# DATA MINING

Pertemuan V

# DISCRIMINATIVE CLASSIFIER

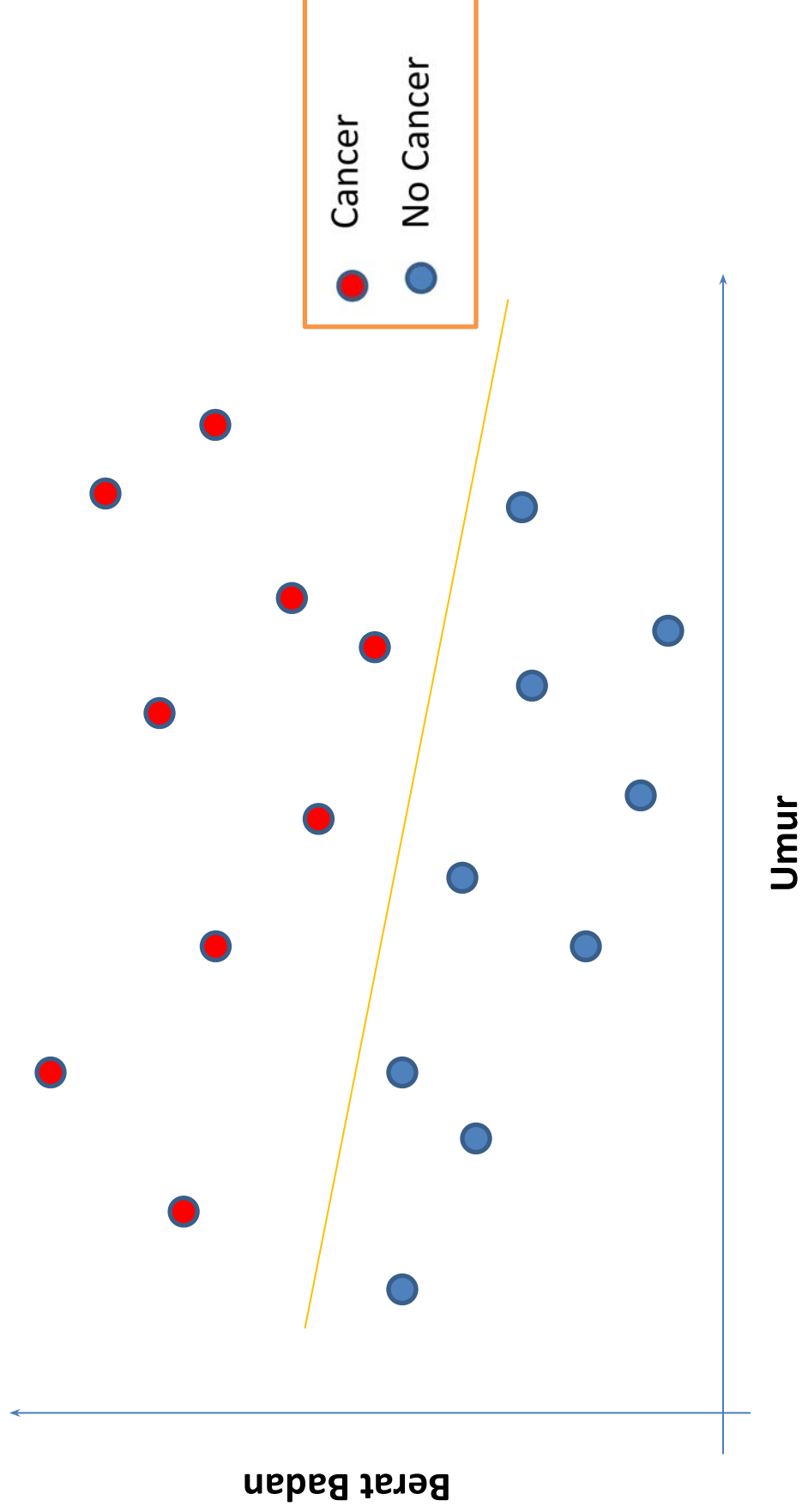# Generative Classifier vs Discriminative Classifier

