

Big Data Science

Support Vector Machine

Incheon Paik

Contents

- ◆ **Classification**
- ◆ **Support Vector Machine (SVM)**

Classification: Definition

- ◆ Given a collection of records (*training set*)
 - Each record contains a set of *attributes*, one of the attributes is the *class*.
- ◆ Find a *model* for class attribute as a function of the values of other attributes.
- ◆ Goal: previously unseen records should be assigned a class as accurately as possible.
 - A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

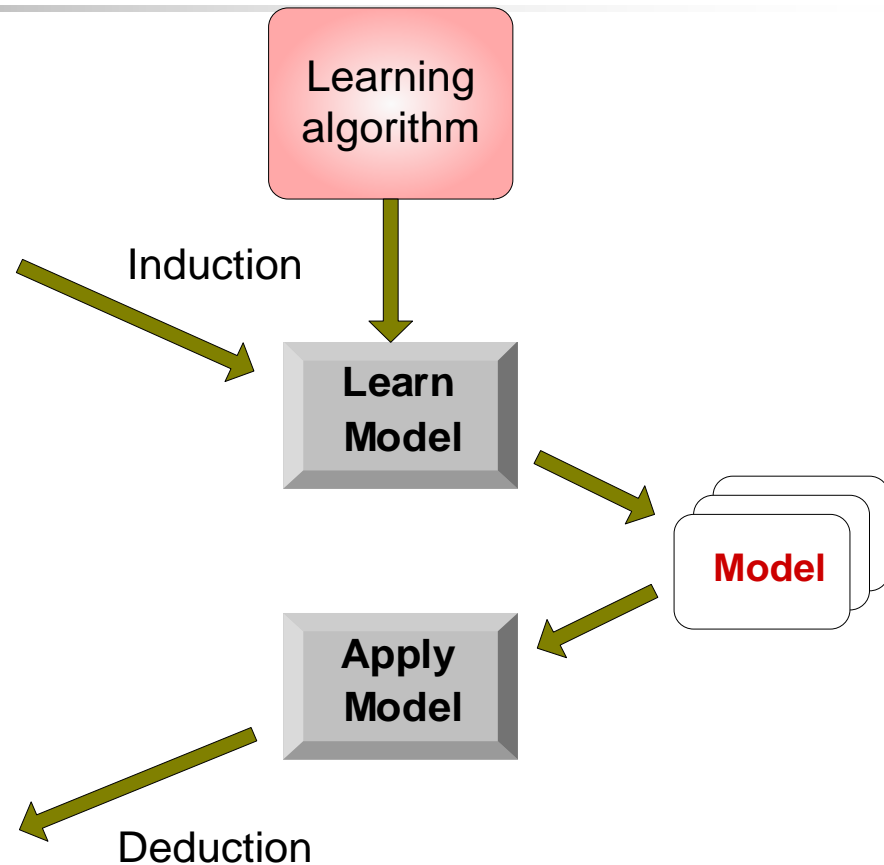
Illustrating Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

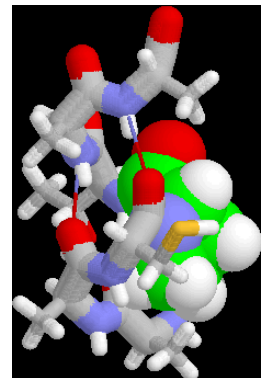
Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



Examples of Classification Task

- ◆ Predicting tumor cells as benign or malignant
- ◆ Classifying credit card transactions as legitimate or fraudulent
- ◆ Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil
- ◆ Categorizing news stories as finance, weather, entertainment, sports, etc



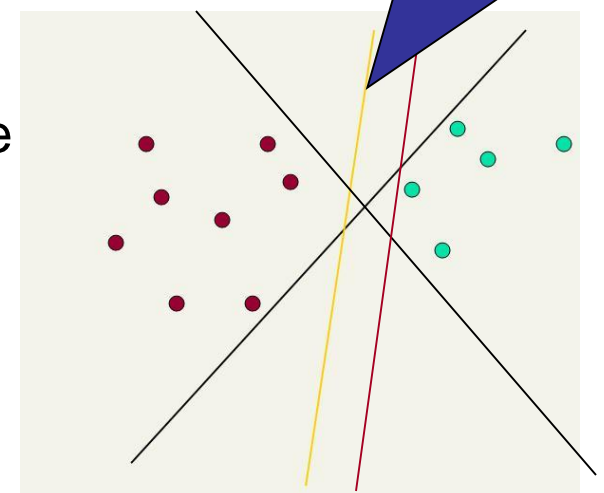
Classification Techniques

- ◆ Decision Tree based Methods
- ◆ Rule-based Methods
- ◆ Memory based reasoning
- ◆ Neural Networks
- ◆ Naïve Bayes and Bayesian Belief Networks
- ◆ Support Vector Machines

Linear classifiers: Which Hyperplane?

