# Supervised Learning
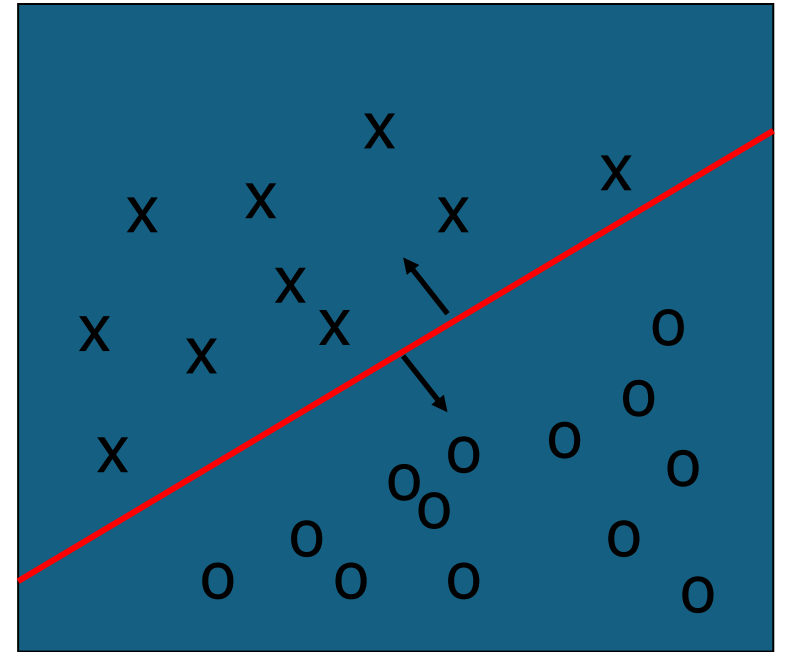
# Classification—A Two-Step Process

- Model construction: describing a set of predetermined classes
  - Each tuple/sample is assumed to belong to a predefined class, as determined by the class label attribute
  - The set of tuples used for model construction is training set
  - The model is represented as classification rules, decision trees, or mathematical formulae
- Model usage: for classifying future or unknown objects
  - Estimate accuracy of the model
    - The known label of test sample is compared with the classified result from the model
    - Accuracy rate is the percentage of test set samples that are correctly classified by the model
    - Test set is independent of training set (otherwise overfitting)
  - If the accuracy is acceptable, use the model to classify new data
- Note: If *the test set* is used to select models, it is called validation (test) set
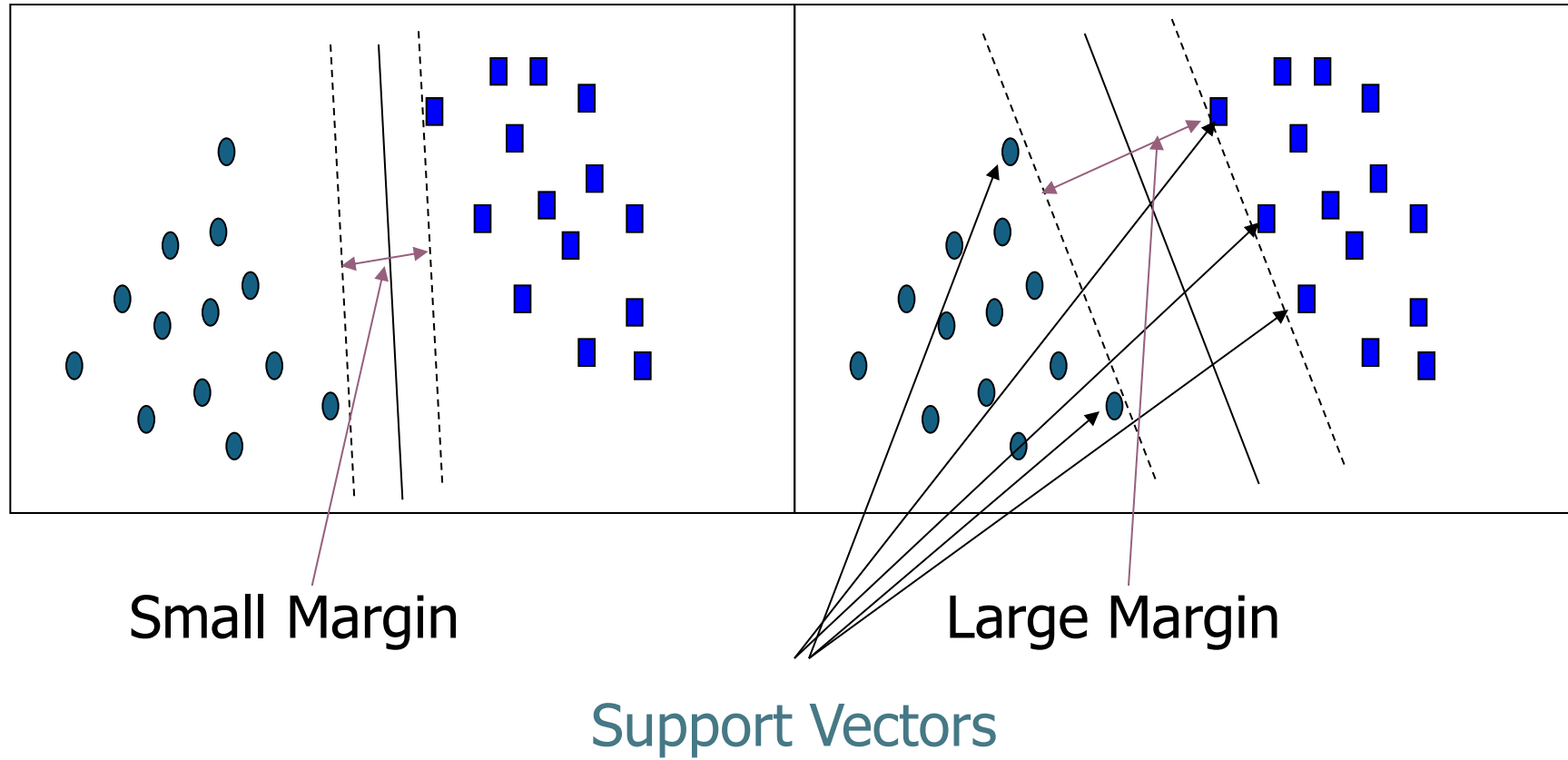
# Classification: A Mathematical Mapping