**Discriminative**

**Generative**

- Generative Classifiers tries to model class, i.e., what are the features of the class. Ex: Naïve Bayes

- Discriminative Classifiers learn what the features in the input are most useful to distinguish between the various possible classes. Ex: SVM

# Example



Berat Badan

Umur

Cancer

No Cancer

# Discriminative Classifiers

- Advantages
  - Prediction accuracy is generally high
  - Robust, works when training examples contain errors
  - Fast evaluation of the learned target function
    - Bayesian networks are normally slow
- Criticism
  - Long training time
  - Difficult to understand the learned function (weights)
    - Bayesian networks can be used easily for pattern discovery
  - Not easy to incorporate domain knowledge
    - Easy in the form of priors on the data or distributions

# SUPPORT VECTOR MACHINE
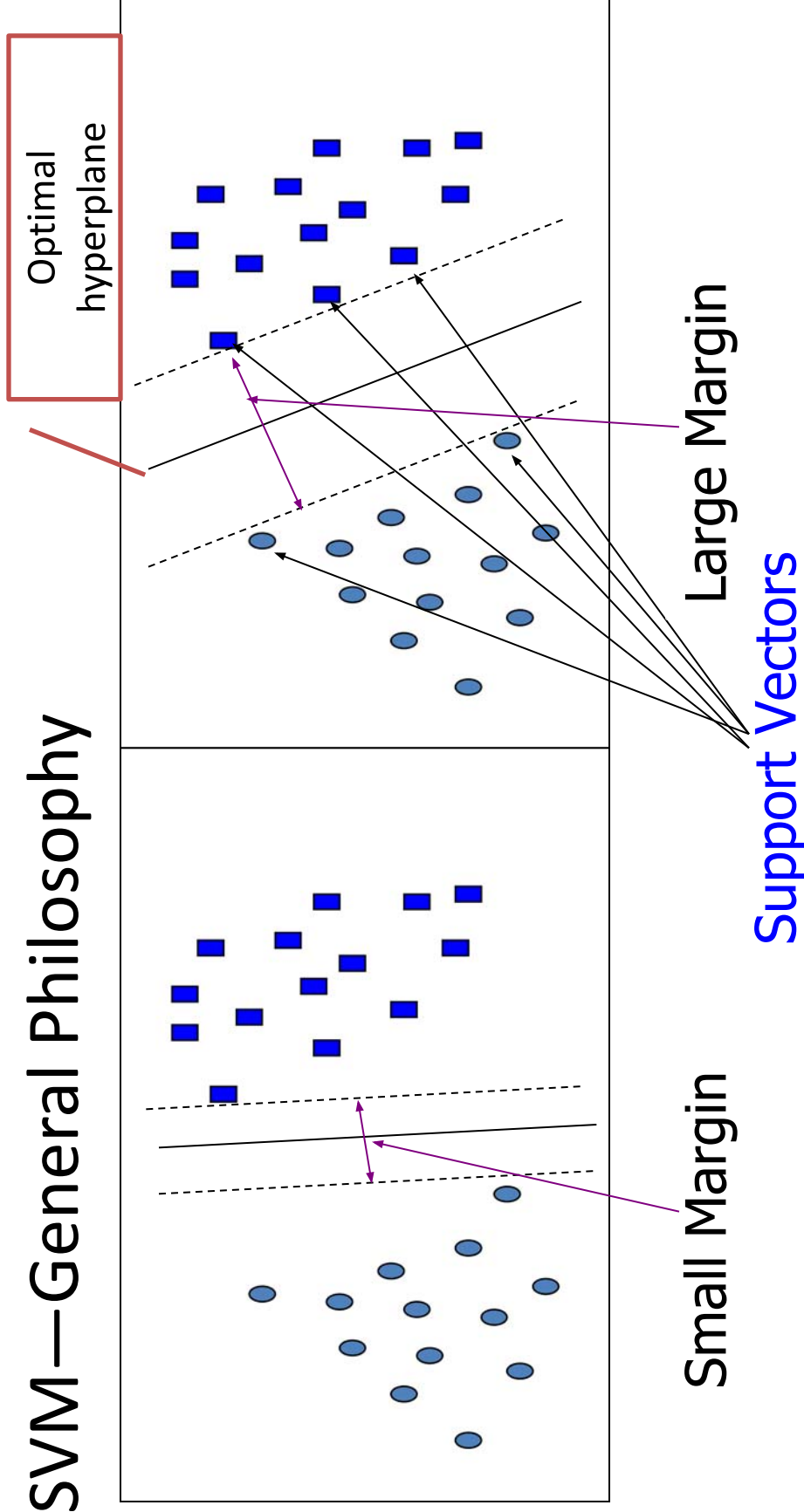