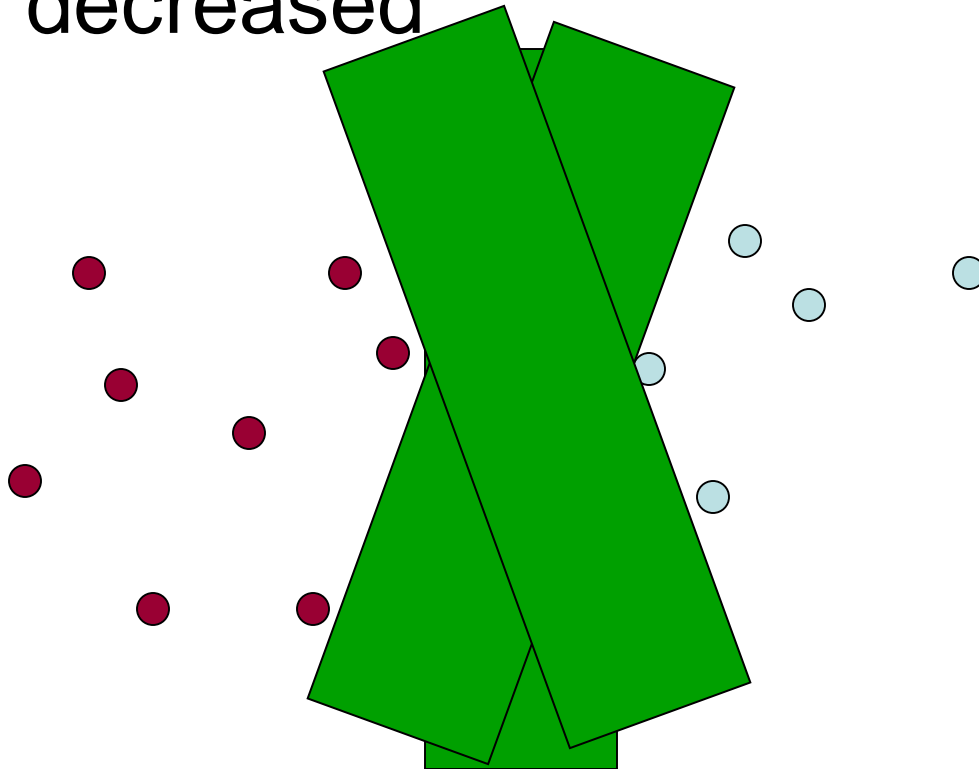
- ◆ Lots of possible solutions for a , b , c .
- ◆ Some methods find a separating hyperplane, but not the optimal one [according to some criterion of expected goodness]
 - E.g., perceptron
- ◆ Support Vector Machine (SVM) finds an optimal* solution.
 - Maximizes the distance between the hyperplane and 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

