

Linear Decision Surfaces and Functions



Binary classification problem

- **Linear decision surface** : a line, plane, or hyperplane that separates the two groups of objects
 - we use **decision functions** rather than decision surfaces to make all this mathematically more convenient
- **From data sets to decision functions**
 - an object x is a position vector in the n -dimensional dot product space \mathbf{R}^n , i.e. $\bar{x} \in \mathbf{R}^n$ (n : the number of attributes)
 - sample set $S \subset \mathbf{R}^n$
 - target function $f : \mathbf{R}^n \rightarrow \{+1, -1\}$
 - training set $D = \{(\bar{x}, y) \mid \bar{x} \in S \text{ and } y = f(\bar{x})\}$
- **Problem**
 - compute a function $\hat{f} : \mathbf{R}^n \rightarrow \{+1, -1\}$ using D such that $\hat{f}(\bar{x}) \cong f(\bar{x})$ for all $\bar{x} \in \mathbf{R}^n$
 - use linear decision surface to construct the model \hat{f}

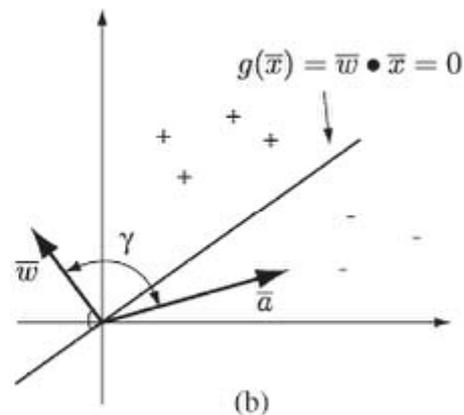
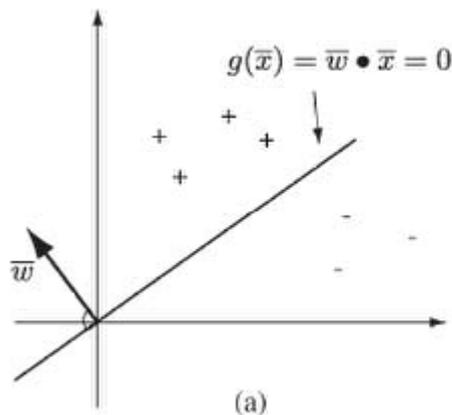
Linear decision surfaces through the origin

■ Assumption

- two-dimensional dot product space
- linearly separable training data set
- guarantee to find a line g that separates the two classes perfectly
- guarantee that the line goes through the origin of dot product space

■ Decision surface

$$\hat{f}(\bar{x}) = \begin{cases} +1 & \text{if } g(\bar{x}) \geq 0 \\ -1 & \text{if } g(\bar{x}) < 0 \end{cases}$$



$$g(\bar{a}) = \bar{w} \bullet \bar{a} = |\bar{w}| |\bar{a}| \cos(\gamma) = k$$

Linear decision surfaces with an offset

■ Relax assumption

- decision surface has not to go through the origin of dot product space

- form of decision surface is $g(\bar{x}) = \bar{w} \bullet \bar{x} = b$

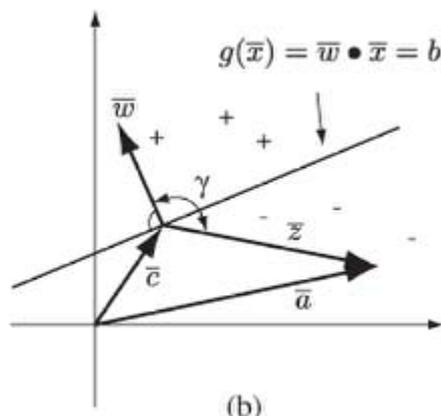
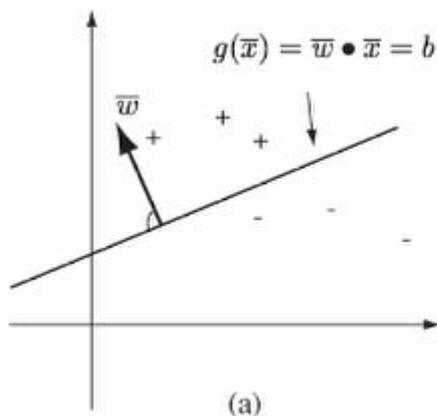
$$g(\bar{c}) = \bar{w} \bullet \bar{c} = b$$

$$\bar{a} = \bar{c} + \bar{z}$$

$$\bar{z} = \bar{a} - \bar{c}$$

■ Decision surface

- \bar{c} : some arbitrary point, say , on the decision surface itself
- if we consider \bar{z} to be the “position vector” of point \bar{a} with respect to point \bar{c} , the point \bar{c} offers almost the same perspective in p.4



$$\bar{w} \bullet \bar{z} = |\bar{w}| |\bar{z}| \cos(\gamma) = k$$

$$\begin{aligned} \bar{w} \bullet \bar{z} &= \bar{w} \bullet (\bar{a} - \bar{c}) \\ &= \bar{w} \bullet \bar{a} - \bar{w} \bullet \bar{c} \\ &= \bar{w} \bullet \bar{a} - b \\ &= g(\bar{a}) - b \end{aligned}$$

$$\hat{f}(\bar{x}) = \begin{cases} +1 & \text{if } g(\bar{x}) - b \geq 0 \\ -1 & \text{if } g(\bar{x}) - b < 0 \end{cases}$$

$$\hat{f}(\bar{x}) = \text{sgn}(\bar{w} \bullet \bar{x} - b)$$

Simple learning algorithm

■ Linearly separable training set

- $D = \{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_l, y_l)\}$
- $y_i \in \{+1, -1\}$

■ Step1

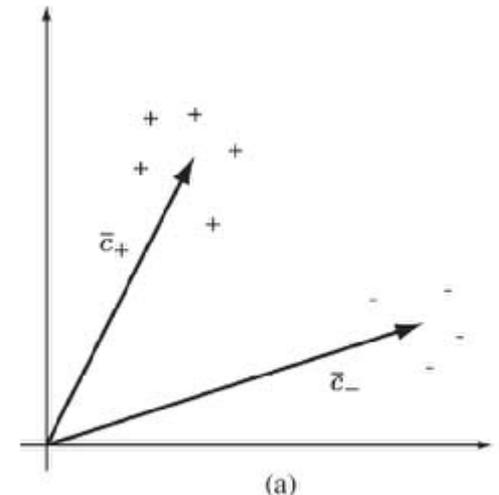
- computing the average objects or means of the two respective classes in D
- mean for the class labeled +1 : \bar{c}_+
- mean for the class labeled -1 : \bar{c}_-