# SVM—Support Vector Machines

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

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

- 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)
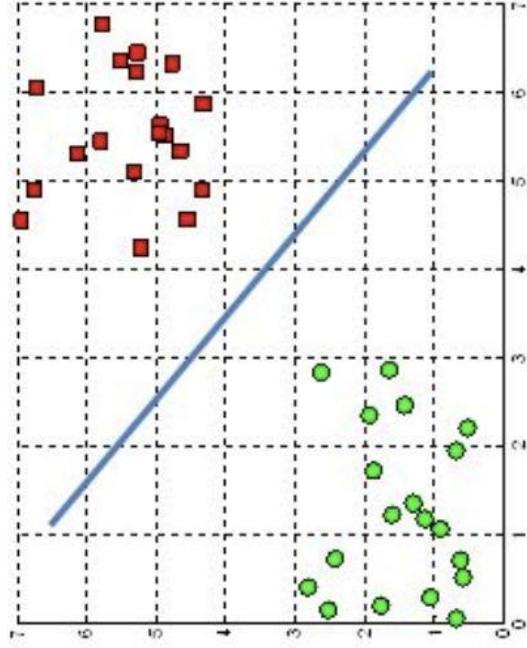
# SVM—History and Applications

- Vapnik and colleagues (1992)—groundwork from Vapnik & Chervonenkis' statistical learning theory in 1960s

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

- Used for: classification and numeric prediction

- Applications:

  – handwritten digit recognition, object recognition, speaker identification, benchmarking time-series prediction tests

# SVM—General Philosophy
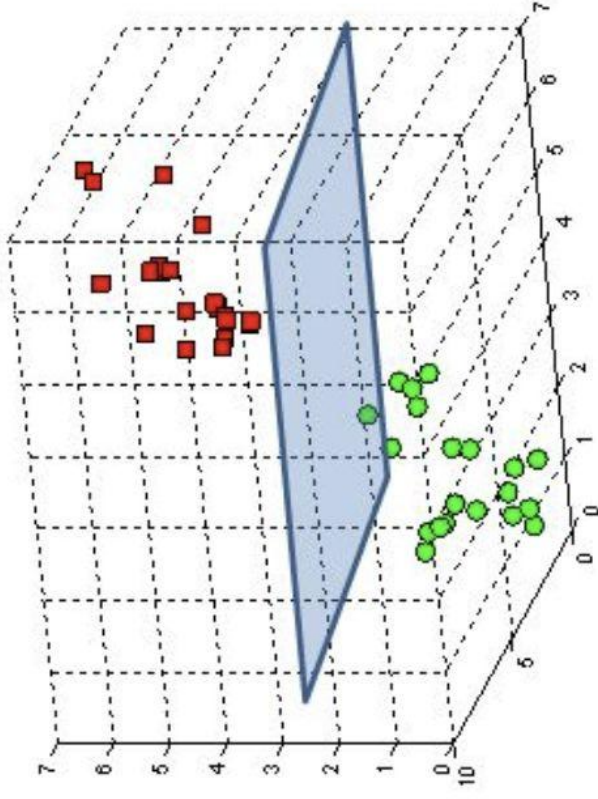
Optimal hyperplane

Support Vectors

Large Margin

Small Margin

The objective of the support vector machine algorithm is to find a hyperplane in an N-dimensional space(N — the number of features) that distinctly classifies the data points