- Classification: predicts categorical class labels
  - E.g., Personal homepage classification
    - $x_i = (x_1, x_2, x_3, \ldots)$, $y_i = +1$ or $-1$
    - $x_1$ : # of word "homepage"
    - $x_2$ : # of word "welcome"
- Mathematically, $x \in X = \mathfrak{R}^n$, $y \in Y = \{+1, -1\}$,
  - We want to derive a function $f: X \rightarrow Y$
- Linear Classification
  - Binary Classification problem
  - Data above the red line belongs to class 'x'
  - Data below red line belongs to class 'o'
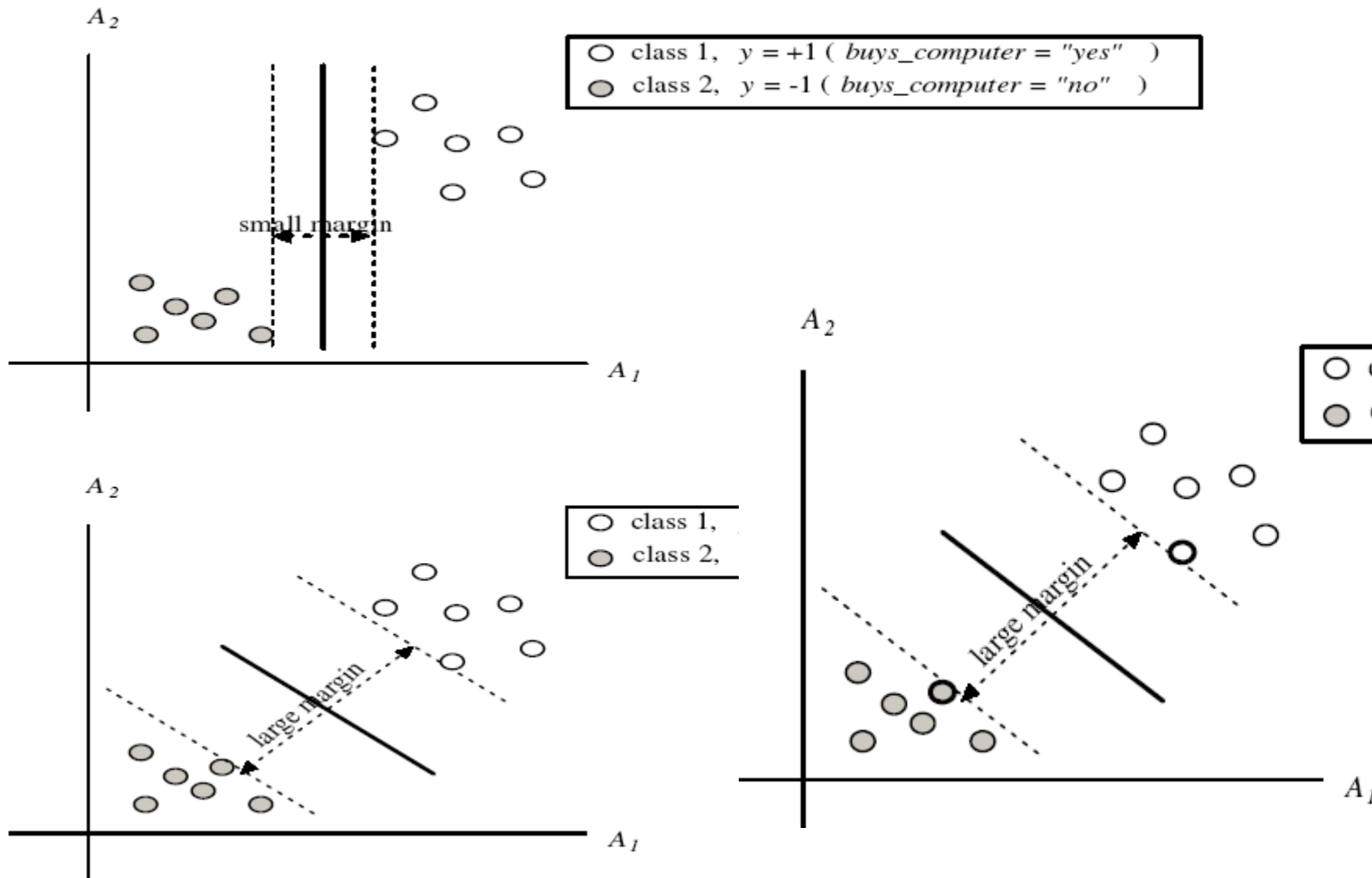  - Examples: SVM, Perceptron, Probabilistic Classifiers

