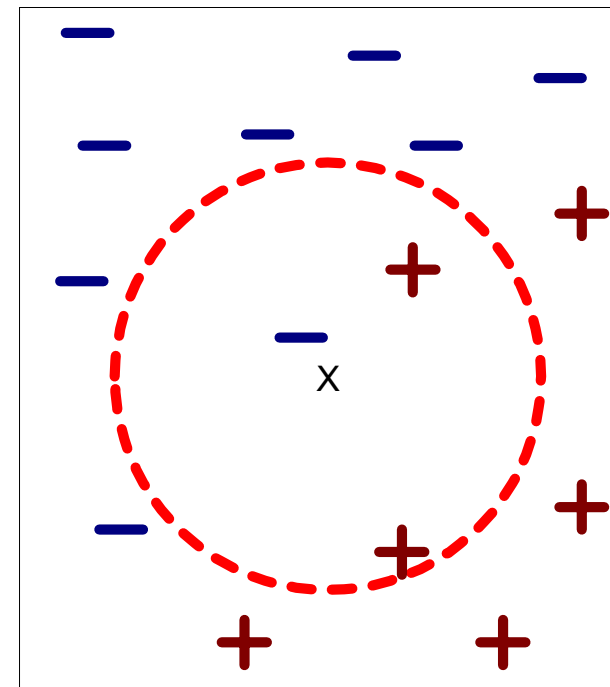(a) 1-nearest neighbor      (b) 2-nearest neighbor      (c) 3-nearest neighbor
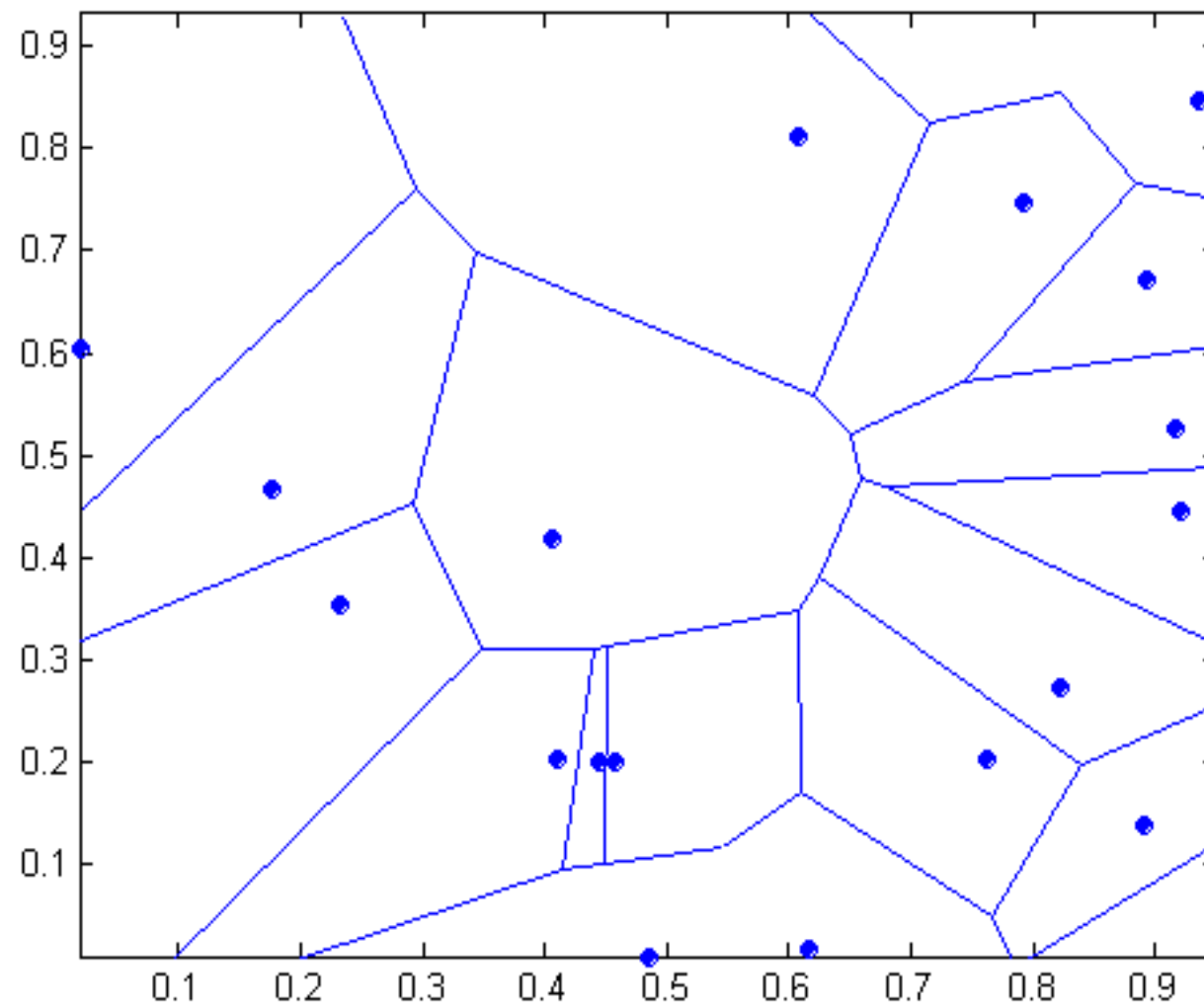
K-nearest neighbors of a record x are data points
that have the k smallest distance to x