A hyperplane in $\mathbb{R}^2$ is a line

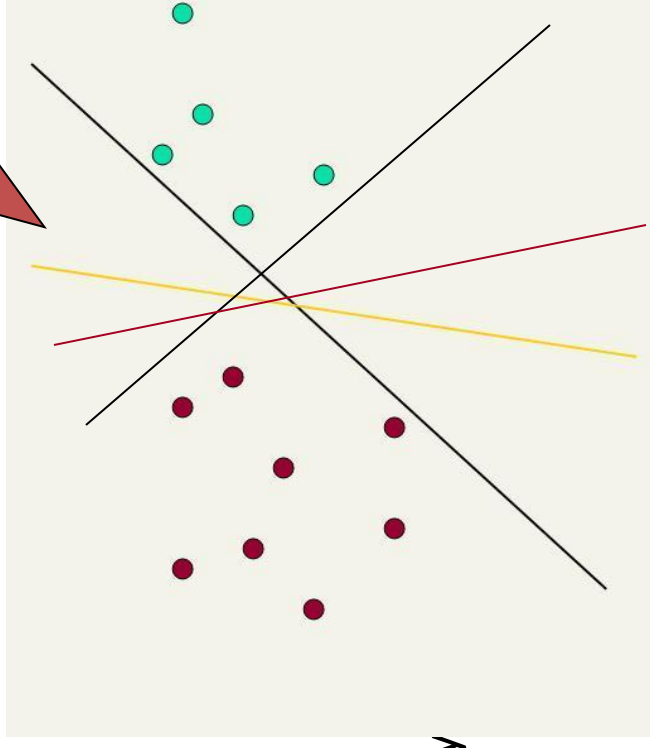A hyperplane in $\mathbb{R}^3$ is a plane

a **hyperplane** is a subspace whose <u>dimension</u> is one less than that of its <u>ambient space.</u>

# Linear classifiers: Which Hyperplane?

To separate the two classes of data points, there are many possible hyperplanes that could be chosen.
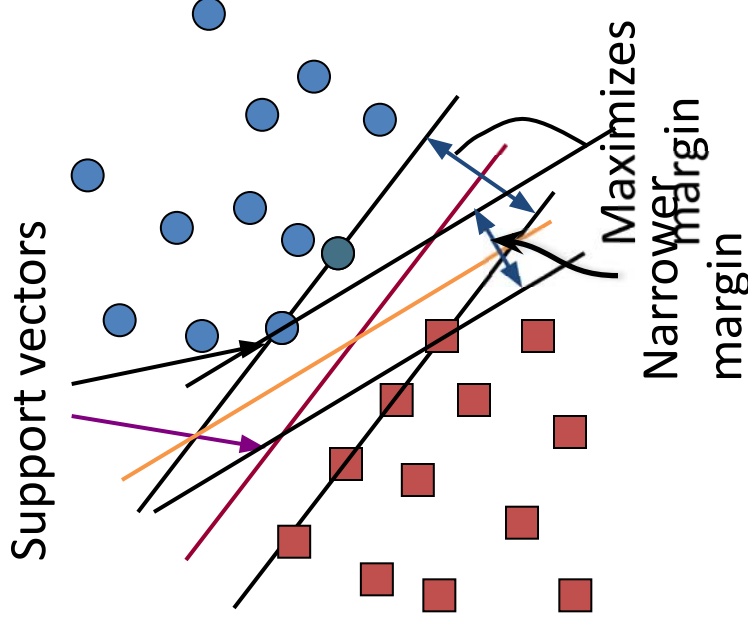
- Lots of possible solutions for *a, b, c.*
- Some methods find a separating hyperplane, but not the optimal one
  - E.g., perceptron
- Support Vector Machine (SVM) finds an optimal* solution.
  - Maximizes the distance between the hyperplane and the "difficult points" close to decision boundary
  - One intuition: if there are no points near the decision surface, then there are no very uncertain classification decisions