# SVM—Support Vector Machines

- A relatively new classification method for both <u>linear and nonlinear</u> data

- It uses a <u>nonlinear mapping</u> to transform the original training data into a higher dimension.

- With the new dimension, it searches for the linear optimal separating **hyperplane** (i.e., "decision boundary").

- With an appropriate nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated by a hyperplane.

- SVM finds this hyperplane using **support vectors** ("essential" training tuples) and **margins** (defined by the support vectors).

- <u>Features</u>: training can be slow but accuracy is high owing to their ability to model complex nonlinear decision boundaries (margin maximization)

- <u>Used for</u>: classification and numeric prediction

- <u>Applications</u>: handwritten digit recognition, object recognition, speaker identification, benchmarking time-series prediction tests
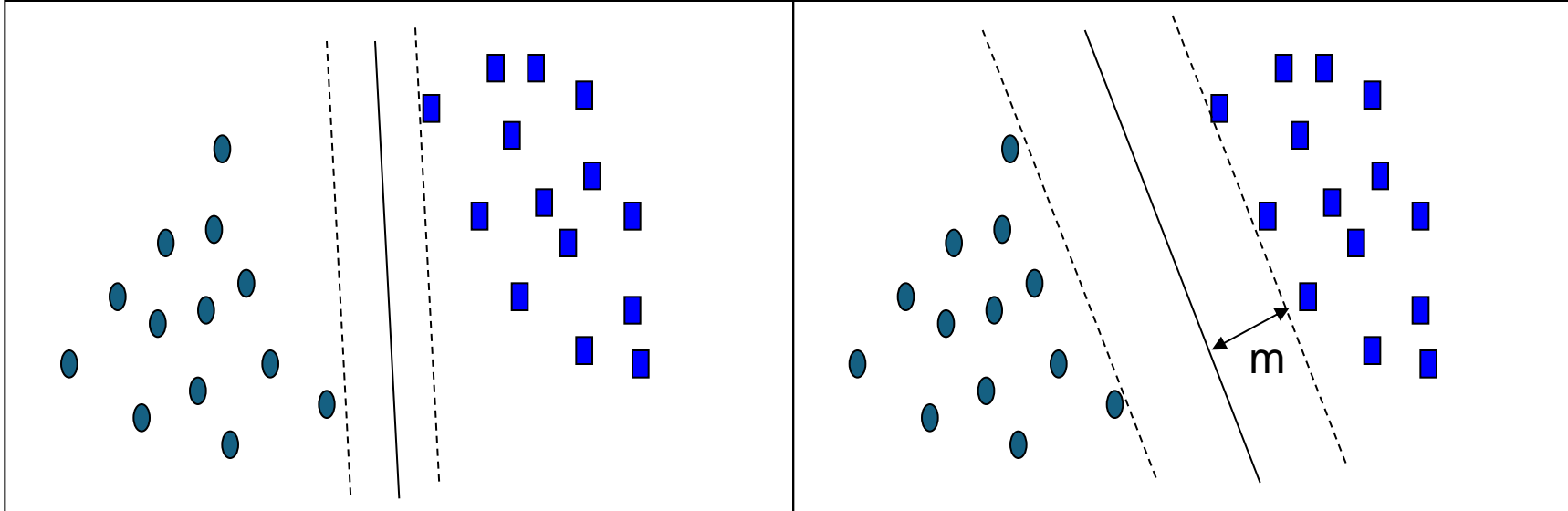
# SVM—General Philosophy



Small Margin

Large Margin