$$\bar{c}_+ = \frac{1}{l_+} \sum_{(\bar{x}_i, +1) \in D} \bar{x}_i$$

$$l_+ = |\{(\bar{x}, y) \mid (\bar{x}, y) \in D \text{ and } y = +1\}|$$

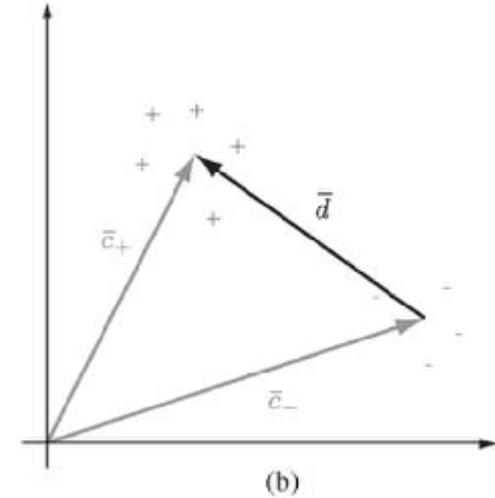
$$\bar{c}_- = \frac{1}{l_-} \sum_{(\bar{x}_i, -1) \in D} \bar{x}_i$$

$$l_- = |\{(\bar{x}, y) \mid (\bar{x}, y) \in D \text{ and } y = -1\}|$$



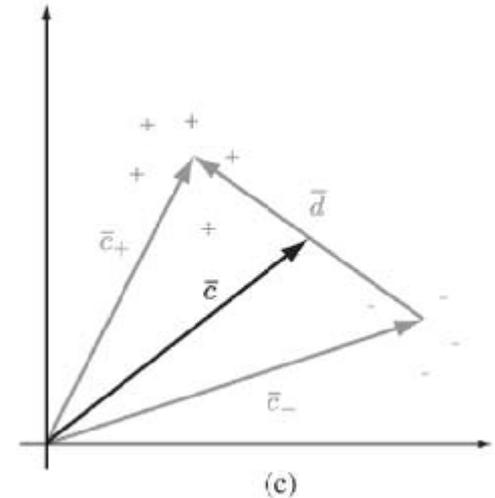
■ Step2

- construct the vector \bar{d} such that $\bar{d} = \bar{c}_+ + \bar{c}_-$



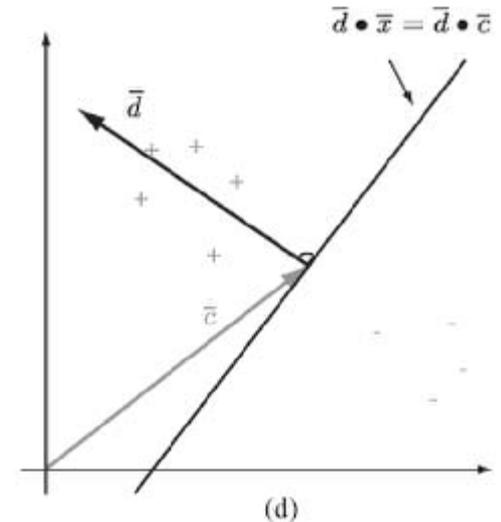
■ Step3

- compute the mean, call it \bar{c} , between the two class means \bar{c}_+ and \bar{c}_-



Step4

- translate the vector \bar{d} so that it is rooted in the average object \bar{c}
- construct a line perpendicular to \bar{d} through \bar{c}
- interpret this line as a decision surface
 - $\bar{w} = \bar{d}$
 - $b = \bar{d} \cdot \bar{c}$



Step5

- plugging equations $\bar{w} = \bar{d}$ and $b = \bar{d} \cdot \bar{c}$ into the equation for a linear decision function $\hat{f}(\bar{x}) = \text{sgn}(\bar{w} \cdot \bar{x} - b)$ gives us the decision function

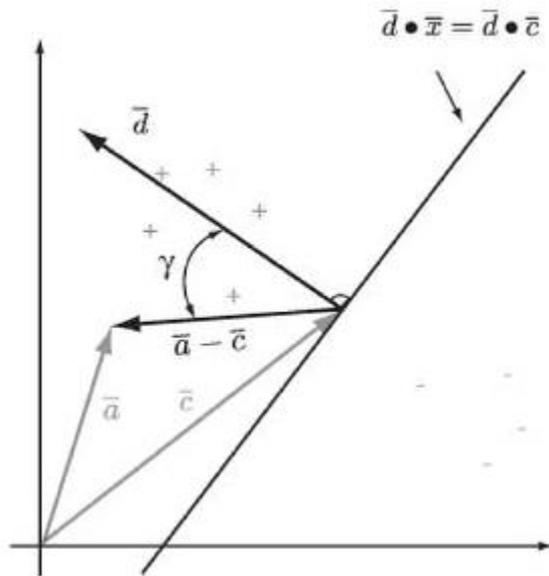
$$\hat{f}(\bar{x}) = \text{sgn}(\bar{d} \cdot \bar{x} - \bar{d} \cdot \bar{c})$$

$$\hat{f}(\bar{x}) = \text{sgn}((\bar{x} - \bar{c}) \cdot \bar{d})$$

$$= \text{sgn}(|\bar{x} - \bar{c}| |\bar{d}| \cos \gamma)$$

■ Step5 (cont.)

- decision function for some point \bar{x} is computed by taking the dot product between
 - the vector $\bar{x}-\bar{c}$, the “position vector” of the point \bar{x} with respect to \bar{c} , and
 - the normal vector of the decision surface \bar{d}
- we derive a purely algebraic expression for the decision function in terms of the class means by plugging equations $\bar{d} = \bar{c}_+ - \bar{c}_-$ and $\bar{c} = \frac{1}{2}(\bar{c}_+ + \bar{c}_-)$ into $\hat{f}(\bar{x}) = \text{sgn}((\bar{x} - \bar{c}) \bullet \bar{d})$



$$\begin{aligned} \hat{f}(\bar{x}) &= \text{sgn}((\bar{x} - \bar{c}) \bullet \bar{d}) \\ &= \text{sgn}\left(\left[\bar{x} - \frac{1}{2}(\bar{c}_+ + \bar{c}_-)\right] \bullet (\bar{c}_+ - \bar{c}_-)\right) \end{aligned}$$

Maximum-margin classifiers

■ New approach

- based on searching for a decision surface that is **equidistant to the class boundaries** where the two classes are closest to each other
- **maximizes** the distances to these class boundaries

■ Reducing the probability of misclassification

- by placing the decision surface right in the middle between the two class boundaries
- by maximizing the distances from the class boundaries

■ Constructing a maximum-margin classifier

- a **convex optimization problem** that can be solved via **quadratic programming techniques**
- training set points: the heaviest constraints on the position of such an optimal decision surface, called **support vectors**

Optimization problems

- **Selecting the best solution from a number of possible or feasible solutions**

- feasible solutions are ranked by an **objective function**
- goal is to find the feasible solution that minimizes (or maximizes) the value of this function

- **Most optimization problems have a set of constraints that limit the solution space**

- **Problem : $\min_{\bar{x}} \phi(\bar{x})$ such that $h_i(\bar{x}) \geq c_i$**

- $\phi(\bar{x})$: objective function
- $h_i(\bar{x})$: constraint with bound c_i
- optimization aim: find feasible solution \bar{x}^* s.t. $\phi(\bar{x}^*) \leq \phi(\bar{q})$ for any other feasible solution \bar{q}

■ Turning a maximization problem into a minimization problem

$$\begin{aligned}\max \phi(\bar{x}) &= \min -\phi(\bar{x}), \\ \max \phi(\bar{x}) &= \min \frac{1}{\phi(\bar{x})}\end{aligned}$$