This line represents the decision boundary:
 $ax + by - c = 0$



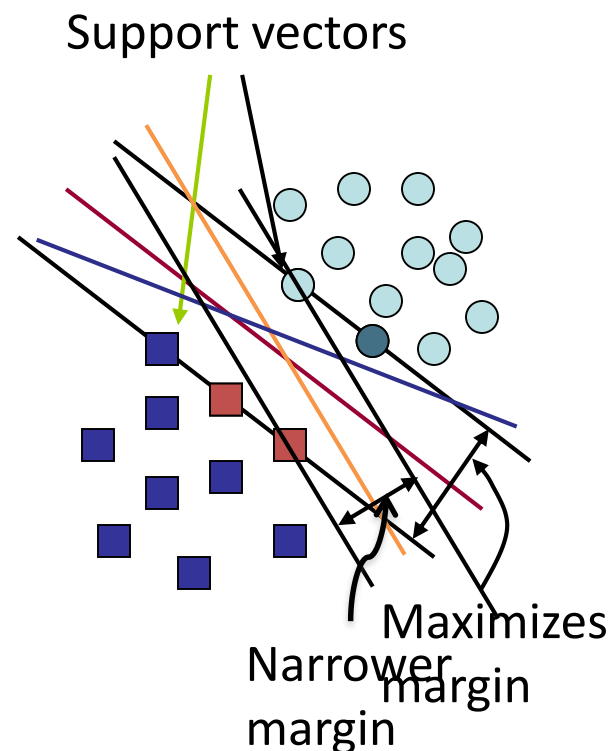
Another intuition

- ◆ If you have to place a fat separator between classes, you have less choices, and so the capacity of the model has been decreased



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*



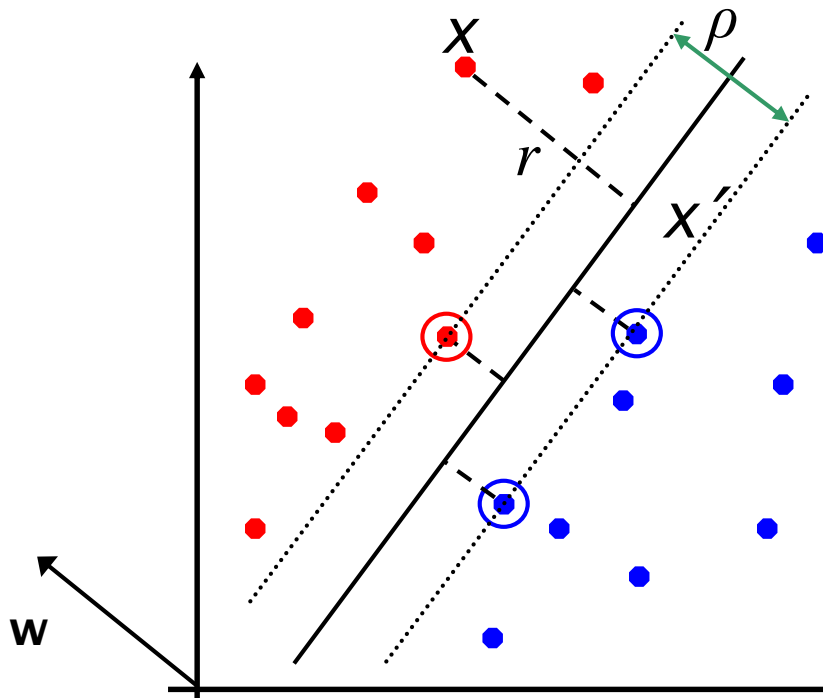
*but other discriminative methods often perform very similarly

Maximum Margin: Formalization

- ◆ \mathbf{w} : decision hyperplane normal vector
- ◆ \mathbf{x}_i : data point i
- ◆ y_i : class of data point i (+1 or -1) NB: Not 1/0
- ◆ Classifier is: $f(\mathbf{x}_i) = \text{sign}(\mathbf{w}^T \mathbf{x}_i + b)$
- ◆ Functional margin of \mathbf{x}_i is: $y_i (\mathbf{w}^T \mathbf{x}_i + b)$
 - But note that we can increase this margin simply by scaling \mathbf{w} , \mathbf{b}
- ◆ Functional margin of dataset is twice the minimum functional margin for any point
 - The factor of 2 comes from measuring the whole width of the margin

Geometric Margin

- ◆ Distance from example to the separator is $r = y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$
- ◆ Examples closest to the hyperplane are **support vectors**.
- ◆ **Margin** ρ of the separator is the width of separation between support vectors of classes.



Derivation of finding r :

Dotted line $\mathbf{x}' - \mathbf{x}$ is perpendicular to decision boundary so parallel to \mathbf{w} .

Unit vector is $\mathbf{w}/\|\mathbf{w}\|$, so line is $r\mathbf{w}/\|\mathbf{w}\|$.
 $\mathbf{x}' = \mathbf{x} - yr\mathbf{w}/\|\mathbf{w}\|$.

\mathbf{x}' satisfies $\mathbf{w}^T \mathbf{x}' + b = 0$.