Support Vectors

# SVM—Margins and Support Vectors

# SVM—When Data Is Linearly Separable



Let data D be $(\mathbf{X}_1, y_1), ..., (\mathbf{X}_{|D|}, y_{|D|})$, where $\mathbf{X}_i$ is the set of training tuples associated with the class labels $y_i$

There are infinite lines (<u>hyperplanes</u>) separating the two classes but we want to <u>find the best one</u> (the one that minimizes classification error on unseen data)

*SVM searches for the hyperplane with the largest margin*, i.e., **maximum marginal hyperplane** (MMH)

# SVM—Linearly Separable

- A separating hyperplane can be written as

    $\mathbf{W} \bullet \mathbf{X} + b = 0$

    where $\mathbf{W}=\{w_1, w_2, ..., w_n\}$ is a weight vector and b a scalar (bias)

- For 2-D it can be written as

    $w_0 + w_1 x_1 + w_2 x_2 = 0$

- The hyperplane defining the sides of the margin:

    $H_1: w_0 + w_1 x_1 + w_2 x_2 \geq 1$     for $y_i = +1$, and

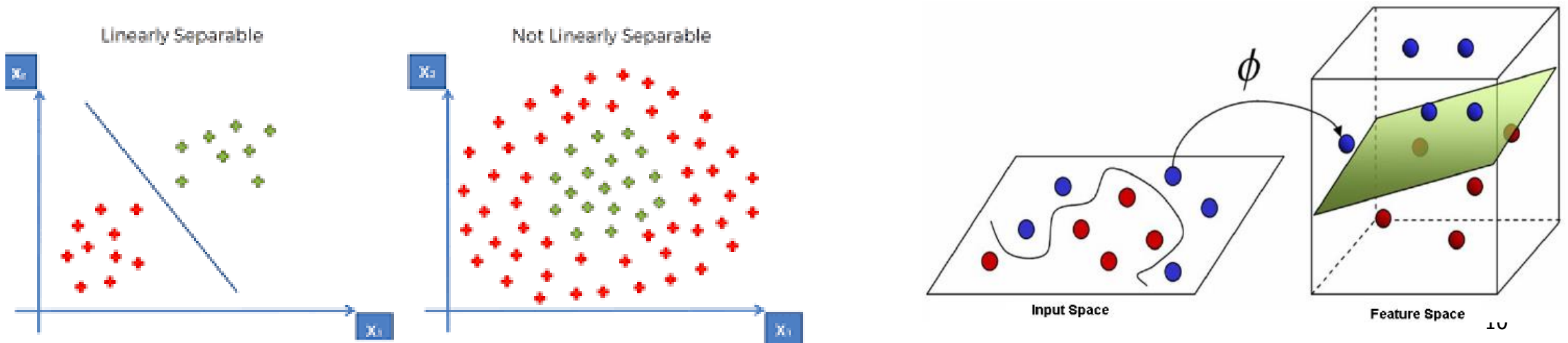    $H_2: w_0 + w_1 x_1 + w_2 x_2 \leq -1$ for $y_i = -1$