- as long as $1/\phi(\bar{x})$ is well-defined

■ Linear optimization problem

- both an objective function and constraints are linear
- both represents lines, planes, or hyperplanes in the appropriate dot product spaces

■ Nonlinear optimization problem

- the objective function **or** the constraints are not linear

Convex optimization problem

■ A **convex objective function** and **linear constraints**

■ Particularly well behaved in that

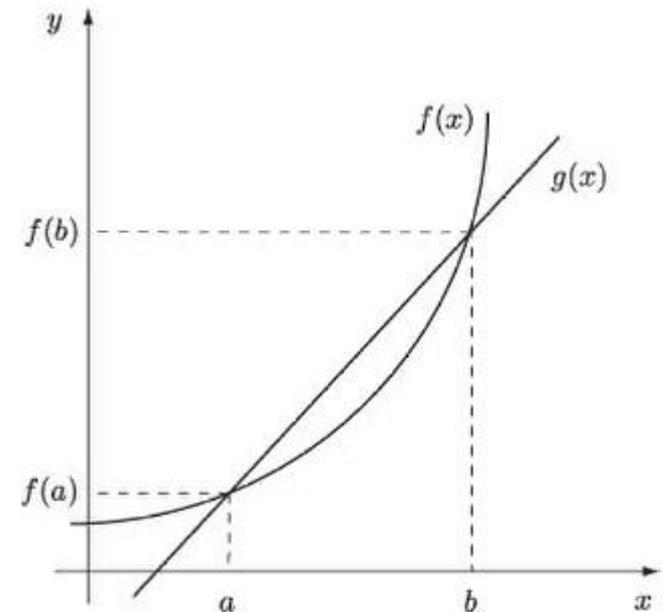
- the objective function has a **global minimum**
- the function surface is **smooth**
 - we can draw a line from one point on the function surface to any other point on the surface without crossing the surface itself

■ **Example**

- $f()$ is convex if $f(x) \leq g(x)$ in $[a, b]$

■ **Quadratic programming technique**

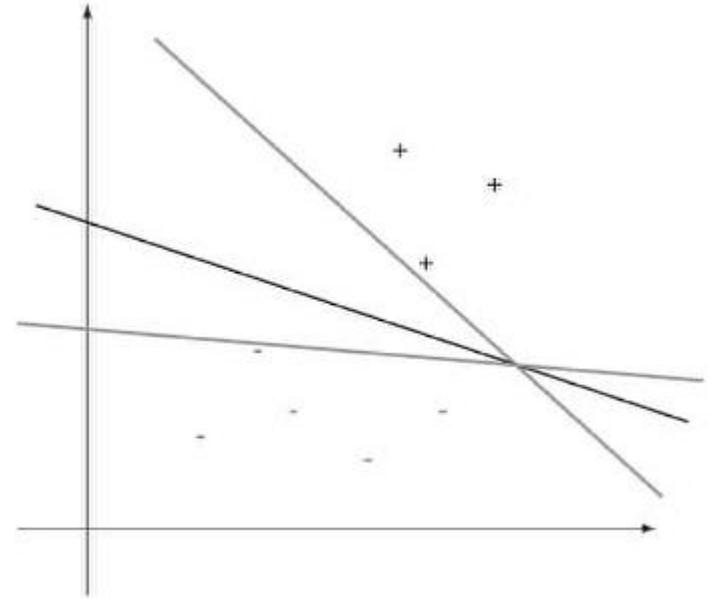
- take advantage of the convexity of an objective function in order to solve a convex optimization problem



Maximum margins

■ Intuition

- **optimal decision surface** is equidistant from the class boundaries given a linearly separable training set for a binary classification problem



■ Definition of **supporting hyperplane**

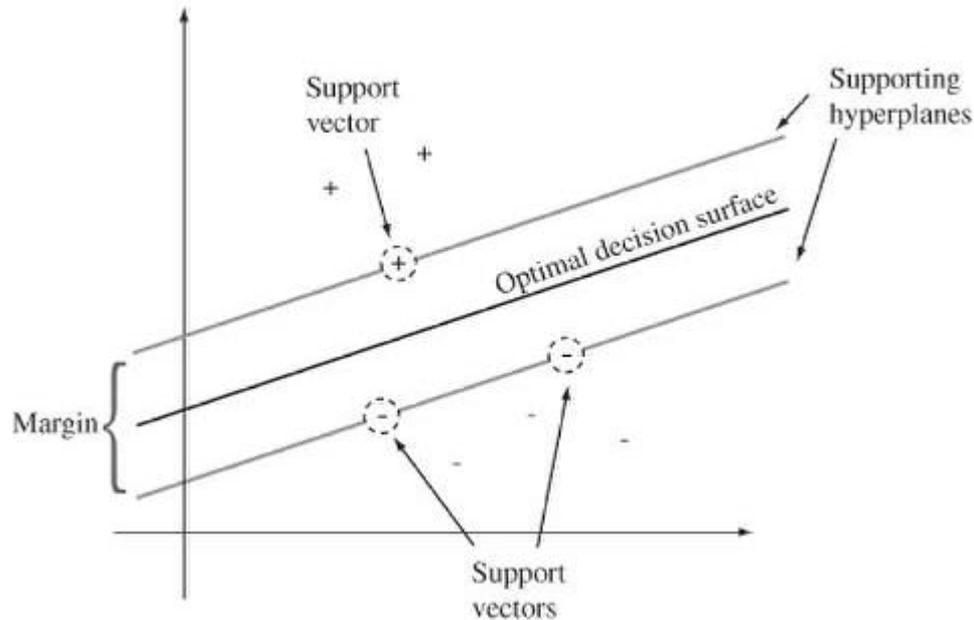
- a hyperplane **supports** a class if it is parallel to a (linear) decision surface and all points of its respective class are either above or below

■ Two supporting hyperplanes

- one is translated in the direction of the class with the + 1 label
- the other is translated in the direction of the class with the -1 label

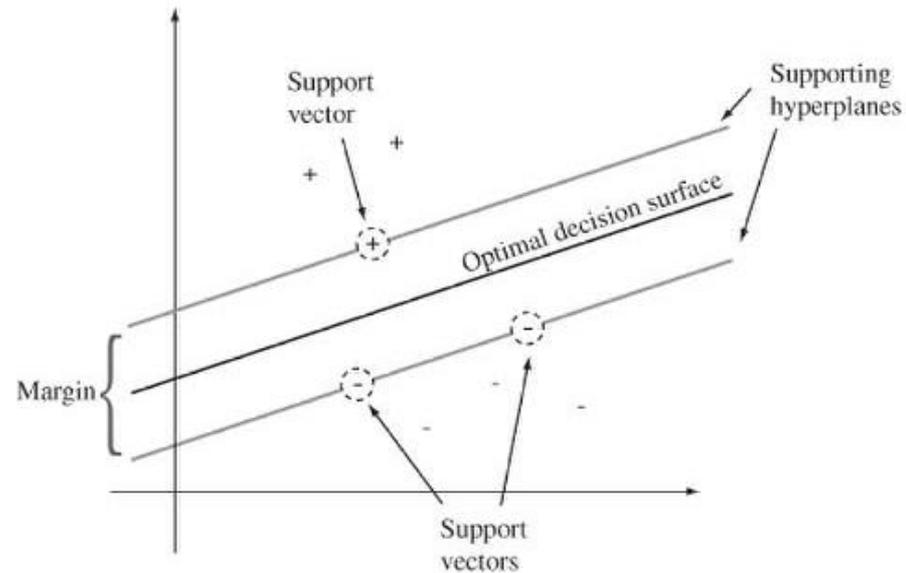
■ Definition of **margin**

- the **distance between the two supporting hyperplanes** is called a margin in a binary classification problem