These lines represent the decision boundary:
$ax + by - c = 0$

12

# Support Vector Machine (SVM)

- SVMs maximize the *margin* around the separating hyperplane.
  - A.k.a. large margin classifiers
- The decision function is fully specified by a subset of training samples, *the support vectors*.
- Solving SVMs is a *quadratic programming* problem
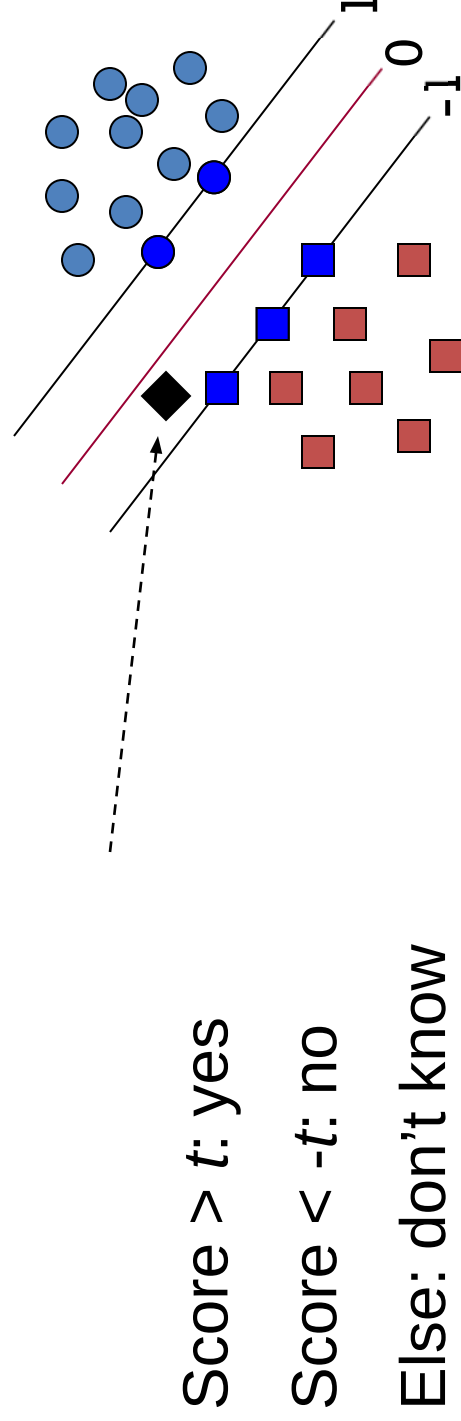- Seen by many as the most successful current text classification method *

Support vectors

Maximizes margin

Narrower margin

# SVM—Linearly Separable

- A separating hyperplane can be written as

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

  where $\mathbf{W}=\{w_1, w_2, \ldots, w_n\}$ 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

# Classification with SVMs

- Given a new point **x**, we can score its projection onto the hyperplane normal:

  – I.e., compute score: $\mathbf{w}^T\mathbf{x} + b = \Sigma \alpha_i y_i \mathbf{x}_i^T\mathbf{x} + b$

    - Decide class based on whether < or > 0

  – Can set confidence threshold $t$.