# 1 Nearest-Neighbor

Voronoi Diagram

# 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

# 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

# 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, $w = 1/d^2$

# Nearest Neighbor Classification

- Choosing the value of k:
  - If k is too small,

# Nearest Neighbor Classification

- Choosing the value of k:
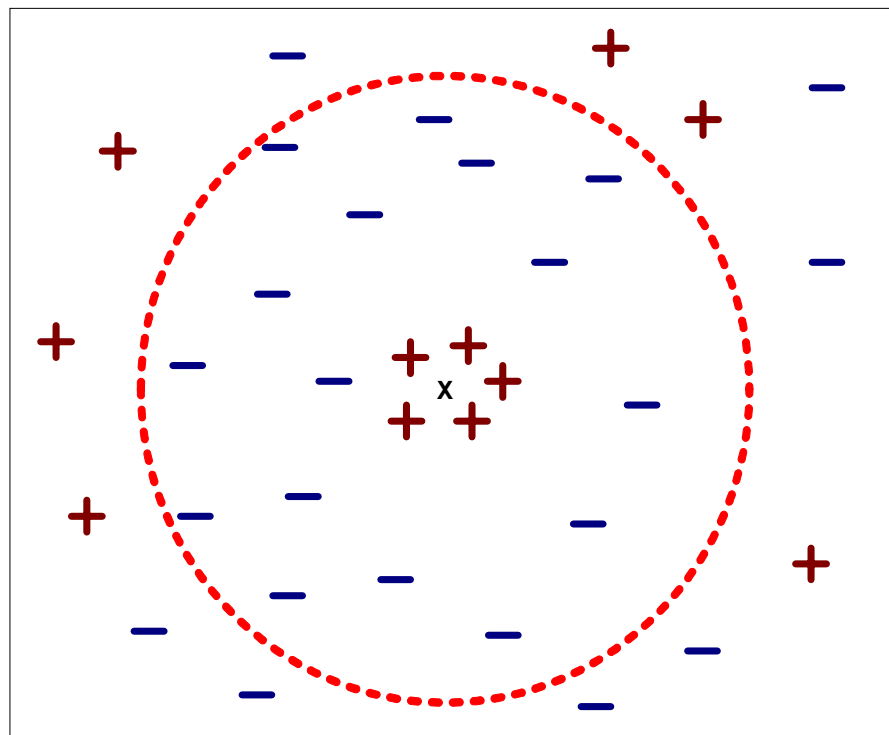  - If k is too small, sensitive to noise points

# Nearest Neighbor Classification

- Choosing the value of k:
  - If k is too small, sensitive to noise points
  - If k is too large,

# 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

# 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

# Nearest Neighbor Classification

- Problem with Euclidean measure:
  - High dimensional data
    - curse of dimensionality

# Nearest Neighbor Classification

- Problem with Euclidean measure:
  - High dimensional data
    - curse of dimensionality
  - Can produce counter-intuitive results

| 1 1 1 1 1 1 1 1 1 1 1 0 |
| --- |

| 0 1 1 1 1 1 1 1 1 1 1 1 |
| --- |

**d = 1.4142**

vs

| 1 0 0 0 0 0 0 0 0 0 0 0 |
| --- |

| 0 0 0 0 0 0 0 0 0 0 0 1 |
| --- |

**d = 1.4142**

# Nearest Neighbor Classification

- Problem with Euclidean measure:
  - High dimensional data
    - curse of dimensionality
  - Can produce counter-intuitive results