■ Definition of optimality

- a decision surface for a binary classification problem is **optimal** if it is **equidistant** from the two supporting hyperplanes and **maximizes their margin**



■ Properties

- we have the two supporting hyperplanes translated so they just touch their respective class boundaries
- the distance between the hyperplanes is the margin
- the optimal decision surface is located at the center of the margin
- the size of the margin is constrained by the circled points in each class, called **support vectors**
- any rotation or translation of the decision surface would result in a smaller margin

Optimizing the margin

■ Optimization problem

- finding a decision surface that **maximizes the margin** between the two supporting hyperplanes
- **feasible solutions** : all possible decision surfaces with their associated supporting hyperplanes
- **objective function** : computes the size of the margin for each decision surface
- **constraints** : the positions of the supporting hyperplanes, which are not allowed to cross their respective class boundaries

■ Formally maximum margin $m^* = \max \phi(\bar{w}, b)$

- objective function $\phi(\bar{w}, b)$: computes the margin of a given decision surface $\bar{w} \cdot \bar{x} = b$
- maximum margin m^* : due to some optimal decision surface $\bar{w}^* \cdot \bar{x} = b^*$

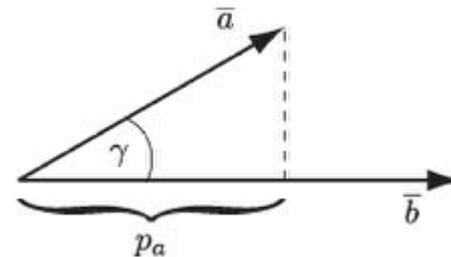
■ To make this optimization problem computable

- we need to derive a suitable expression for the objective function $\phi()$

■ Definition of projection

- let \bar{a} and \bar{b} be vectors in \mathbf{R}^n that form an angle γ between them; then we say that p_a is the projection of \bar{a} in the direction of \bar{b} s.t.

$$p_a = |\bar{a}| \cos \gamma = \frac{\bar{a} \bullet \bar{b}}{|\bar{b}|}$$



■ Our assumption

- linearly separable training set
- optimal decision surface

$$D = \{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_l, y_l)\} \subseteq \mathbb{R}^n \times \{+1, -1\}$$

$$\bar{w}^* \bullet \bar{x} = b^*$$

- The following identities hold:

$$m^* = \phi(\bar{w}^*, b^*) = \max \phi(\bar{w}, b)$$

- Maximum margin m^*

➤ can be computed by the objective function $\phi()$ given the parameters \bar{w}^* and b^* of the optimal decision surface

- Two supporting hyperplanes equidistant from the optimal decision surface

$$\bar{w}^* \bullet \bar{x} = b^*$$

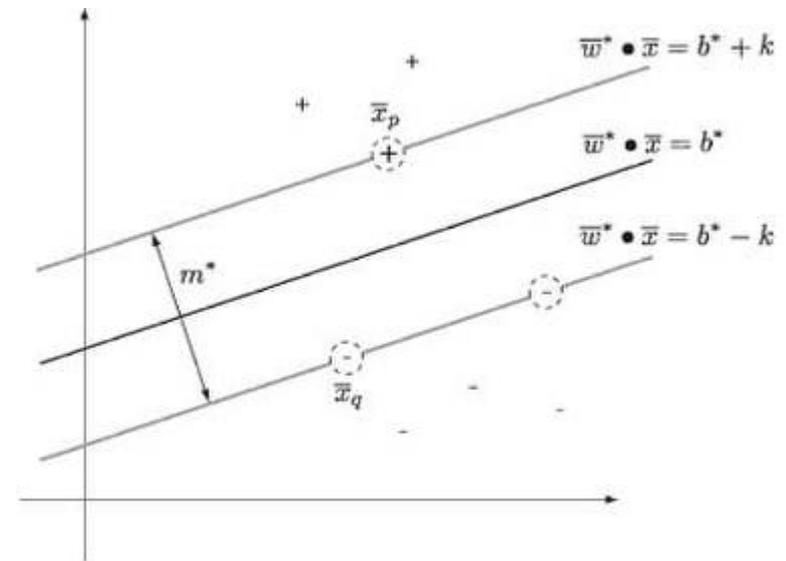
$$\bar{w}^* \bullet \bar{x} = b^* + k$$

$$\bar{w}^* \bullet \bar{x} = b^* - k$$

- Support vectors lie on the supporting hyperplane

$$\bar{w}^* \bullet \bar{x}_p = b^* + k$$

$$\bar{w}^* \bullet \bar{x}_q = b^* - k$$

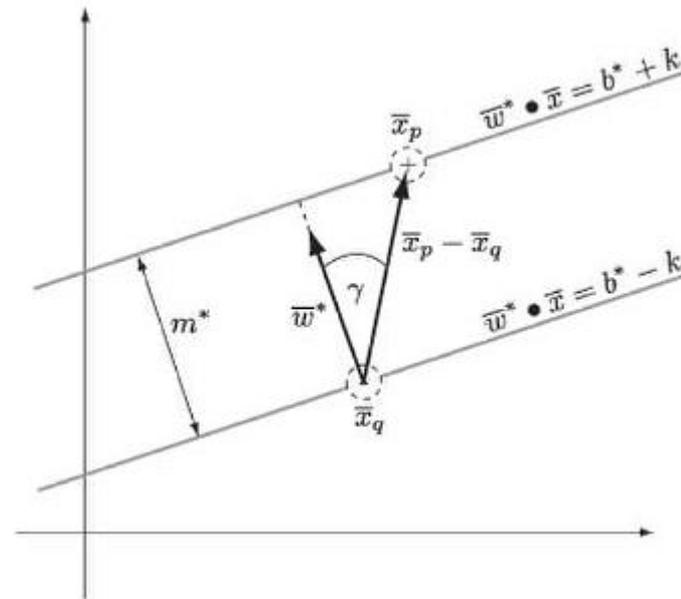


Distance between two supporting planes

- the distance is equal to the margin m^*
- it can be computed as the projection of the vector $\bar{x}_p - \bar{x}_q$ in the direction of \bar{w}^*
 - i.e., we can compute the margin as the projection of the difference between the two support vectors in the direction of the normal vector of the decision surface

$$\begin{aligned} m^* &= |\bar{x}_p - \bar{x}_q| \cos \gamma \\ &= \frac{\bar{w}^* \bullet (\bar{x}_p - \bar{x}_q)}{|\bar{w}^*|} \\ &= \frac{\bar{w}^* \bullet \bar{x}_p - \bar{w}^* \bullet \bar{x}_q}{|\bar{w}^*|} \\ &= \frac{(b^* + k) - (b^* - k)}{|\bar{w}^*|} \\ &= \frac{2k}{|\bar{w}^*|}. \end{aligned}$$

$$m^* = \max \frac{2k}{|\bar{w}|}$$



■ We want to express the maximization problem as a minimization problem

$$\begin{aligned} m^* &= \max \frac{2k}{|\bar{w}|} \\ &= \min \frac{|\bar{w}|}{2k} \\ &= \min \frac{|\bar{w}|^2}{2k} \\ &= \min \frac{1}{2k} \bar{w} \bullet \bar{w} \\ &= \min \frac{1}{2} \bar{w} \bullet \bar{w}. \end{aligned}$$