- Any training tuples that fall on hyperplanes $H_1$ or $H_2$ (i.e., the sides defining the margin) are **support vectors**

- This becomes a **constrained (convex) quadratic optimization** problem: Quadratic objective function and linear constraints → *Quadratic Programming (QP)* → Lagrangian multipliers

# Why Is SVM Effective on High Dimensional Data?

- The **complexity** of trained classifier is characterized by the # of support vectors rather than the dimensionality of the data

- The **support vectors** are the essential or critical training examples —they lie closest to the decision boundary (MMH)

- If all other training examples are removed and the training is repeated, the same separating hyperplane would be found

- The number of support vectors found can be used to compute an (upper) bound on the expected error rate of the SVM classifier, which is independent of the data dimensionality

- Thus, an SVM with a small number of support vectors can have good generalization, even when the dimensionality of the data is high

# SVM—Linearly Inseparable

- Transform the original input data into a higher dimensional space.

- Non-linear SVMs use kernel functions to transform data into higher-dimensional spaces, allowing for the creation of complex, non-linear decision boundaries that are not possible with linear SVMs, enabling effective classification of non-linearly separable data.

- Instead of computing the dot product on the transformed data, it is math. equivalent to applying a kernel function $K(\mathbf{X_i}, \mathbf{X_j})$ to the original data, i.e., $K(\mathbf{X_i}, \mathbf{X_j}) = \Phi(\mathbf{X_i})\,\Phi(\mathbf{X_j})$

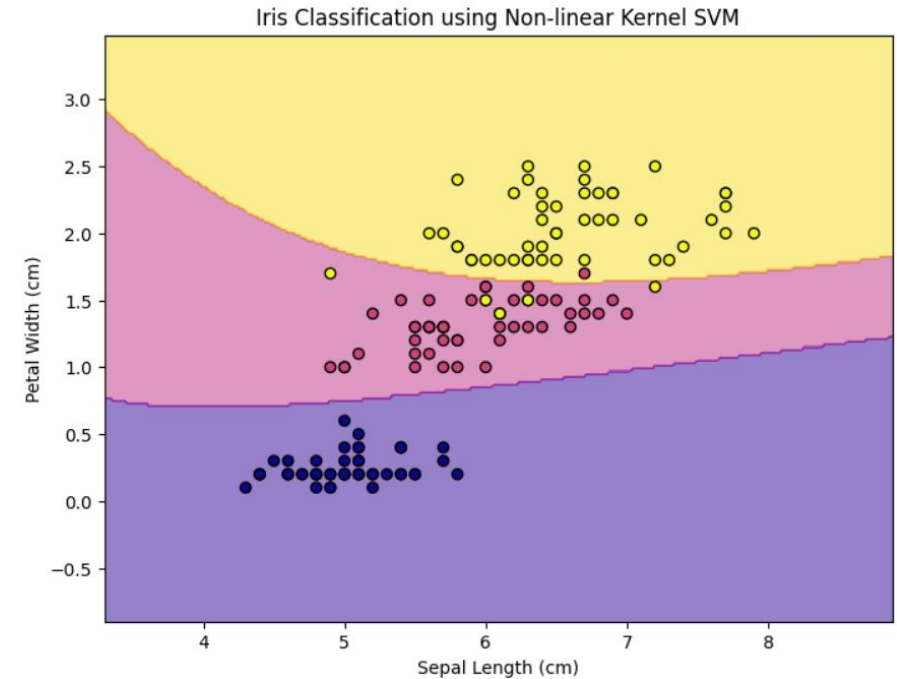- Search for a linear separating hyperplane in the new space.