Score > $t$: yes
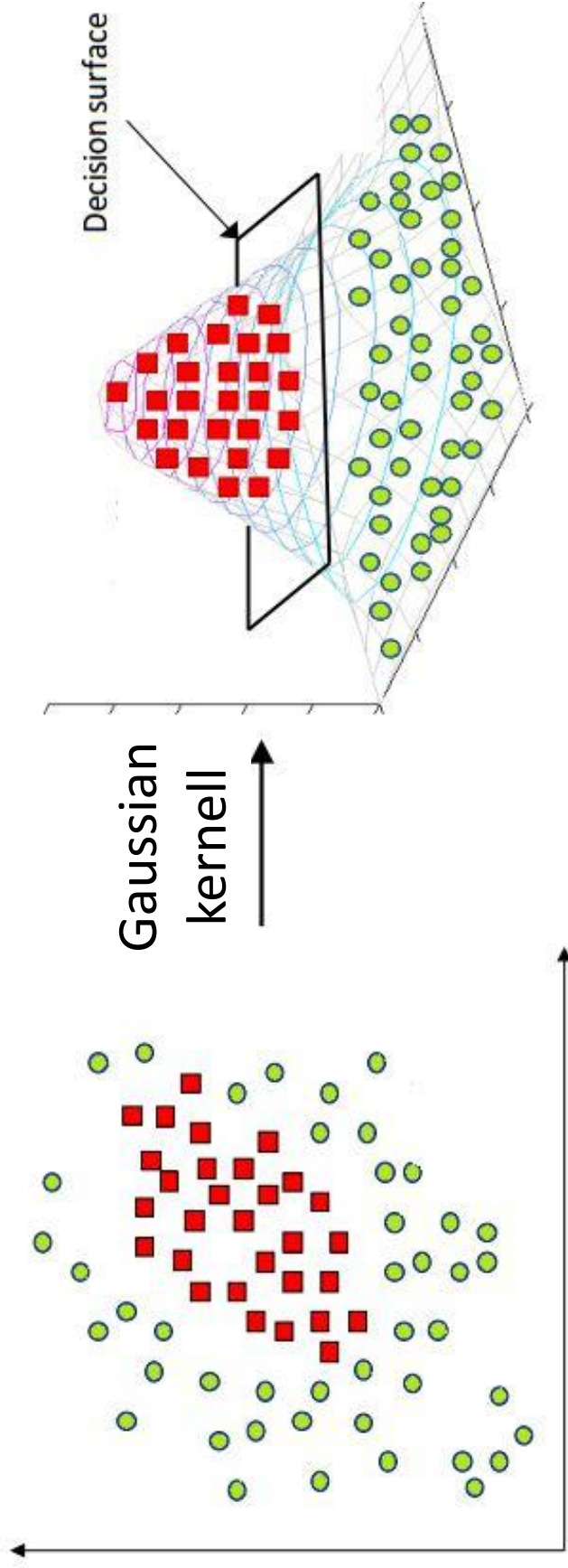
Score < -$t$: no

Else: don't know

# 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

- Search for a linear separating hyperplane in the new space

Gaussian kernell

Decision surface

# SVM: Different Kernel functions

- 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., $K(\mathbf{X_i}, \mathbf{X_j}) = \Phi(\mathbf{X_i}) \, \Phi(\mathbf{X_j})$