- the second step is justified since optimization over the positive values $|\bar{w}|$ is invariant under the transformation with the square function
 - square function preserves the order-theoretic properties of its domain
- the last step is justified since optimization is invariant under scaling with a constant
 - free to pick a convenient value for k, and in our case we chose $k = 1$

■ Our objective function is now $\phi(\bar{w}, b) = \frac{1}{2}\bar{w} \bullet \bar{w}$

- The function does not have a term b to be optimized
- The offset term b plays a role in the constraints

■ Shifting the focus to the constraints

- supporting hyperplanes are not allowed to cross their respective class boundaries

$$\bar{w}^* \bullet \bar{x}_i \geq b^* + k \quad \text{for all } (\bar{x}_i, y_i) \in D \text{ s.t. } y_i = +1$$

$$\bar{w}^* \bullet \bar{x}_i \leq b^* - k \quad \text{for all } (\bar{x}_i, y_i) \in D \text{ s.t. } y_i = -1$$

- generalizing and taking to our choice of $k = 1$ into account gives us

$$\bar{w} \bullet \bar{x}_i \geq 1 + b \quad \text{for all } (\bar{x}_i, y_i) \in D \text{ s.t. } y_i = +1,$$

$$\bar{w} \bullet (-\bar{x}_i) \geq 1 - b \quad \text{for all } (\bar{x}_i, y_i) \in D \text{ s.t. } y_i = -1.$$

- more compact form is $\bar{w} \bullet (y_i \bar{x}_i) \geq 1 + y_i b \quad \text{for all } (\bar{x}_i, y_i) \in D.$

■ Proposition : Maximum-Margin Classifier

- given a linearly separable training set

$$D = \{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_l, y_l)\} \subseteq \mathbb{R}^n \times \{+1, -1\},$$

- we can compute a maximum-margin decision surface with an optimization

$$\bar{w}^* \bullet \bar{x} = b^*$$

$$\min \phi(\bar{w}, b) = \min_{\bar{w}, b} \frac{1}{2} \bar{w} \bullet \bar{w}$$

- subject to the constraints

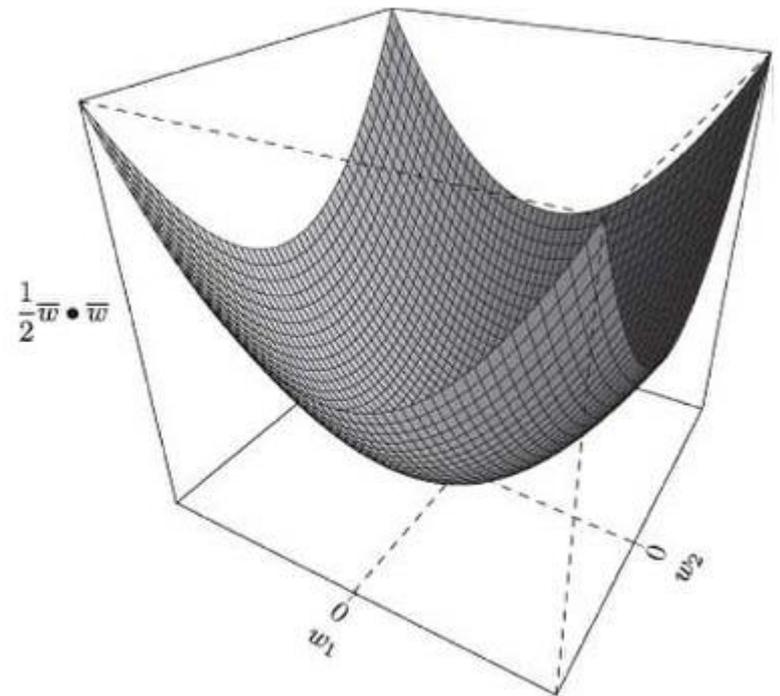
$$\bar{w} \bullet (y_i \bar{x}_i) \geq 1 + y_i b \quad \text{for all } (\bar{x}_i, y_i) \in D$$

Quadratic programming

- Our objective function is a convex function

$$\phi(\bar{w}, b) = \frac{1}{2} \bar{w} \bullet \bar{w} = \frac{1}{2} (w_1^2 + \dots + w_n^2)$$

- convexity implies that we are able to find the global minimum of our objective function
- an efficient way to solve convex optimization problems of the form given is via **quadratic programming**



$$\bar{w}^* = \text{solve}(\mathbf{Q}, \bar{q}, \mathbf{X}, \bar{c})$$

$$\bar{w}^* = \underset{\bar{w}}{\text{argmin}} \left(\frac{1}{2} \bar{w}^T \mathbf{Q} \bar{w} - \bar{q} \bullet \bar{w} \right) \text{ subject to the constraints } \mathbf{X}^T \bar{w} \geq \bar{c}$$

$$\bar{w}^* = \underset{\bar{w}}{\operatorname{argmin}} \left(\frac{1}{2} \bar{w}^T \mathbf{Q} \bar{w} - \bar{q} \bullet \bar{w} \right) \text{ subject to the constraints } \mathbf{X}^T \bar{w} \geq \bar{c}$$

■ Parameters

- $\mathbf{Q} : n \times n$ matrix
- $\mathbf{X} : l \times n$ matrix
- \bar{w}^* , \bar{w} , and $\bar{q} : n$ -dimensional vectors
- $\bar{c} : l$ -dimensional vector

■ If we let \mathbf{Q} be the identity matrix \mathbf{I} and $\bar{q} = \bar{0}$

$$\bar{w}^* = \underset{\bar{w}}{\operatorname{argmin}} \left(\frac{1}{2} \bar{w}^T \mathbf{I} \bar{w} - \bar{0} \bullet \bar{w} \right) = \underset{\bar{w}}{\operatorname{argmin}} \left(\frac{1}{2} \bar{w} \bullet \bar{w} \right) \quad \bar{w}^T \mathbf{I} \bar{w} = \bar{w} \bullet \bar{w}$$

■ Big difference between margin optimization and quadratic program solver

- quadratic program solvers return the **argument that minimizes the objective function**
 - rather than the minimized value of the objective function

■ Constraints in quadratic program solvers are expressed in matrix form

$$\mathbf{X}^T \bar{w} \geq \bar{c}$$

$$\bar{w} \bullet (y_i \bar{x}_i) \geq 1 + y_i b \quad \text{for all } (\bar{x}_i, y_i) \in D$$

$$(y_i \bar{x}_i) \bullet \bar{w} \geq 1 + y_i b$$

$$\mathbf{X} = \begin{pmatrix} y_1 x_1^1 & \cdots & y_i x_i^1 & \cdots & y_l x_l^1 \\ \vdots & & \vdots & & \vdots \\ y_1 x_1^n & \cdots & y_i x_i^n & \cdots & y_l x_l^n \end{pmatrix}$$

$$\bar{c} = \begin{pmatrix} 1 + y_1 b \\ 1 + y_2 b \\ \vdots \\ 1 + y_l b \end{pmatrix}$$