So $\mathbf{w}^T (\mathbf{x} - yr\mathbf{w}/\|\mathbf{w}\|) + b = 0$

Recall that $\|\mathbf{w}\| = \sqrt{\mathbf{w}^T \mathbf{w}}$.

So $\mathbf{w}^T \mathbf{x} - yr\|\mathbf{w}\| + b = 0$

So, solving for r gives:

$$r = y(\mathbf{w}^T \mathbf{x} + b)/\|\mathbf{w}\|$$

Linear SVM Mathematically

The linearly separable case

- ◆ Assume that all data is at least distance 1 from the hyperplane, then the following two constraints follow for a training set $\{(\mathbf{x}_i, y_i)\}$

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \quad \text{if } y_i = 1$$

$$\mathbf{w}^T \mathbf{x}_i + b \leq -1 \quad \text{if } y_i = -1$$

- ◆ For support vectors, the inequality becomes an equality
- ◆ Then, since each example's distance from the hyperplane is

$$r = y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$$

- ◆ The margin is:

$$r = \frac{2}{\|\mathbf{w}\|}$$

Linear Support Vector Machine (SVM)

- ◆ **Hyperplane**

$$\mathbf{w}^T \mathbf{x} + b = 0$$

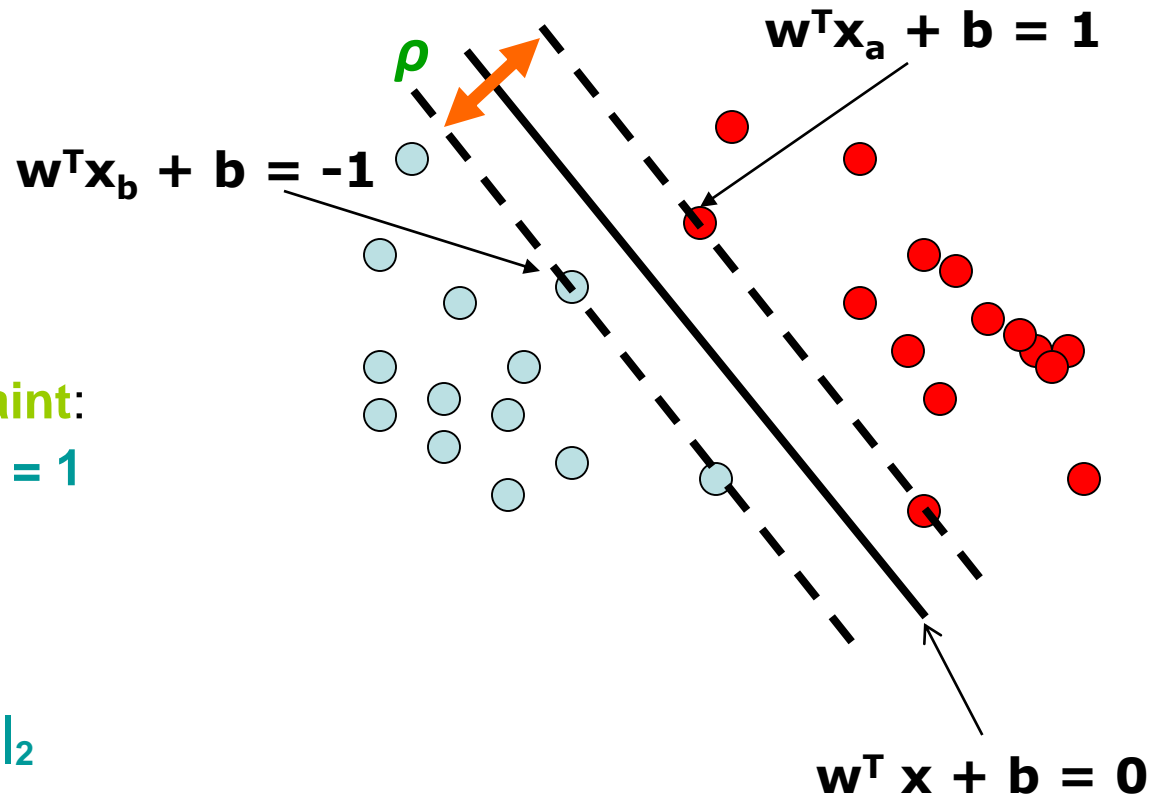
- ◆ **Extra scale constraint:**

$$\min_{i=1,\dots,n} |\mathbf{w}^T \mathbf{x}_i + b| = 1$$

- ◆ This implies:

$$\mathbf{w}^T (\mathbf{x}_a - \mathbf{x}_b) = 2$$

$$\rho = \|\mathbf{x}_a - \mathbf{x}_b\|_2 = 2 / \|\mathbf{w}\|_2$$



Linear SVMs Mathematically (cont.)