Linearly Separable

Not Linearly Separable

Input Space

Feature Space

# SVM: Different Kernel functions

Iris Classification using Non-linear Kernel SVM

Typical Kernel Functions

Polynomial kernel of degree $h$ : $\quad K(X_i, X_j) = (X_i \cdot X_j + 1)^h$

Gaussian radial basis function kernel : $\quad K(X_i, X_j) = e^{-\|X_i - X_j\|^2 / 2\sigma^2}$

Sigmoid kernel : $\quad K(X_i, X_j) = \tanh(\kappa X_i \cdot X_j - \delta)$

SVM can also be used for

- classifying multiple (> 2) classes and
- regression analysis

Regression

Classification

# Model Evaluation and Selection

- Evaluation metrics: How can we measure accuracy?  Other metrics to consider?

- Use **validation/test set** of class-labeled tuples instead of training set when assessing accuracy

- Methods for estimating a classifier's accuracy:
  - Holdout method, random subsampling
  - Cross-validation
  - Bootstrap

# Classifier Evaluation Metrics: Confusion Matrix

**Confusion Matrix:**

| Actual class\Predicted class | $C_1$ | $\neg C_1$ |
|---|---|---|
| $C_1$ | **True Positives (TP)** | **False Negatives (FN)** |
| $\neg C_1$ | **False Positives (FP)** | **True Negatives (TN)** |

**Example of Confusion Matrix:**

| Actual class\Predicted class | buy_computer = yes | buy_computer = no | Total |
|---|---|---|---|
| buy_computer = yes | **6954** | **46** | 7000 |
| buy_computer = no | **412** | **2588** | 3000 |
| Total | 7366 | 2634 | 10000 |

- Given *m* classes, an entry, ***CM**_{i,j}* in a **confusion matrix** indicates # of tuples in class *i* that were labeled by the classifier as class *j*

- May have extra rows/columns to provide totals

# Classifier Evaluation Metrics

- **Classifier Accuracy,** or recognition rate: percentage of test set tuples that are correctly classified.

  **Accuracy = (TP + TN)/All**

- **Error rate = *1 – accuracy* = (FP + FN)/All**

- **Class Imbalance Problem**:
  - One class may be *rare*, e.g. fraud, or HIV-positive
  - Significant *majority of the negative class* and minority of the positive class

- **Precision**: exactness – what % of tuples that the classifier labeled as positive are actually positive.

- **Recall:** completeness – what % of positive tuples did the classifier label as positive?

- Perfect score is 1.0

$$precision = \frac{TP}{TP + FP} \qquad recall = \frac{TP}{TP + FN}$$

- Inverse relationship between precision & recall
- ***F* measure (*F₁* - score)**: harmonic mean of precision and recall.

$$F = \frac{2 \times precision \times recall}{precision + recall}$$

# Classifier Evaluation Metrics

## Classification Metrics Examples

| Actual | Predicted | | Row Totals |
|---|---|---|---|
| | Positive | Negative | Row Totals |
| Positive | 60 | 10 | 70 |
| Negative | 5 | 25 | 30 |
| Col Totals | 65 | 35 | 100 |