- notice : free variable b in the quadratic program
 - need to take this free variable into account in our optimization problem
- optimization has to minimize the objective function over both \bar{w} and b

■ Proposition

- given the linearly separable training set

$$D = \{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_l, y_l)\} \subseteq \mathbb{R}^n \times \{+1, -1\}.$$

- we can compute a maximum-margin decision surface with a quadratic programming approach that solves the generalized optimization problem

$$(\bar{w}^*, b^*) = \underset{\bar{w}, b}{\operatorname{argmin}} \left(\frac{1}{2} \bar{w}^T \mathbf{Q} \bar{w} - \bar{q} \bullet \bar{w} \right)$$

- subject to the constraints $\mathbf{X}^T \bar{w} \geq \bar{c}$

- with $\mathbf{Q} = \mathbf{I}$, $\bar{q} = \bar{\mathbf{0}}$, and

$$\mathbf{X} = \begin{pmatrix} y_1 x_1^1 & \cdots & y_i x_i^1 & \cdots & y_l x_l^1 \\ \vdots & & \vdots & & \vdots \\ y_1 x_1^n & \cdots & y_i x_i^n & \cdots & y_l x_l^n \end{pmatrix}$$

$$\bar{c} = \begin{pmatrix} 1 + y_1 b \\ 1 + y_2 b \\ \vdots \\ 1 + y_l b \end{pmatrix}$$

Algorithm

■ Computing a decision surface with a maximum margin using a quadratic program solver

- repeated tests in the algorithm as to whether or not the solver was successful in finding a solution
- r : radius of the training set D
- q : size of the search interval for offset term values (we set it to 1000)
 - precise value is highly data dependent
 - needs to be determined experimentally

```
let  $D = \{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_l, y_l)\} \subset \mathbb{R}^n \times \{+1, -1\}$   
 $r \leftarrow \max\{|\bar{x}| \mid (\bar{x}, y) \in D\}$   
 $q \leftarrow 1000$   
let  $\bar{w}^*$  and  $b^*$  be undefined  
Construct  $\mathbf{X}$  according to (6.32) using  $D$ .  
for each  $b \in [-q, q]$  do  
  Construct  $\bar{c}$  according to (6.33) using  $b$ .  
   $\bar{w} \leftarrow \text{solve}(\mathbf{I}, \bar{0}, \mathbf{X}, \bar{c})$   
  if ( $\bar{w}$  is defined and  $\bar{w}^*$  is undefined) or  
    ( $\bar{w}$  is defined and  $|\bar{w}| < |\bar{w}^*|$ ) then  
     $\bar{w}^* \leftarrow \bar{w}$   
     $b^* \leftarrow b$   
  end if  
end for  
if  $\bar{w}^*$  is undefined then  
  stop constraints not satisfiable  
else if  $|\bar{w}^*| > q/r$  then  
  stop bounding assumption of  $|\bar{w}|$  violated  
end if  
return  $(\bar{w}^*, b^*)$ 
```

■ To determine a reasonable interval of values for b during optimization

$$b = \bar{w} \bullet \bar{x}$$

$$b = |\bar{w}| |\bar{x}| \cos \gamma$$

$$-|\bar{w}| |\bar{x}| \leq b \leq |\bar{w}| |\bar{x}|$$

$$-|\bar{w}| r \leq b \leq |\bar{w}| r$$

$$\begin{aligned} m^* &= |\bar{x}_p - \bar{x}_q| \cos \gamma \\ &= \frac{\bar{w}^* \bullet (\bar{x}_p - \bar{x}_q)}{|\bar{w}^*|} \\ &= \frac{\bar{w}^* \bullet \bar{x}_p - \bar{w}^* \bullet \bar{x}_q}{|\bar{w}^*|} \\ &= \frac{(b^* + k) - (b^* - k)}{|\bar{w}^*|} \\ &= \frac{2k}{|\bar{w}^*|} \end{aligned}$$

➤ largest possible margin

$$\frac{2}{|\bar{w}|} \leq 2r$$

$$\frac{1}{r} \leq |\bar{w}|$$

➤ we have $|\bar{w}| \rightarrow \infty$ for infinitesimally small margins

➤ q : a bounding constant that bounds the value of $|\bar{w}|$ to a multiple of the maximum margin $1/r$

$$\frac{1}{r} \leq |\bar{w}| \leq \frac{q}{r}$$

$$-q \leq b \leq q$$

Duality

- **Support vector machines can be viewed as the **dual** to the maximum-margin classifiers**
 - the dual is obtained by applying *Lagrangian optimization theory* to the maximum-margin classifier optimization problem
- **Linear classifiers based on support vector machines can easily be extended to **nonlinear classifiers****
 - broadening the applicability of support vector machines tremendously
 - by applying what is referred to as the *kernel trick* to a linear support vector machine
 - *retaining the efficiency* of finding linear decision surfaces for training sets that are not linearly separable

Slack variables

- **Generalizing support vector machines by allowing the underlying maximum-margin classifier to **make mistakes** on the training set**

- real-world training sets are not perfect and contain noise
- noise might give rise to an **extremely complicated boundary** (decision surface) between the classes of a classification problem
- allowing the classifier to ignore these points can compute a **much simpler decision surface**
- simple decision surfaces have a **much higher probability of classifying points correctly**

- **Soft-margin classifiers**

- maximum-margin classifiers that incorporate slack variables
- c.f., hard-margin classifiers

Lagrangian dual

- Convenient technique to derive the dual of an optimization problem
- **Primal optimization problem** (*similar with p.12 in Part I*)
 - optimization problem of the form $\min_{\bar{x}} \phi(\bar{x})$ such that $g_i(\bar{x}) \geq 0$
 - $\phi(\bar{x})$ is a convex objective function
 - constraints $g_i(\bar{x})$ are linear
- **Lagrangian optimization problem** (new problem)

$$\max_{\bar{\alpha}} \min_{\bar{x}} L(\bar{\alpha}, \bar{x}) = \max_{\bar{\alpha}} \min_{\bar{x}} \left(\phi(\bar{x}) - \sum_{i=1}^l \alpha_i g_i(\bar{x}) \right) \quad \alpha_i \geq 0$$

$$\max_{\bar{\alpha}} \min_{\bar{x}} L(\bar{\alpha}, \bar{x}) = \max_{\bar{\alpha}} \min_{\bar{x}} \left(\phi(\bar{x}) - \sum_{i=1}^l \alpha_i g_i(\bar{x}) \right) \quad \alpha_i \geq 0$$

■ **New objective function $L(\bar{\alpha}, \bar{x})$ is called **Lagrangian****

- it incorporates the original objective function $\phi()$ together with a linear combination of the constraints $g_i()$
- the values $\alpha_1, \dots, \alpha_l$ are called the **Lagrangian multipliers**

$$\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_l)$$

- we have exactly one Lagrangian multiplier α_i for each constraint g_i
- we call \bar{x} the **primal variable** and $\bar{\alpha}$ the **dual variable**