- ◆ Then we can formulate the *quadratic optimization problem*:

Find \mathbf{w} and b such that

$$r = \frac{2}{\|\mathbf{w}\|} \text{ is maximized; and for all } \{(\mathbf{x}_i, y_i)\}$$

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \text{ if } y_i = 1; \quad \mathbf{w}^T \mathbf{x}_i + b \leq -1 \text{ if } y_i = -1$$

- ◆ A better formulation ($\min \|\mathbf{w}\| = \max 1/\|\mathbf{w}\|$):

Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ is minimized;}$$

$$\text{and for all } \{(\mathbf{x}_i, y_i)\}: \quad y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Solving the Optimization Problem

Find \mathbf{w} and b such that
 $\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$ is minimized;
and for all $\{(\mathbf{x}_i, y_i)\}$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

- ◆ This is now optimizing a *quadratic* function subject to *linear* constraints
- ◆ Quadratic optimization problems are a well-known class of mathematical programming problem, and many (intricate) algorithms exist for solving them (with many special ones built for SVMs), ex) SMO
- ◆ The solution involves constructing a *dual problem* where a *Lagrange multiplier* α_i is associated with every constraint in the primary problem:

Find $\alpha_1 \dots \alpha_N$ such that
 $Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ is maximized and
(1) $\sum \alpha_i y_i = 0$
(2) $\alpha_i \geq 0$ for all α_i

The Optimization Problem Solution

- ◆ The solution has the form:

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \quad b = y_k - \mathbf{w}^T \mathbf{x}_k \text{ for any } \mathbf{x}_k \text{ such that } \alpha_k \neq 0$$

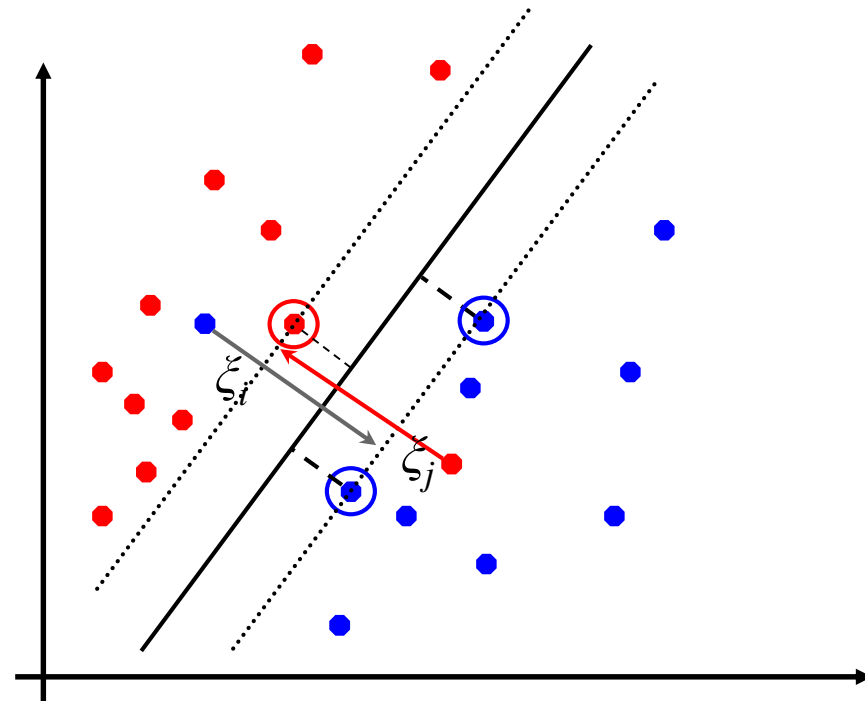
- ◆ Each non-zero α_i indicates that corresponding \mathbf{x}_i is a support vector.
- ◆ Then the classifying function will have the form:

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

- ◆ Notice that it relies on an *inner product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i
 - We will return to this later.
- ◆ Also keep in mind that solving the optimization problem involved computing the inner products $\mathbf{x}_i^T \mathbf{x}_j$ between all pairs of training points.