- Typical Kernel Functions

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

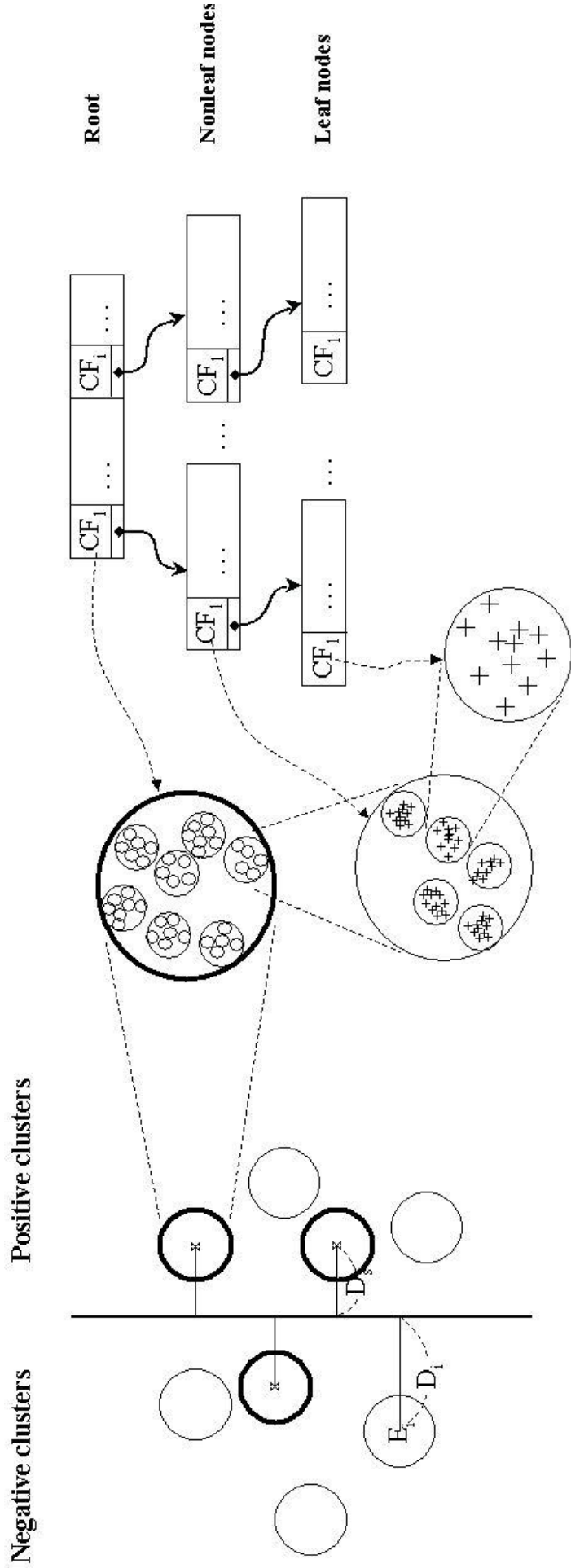Gaussian radial basis function kernel : $\quad K(X_i, X_j) = e^{-\|x_i - x_j\|^2 / 2\sigma^2}$

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 for regression analysis (with additional parameters)

18

# Scaling SVM by Hierarchical Micro-Clustering

- SVM is not scalable to the number of data objects in terms of training time and memory usage

- H. Yu, J. Yang, and J. Han, "Classifying Large Data Sets Using SVM with Hierarchical Clusters", KDD'03)

- CB-SVM (Clustering-Based SVM)

  - Given limited amount of system resources (e.g., memory), maximize the SVM performance in terms of accuracy and the training speed

  - Use micro-clustering to effectively reduce the number of points to be considered

  - At deriving support vectors, de-cluster micro-clusters near "candidate vector" to ensure high classification accuracy

# CF-Tree: Hierarchical Micro-cluster



- Read the data set once, construct a statistical summary of the data (i.e., hierarchical clusters) given a limited amount of memory
- Micro-clustering: Hierarchical indexing structure
  - provide finer samples closer to the boundary and coarser samples farther from the boundary

# Multiclass Classification

- Classification involving more than two classes (i.e., > 2 Classes)

- Method 1. **One-vs.-all** (OVA): Learn a classifier one at a time
  - Given m classes, train m classifiers: one for each class
  - Classifier j: treat tuples in class j as *positive* & all others as *negative*
  - To classify a tuple **X**, the set of classifiers vote as an ensemble

- Method 2. **All-vs.-all** (AVA): Learn a classifier for each pair of classes
  - Given m classes, construct m(m-1)/2 binary classifiers
  - A classifier is trained using tuples of the two classes
  - To classify a tuple **X**, each classifier votes. X is assigned to the class with maximal vote

- Comparison
  - All-vs.-all tends to be superior to one-vs.-all
  - Problem: Binary classifier is sensitive to errors, and errors affect vote count

# NEXT: KNN