■ Two nested optimization operators with opposing optimization objectives

- each optimization operator returns a partially evaluated function to be optimized by the other optimization operator
- if we fix the vector \bar{x} to the value \bar{x}^* , then the optimization problem becomes **the maximization problem**

$$\max_{\bar{\alpha}} L(\bar{\alpha}, \bar{x}^*) = \max_{\bar{\alpha}} \left(\phi(\bar{x}^*) - \sum_{i=1}^l \alpha_i g_i(\bar{x}^*) \right)$$

- if we fix the vector $\bar{\alpha}$ to the value $\bar{\alpha}^*$, then the optimization problem becomes **the minimization problem**

$$\min_{\bar{x}} L(\bar{\alpha}^*, \bar{x}) = \min_{\bar{x}} \left(\phi(\bar{x}) - \sum_{i=1}^l \alpha_i^* g_i(\bar{x}) \right)$$

■ Solutions to the Lagrangian optimization

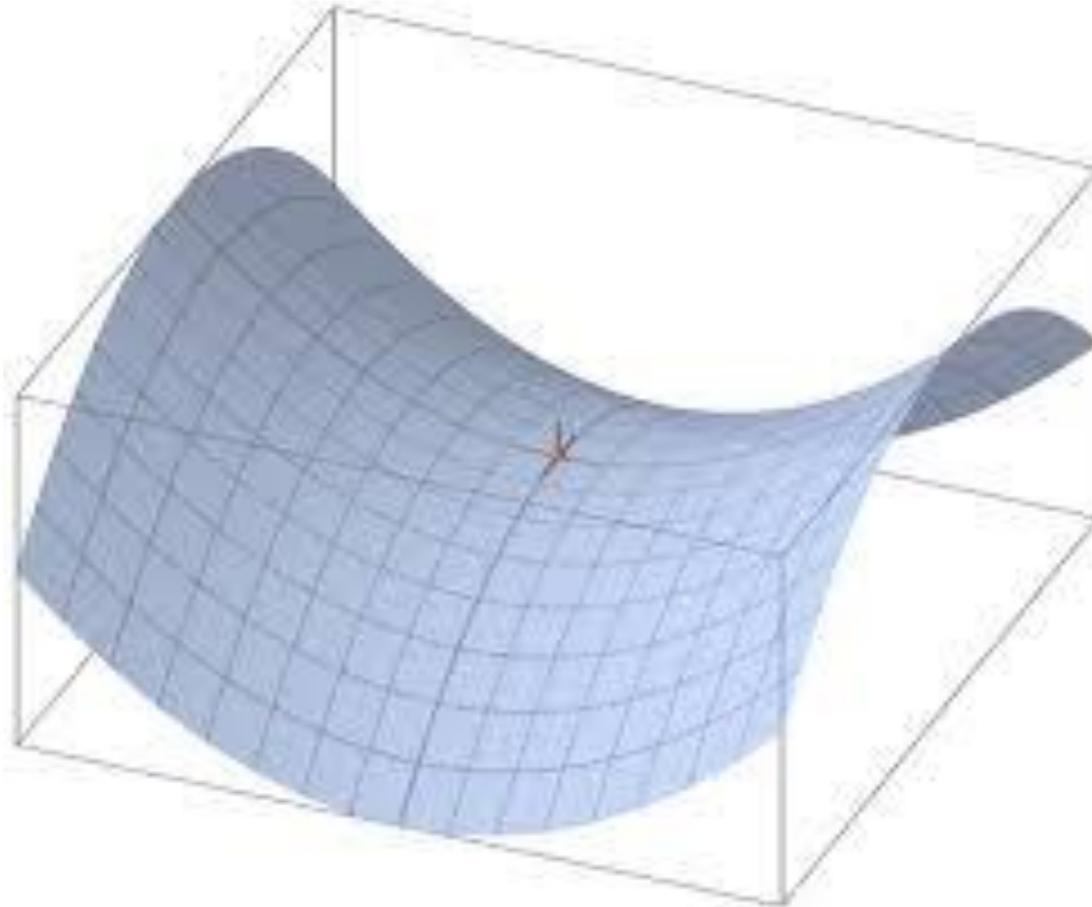
- points that both **maximize** the function $L(\bar{\alpha}, \bar{x})$ with respect to the dual variable $\bar{\alpha}$ and **minimize** it with respect to the primal variable \bar{x}
- the solutions are **saddle points** on the graph of the function $L(\bar{\alpha}, \bar{x})$
- we have a unique saddle point
 - since we assume that the primal objective function $\phi(\bar{x})$ is convex and the constraints $g_i(\bar{x})$ are linear
- the partial derivative of $L(\bar{\alpha}, \bar{x})$ with respect to \bar{x} at the saddle point has to be zero

$$\frac{\partial L}{\partial \bar{x}} = \bar{0}$$

$$\frac{\partial L}{\partial \bar{x}}(\bar{\alpha}, \bar{x}^*) = \bar{0}$$

- the point \bar{x}^* represents an optimum of $L(\bar{\alpha}, \bar{x})$ with respect to \bar{x}

Saddle point



- Let $\bar{\alpha}^*$ and \bar{x}^* be a solution to the Lagrangian such that

$$\max_{\bar{\alpha}} \min_{\bar{x}} L(\bar{\alpha}, \bar{x}) = L(\bar{\alpha}^*, \bar{x}^*) = \phi(\bar{x}^*) - \sum_{i=1}^l \alpha_i^* g_i(\bar{x}^*)$$

- Then, \bar{x}^* is a solution to the primal objective function if and only if the following conditions hold:

$$\frac{\partial L}{\partial \bar{x}}(\bar{\alpha}^*, \bar{x}^*) = \bar{0},$$

$$\alpha_i^* g_i(\bar{x}^*) = 0,$$

$$g_i(\bar{x}^*) \geq 0,$$

$$\alpha_i^* \geq 0$$

- each constraint $g_i()$ evaluated at \bar{x}^* and multiplied by its corresponding Lagrangian multiplier α_i^* has to result in a value zero
- the term $\sum_{i=1}^l \alpha_i^* g_i(\bar{x}^*)$ has to be vanish so that $L(\bar{\alpha}^*, \bar{x}^*) = \phi(\bar{x}^*)$

■ Karush-Kuhn-Tucker(KKT) conditions

- condition1: the value \bar{x}^* lies on the saddle point

$$\frac{\partial L}{\partial \bar{x}}(\bar{\alpha}^*, \bar{x}^*) = \bar{0}$$

- condition2, 3: original constraints of the primal and Lagrangian optimization problems
 - ensure that the points $\bar{\alpha}^*$ and \bar{x}^* lie in the respective feasible regions

$$g_i(\bar{x}^*) \geq 0$$
$$\alpha_i^* \geq 0$$

■ KKT complementarity condition

$$\alpha_i^* g_i(\bar{x}^*) = 0$$

■ Solving equation $\frac{\partial L}{\partial \bar{x}}(\bar{\alpha}, \bar{x}^*) = \bar{0}$ for \bar{x}^*