| 1 1 1 1 1 1 1 1 1 1 1 0 |
|---|

| 0 1 1 1 1 1 1 1 1 1 1 1 |
|---|

vs

| 1 0 0 0 0 0 0 0 0 0 0 0 |
|---|

| 0 0 0 0 0 0 0 0 0 0 0 1 |
|---|

**d = 1.4142**                    **d = 1.4142**

- Solution: Normalize the vectors to unit length

# Nearest Neighbor Classification

- k-NN classifiers are lazy learners
  - It does not build models explicitly
  - Unlike eager learners such as decision tree induction and rule-based systems
  - Classifying unknown records are relatively expensive
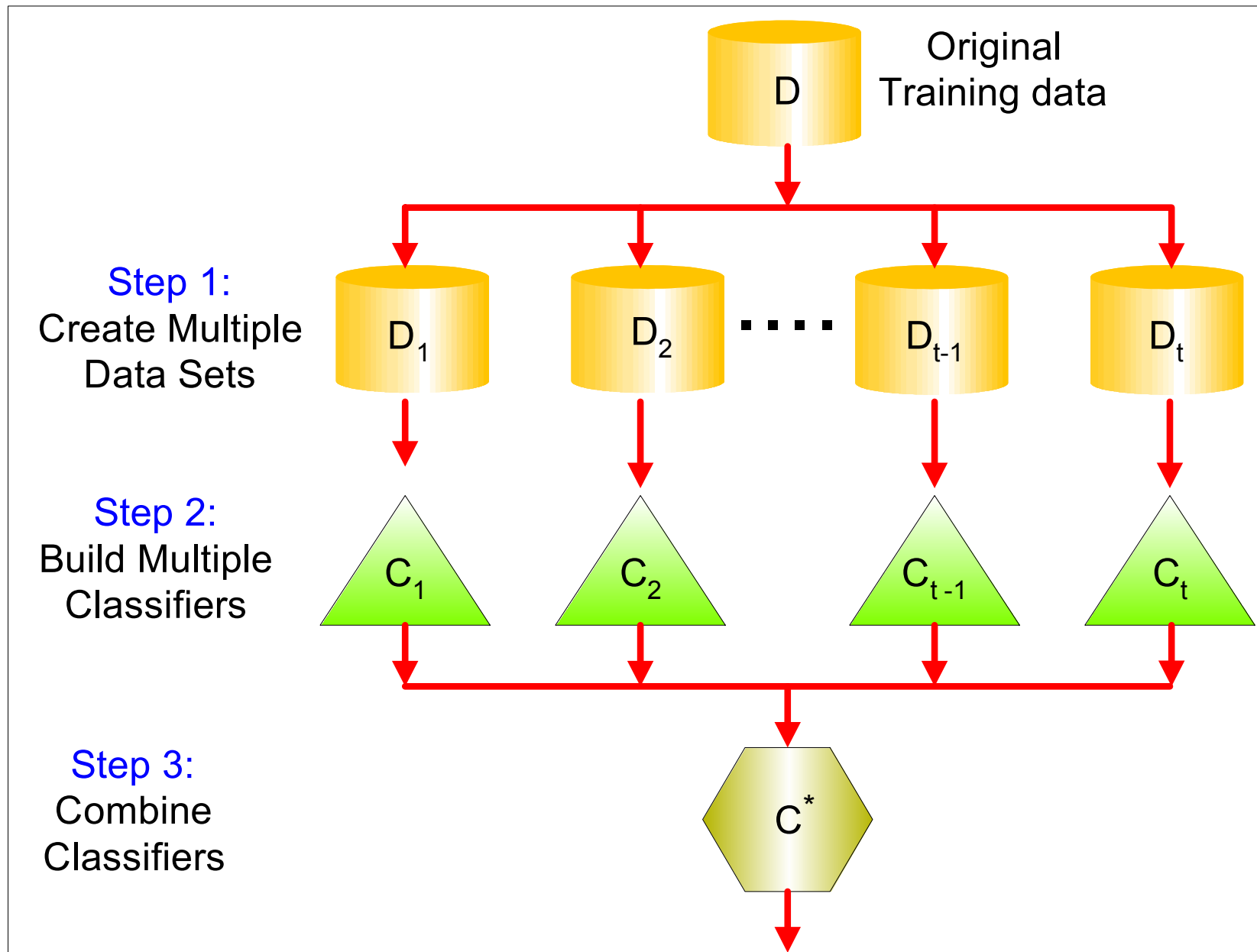
# Ensemble Methods

# Ensemble Methods

- Construct a set of classifiers from the training data

- Predict class label of previously unseen records by aggregating predictions made by multiple classifiers

# General Idea

# Why does it work?

- Suppose there are 25 base classifiers
  - Each classifier has error rate, $\varepsilon$ = 0.35
  - Assume classifiers are independent
  - Probability that the ensemble classifier makes a wrong prediction:

$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1-\varepsilon)^{25-i} = 0.06$$

# Many Approaches to Step 1 (Creating Multiple Datasets)

- Copy the dataset multiple times

- Partitioning the dataset

- Bagging

- Boosting

# Bagging

- Sampling with replacement

| Original Data | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Bagging (Round 1) | 7 | 8 | 10 | 8 | 2 | 5 | 10 | 10 | 5 | 9 |
| Bagging (Round 2) | 1 | 4 | 9 | 1 | 2 | 3 | 2 | 7 | 3 | 2 |
| Bagging (Round 3) | 1 | 8 | 5 | 10 | 5 | 5 | 9 | 6 | 3 | 7 |

- Build classifier on each bootstrap sample