Soft Margin Classification

- ◆ If the training data is not linearly separable, *slack variables* ξ_i can be added to allow misclassification of difficult or noisy examples.
- ◆ Allow some errors
 - Let some points be moved to where they belong, at a cost
- ◆ Still, try to minimize training set errors, and to place hyperplane “far” from each class (large margin)



Soft Margin Classification Mathematically

- ◆ The old formulation:

Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ is minimized and for all } \{(\mathbf{x}_i, y_i)\}$$
$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

- ◆ The new formulation incorporating slack variables:

Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i \text{ is minimized and for all } \{(\mathbf{x}_i, y_i)\}$$
$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for all } i$$

- ◆ Parameter C can be viewed as a way to control overfitting
 - A regularization term

Soft Margin Classification – Solution

- ◆ The dual problem for soft margin classification:

Find $\alpha_1 \dots \alpha_N$ such that

$Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ is maximized and

(1) $\sum \alpha_i y_i = 0$

(2) $0 \leq \alpha_i \leq C$ for all α_i

- ◆ Neither slack variables ξ_i nor their Lagrange multipliers appear in the dual problem!
- ◆ Again, \mathbf{x}_i with non-zero α_i will be support vectors.
- ◆ Solution to the dual problem is:

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i$$

$$b = y_k (1 - \xi_k) - \mathbf{w}^T \mathbf{x}_k \text{ where } k = \arg\max_{k'} \alpha_{k'}$$

\mathbf{w} is not needed explicitly for classification!

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

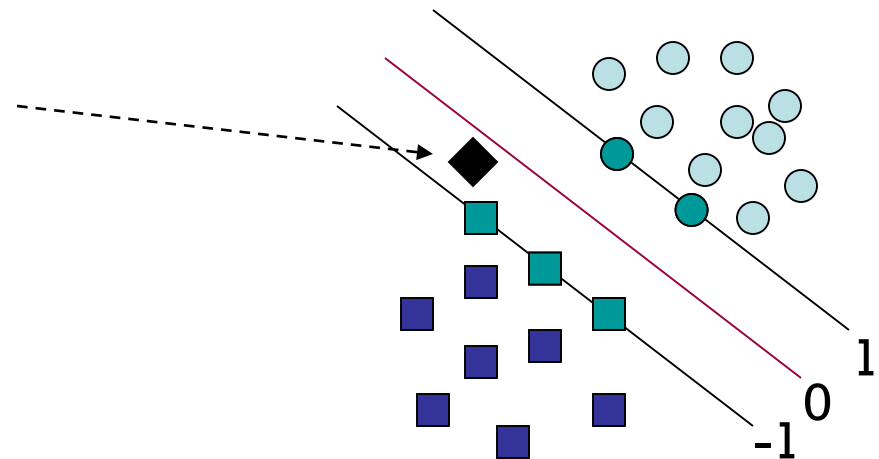
Classification with SVMs

- ◆ Given a new point \mathbf{x} , we can score its projection onto the hyperplane normal:
 - I.e., compute score: $\mathbf{w}^T \mathbf{x} + b = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$
 - Decide class based on whether $<$ or > 0
 - Can set confidence threshold t .

Score $> t$: yes

Score $< -t$: no

Else: don't know



Linear SVMs: Summary

- ◆ The classifier is a *separating hyperplane*.
- ◆ The most “important” training points are the support vectors; they define the hyperplane.
- ◆ Quadratic optimization algorithms can identify which training points \mathbf{x}_i are support vectors with non-zero Lagrangian multipliers α_i .
- ◆ Both in the dual formulation of the problem and in the solution, training points appear only inside inner products:

Find $\alpha_1 \dots \alpha_N$ such that

$Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ is maximized and

(1) $\sum \alpha_i y_i = 0$

(2) $0 \leq \alpha_i \leq C$ for all α_i

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

Linear SVMs: Example

For large data sets, the dual optimization problem can be solved using numerical techniques such as quadratic programming, a topic that is beyond the scope of this book. Once the λ_i 's are found, we can use Equations 5.39 and 5.42 to obtain the feasible solutions for w and b . The decision boundary can be expressed as follows:

$$\left(\sum_{i=1}^N \lambda_i y_i x_i \cdot x \right) + b = 0. \quad (5.44)$$

b is obtained by solving Equation 5.42 for the support vectors. Because the λ_i 's are calculated numerically and can have numerical errors, the value computed for b may not be unique. Instead it depends on the support vector used in Equation 5.42. In practice, the average value for b is chosen to be the parameter of the decision boundary.