$$\text{Recall} = \frac{60}{70} = 0.857$$

$$\text{Specificity} = \frac{25}{30} = 0.833$$

$$\text{Error} = \frac{15}{100} = 15\%$$

$$\text{Accuracy} = \frac{85}{100} = 85\%$$

$$\text{Precision} = \frac{60}{65} = 0.923$$

$$F = 2 * \frac{0.857 * 0.923}{0.857 + 0.923} = 0.889$$

# Holdout & Cross-Validation Methods
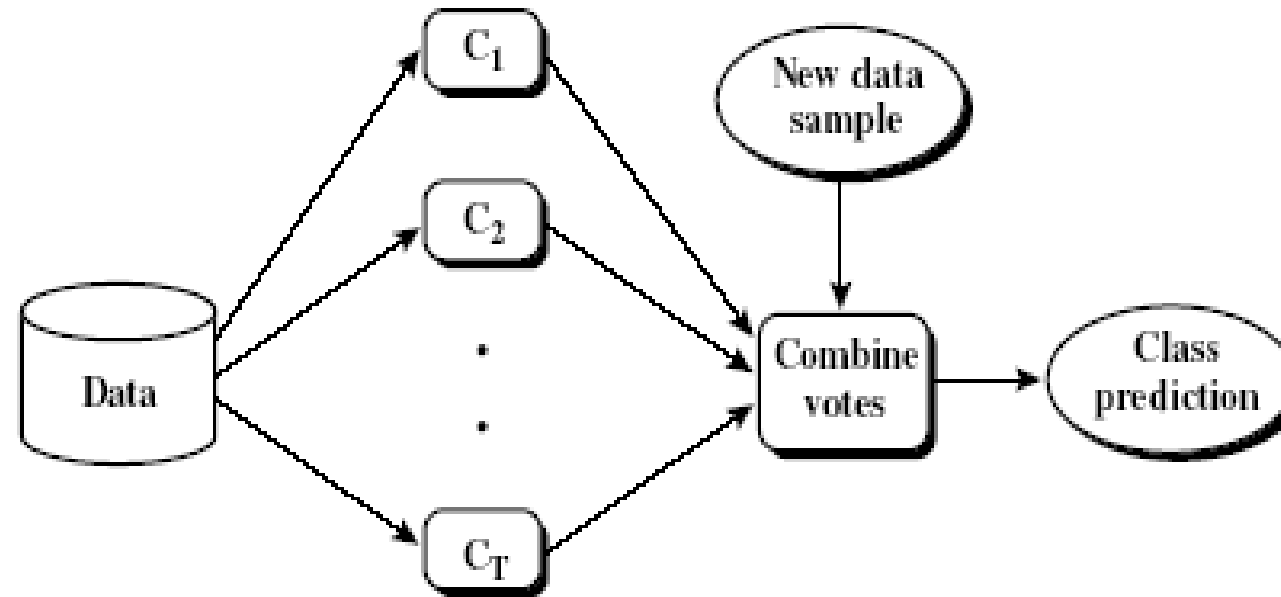
- **Holdout method**

  - Given data is randomly partitioned into two independent sets

    - Training set (e.g., 2/3) for model construction

    - Test set (e.g., 1/3) for accuracy estimation

  - <u>Random sampling</u>: a variation of holdout

    - Repeat holdout k times, accuracy = avg. of the accuracies obtained

- **Cross-validation** (*k*-fold, where k = 10 is most popular)

  - Randomly partition the data into *k mutually exclusive* subsets, each approximately equal size

  - At *i*-th iteration, use $D_i$ as test set and others as training set

  - <u>Leave-one-out</u>: *k* folds where *k* = # of tuples, for small sized data

  - <u>**\*Stratified cross-validation\***</u>: folds are stratified so that class dist. in each fold is approx. the same as that in the initial data

# Ensemble Methods: Increasing the Accuracy



- Ensemble methods
  - Use a combination of models to increase accuracy
  - Combine a series of k learned models, $M_1$, $M_2$, …, $M_k$, with the aim of creating an improved model M*
- Popular ensemble methods
  - Bagging: averaging the prediction over a collection of classifiers
  - Boosting: weighted vote with a collection of classifiers
  - Ensemble: combining a set of heterogeneous classifiers

# Random Forest (Breiman 2001)

- Random Forest:

    - Each classifier in the ensemble is a *decision tree* classifier and is generated using a random selection of attributes at each node to determine the split

    - During classification, each tree votes and the most popular class is returned

- Two Methods to construct Random Forest:

    - Forest-RI (*random input selection*):  Randomly select, at each node, F attributes as candidates for the split at the node. The CART methodology is used to grow the trees to maximum size

    - Forest-RC (*random linear combinations*):  Creates new attributes (or features) that are a linear combination of the existing attributes (reduces the correlation between individual classifiers)

- Comparable in accuracy to Adaboost, but more robust to errors and outliers

- Insensitive to the number of attributes selected for consideration at each split, and faster than bagging or boosting