Example 5.5. Consider the two-dimensional data set shown in Figure 5.24, which contains eight training instances. Using quadratic programming, we can solve the optimization problem stated in Equation 5.43 to obtain the Lagrange multiplier λ_i for each training instance. The Lagrange multipliers are depicted in the last column of the table. Notice that only the first two instances have non-zero Lagrange multipliers. These instances correspond to the support vectors for this data set.

Let $w = (w_1, w_2)$ and b denote the parameters of the decision boundary. Using Equation 5.39, we can solve for w_1 and w_2 in the following way:

$$w_1 = \sum_i \lambda_i y_i x_{i1} = 65.5621 \times 1 \times 0.3858 + 65.5621 \times -1 \times 0.4871 = -6.64$$

$$w_2 = \sum_i \lambda_i y_i x_{i2} = 65.5621 \times 1 \times 0.4687 + 65.5621 \times -1 \times 0.611 = -9.32$$

The bias term b can be computed using Equation 5.42 for each support vector:

$$b^{(1)} = 1 - w \cdot x_1 = 1 - (-6.64)(0.3858) - (-9.32)(0.4687) = 7.9300.$$

$$b^{(2)} = -1 - w \cdot x_2 = -1 - (-6.64)(0.4871) - (-9.32)(0.611) = 7.9289.$$

Averaging these values, we obtain $b = 7.93$. The decision boundary corresponding to these parameters is shown in Figure 5.24. ■

Linear SVMs: Example

x_1	x_2	y	Lagrange Multiplier
0.3858	0.4687	1	65.5261
0.4871	0.611	-1	65.5261
0.9218	0.4103	-1	0
0.7382	0.8936	-1	0
0.1763	0.0579	1	0
0.4057	0.3529	1	0
0.9355	0.8132	-1	0
0.2146	0.0099	1	0

Once the parameters of the decision boundary are found, a test instance z is classified as follows:

$$f(z) = \text{sign}(w \cdot z + b) = \text{sign}\left(\sum_{i=1}^N \lambda_i y_i x_i \cdot z + b\right).$$

If $f(z) = 1$, then the test instance is classified as a positive class; otherwise, it is classified as a negative class.

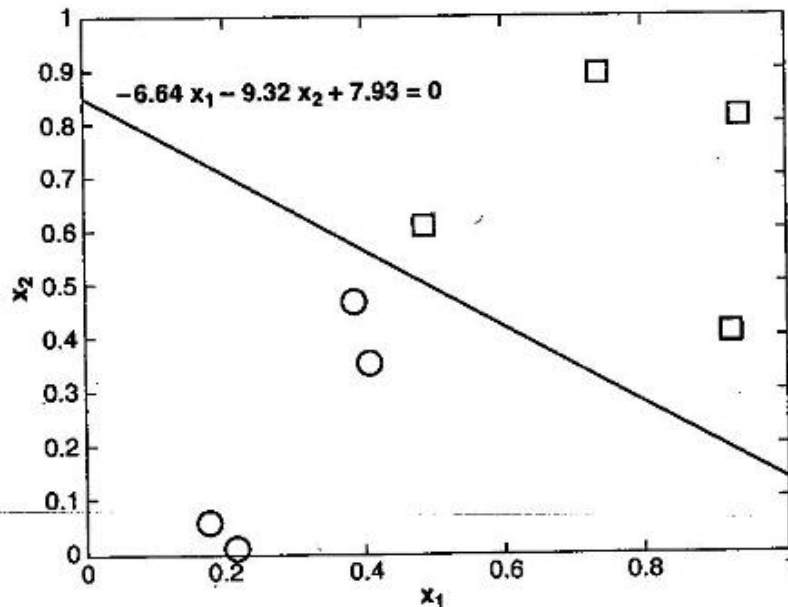


Figure 5.24. Example of a linearly separable data set.