- reformulate the original optimization problem in terms of its dual variable only

$$L(\bar{\alpha}, \bar{x}^*) = \phi'(\bar{\alpha})$$

- find the optimum with respect to the dual variable as the Lagrangian optimization

$$\max_{\bar{\alpha}} \phi'(\bar{\alpha})$$

$$\alpha_i \geq 0$$

■ Lagrangian dual function ϕ'

- we can solve our primal optimization problem using the Lagrangian dual

$$\max_{\bar{\alpha}} \phi'(\bar{\alpha}) = \phi'(\bar{\alpha}^*) = L(\bar{\alpha}^*, \bar{x}^*) = \phi(\bar{x}^*)$$

- where $\bar{\alpha}^*$ and \bar{x}^* have to satisfy the KKT conditions

Example

- Consider the convex optimization problem $\min \phi(x) = \min \frac{1}{2}x^2$

➤ subject to the linear constraint

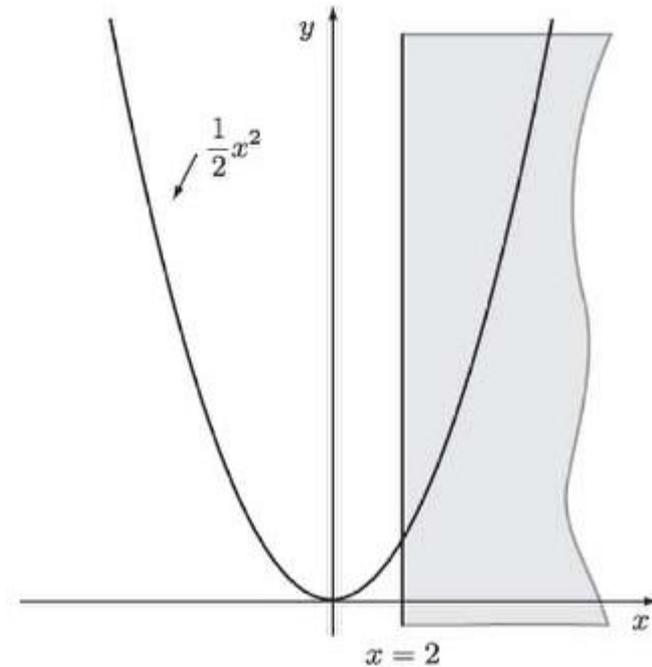
$$g(x) = x - 2 \geq 0$$

- Standard technique of finding the minimum

➤ by taking the derivative of ϕ with respect to x and setting it to zero

$$\frac{d\phi}{dx} = 0$$

fails because the value $x = 0$ that gives rise to the minimum is not part of the feasible region



Lagrangian

$$L(\alpha, x) = \frac{1}{2}x^2 - \alpha(x - 2)$$

- saddle point has to occur where the gradient of the Lagrangian with respect to the variable x is equal to zero

$$\frac{\partial L}{\partial x}(\alpha, x^*) = x^* - \alpha = 0$$

- solving for x^* gives us

$$x^* = \alpha$$

$$L(\alpha, x^*) = \frac{1}{2}\alpha^2 - \alpha^2 + 2\alpha = 2\alpha - \frac{1}{2}\alpha^2$$

- Lagrangian dual optimization with $\phi'(\alpha) = L(\alpha, x^*)$

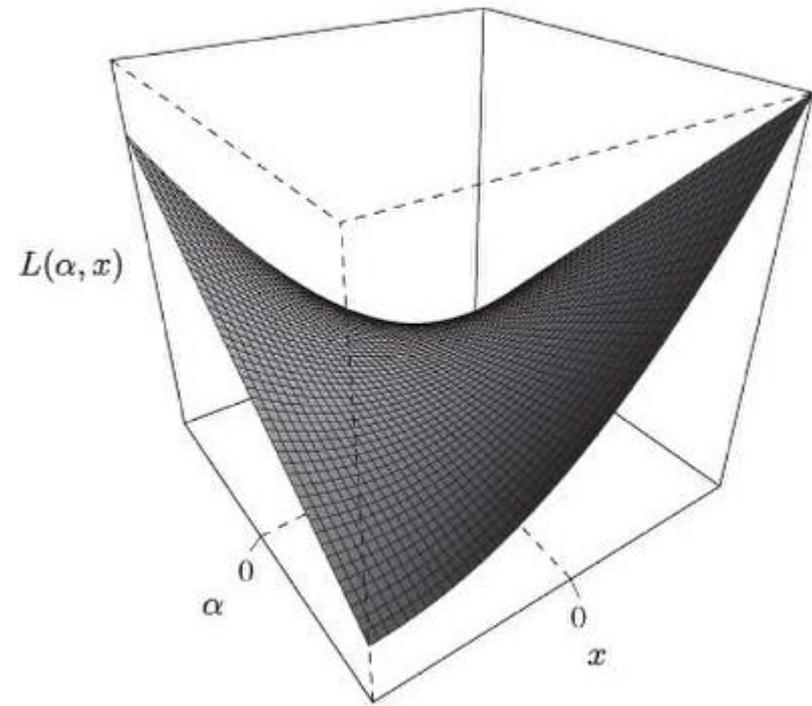
$$\max_{\alpha} \phi'(\alpha) = \max_{\alpha} \left(2\alpha - \frac{1}{2}\alpha^2 \right) \quad \alpha \geq 0$$

- compute the value α^*

$$\frac{d\phi'}{d\alpha}(\alpha^*) = 2 - \alpha^* = 0$$

$$\alpha^* = 2$$

$$\alpha^* g(x^*) = \alpha^*(x^* - 2) = 2(2 - 2) = 0$$



Dual Maximum-margin optimization

- **Support vector machines can be seen as the dual to maximum-margin classifiers**
 - derive this dual by applying the technique of the Lagrangian dual to maximum-margin classifiers
 - rewrite maximum-margin optimization problem in a form appropriate for Lagrangian optimization

$$\min_{\bar{w}, b} \phi(\bar{w}, b) = \min_{\bar{w}, b} \frac{1}{2} \bar{w} \bullet \bar{w}$$

$$g_i(\bar{w}, b) = y_i(\bar{w} \bullet \bar{x}_i - b) - 1 \geq 0$$

$$\begin{aligned} L(\bar{\alpha}, \bar{w}, b) &= \phi(\bar{w}, b) - \sum_{i=1}^l \alpha_i g_i(\bar{w}, b) \\ &= \frac{1}{2} \bar{w} \bullet \bar{w} - \sum_{i=1}^l \alpha_i (y_i(\bar{w} \bullet \bar{x}_i - b) - 1) \\ &= \frac{1}{2} \bar{w} \bullet \bar{w} - \sum_{i=1}^l \alpha_i y_i \bar{w} \bullet \bar{x}_i + b \sum_{i=1}^l \alpha_i y_i + \sum_{i=1}^l \alpha_i \end{aligned}$$

■ Lagrangian optimization problem for maximum-margin classifiers

$$\max_{\bar{\alpha}} \min_{\bar{w}, b} L(\bar{\alpha}, \bar{w}, b)$$

$$\alpha_i \geq 0$$

- let $\bar{\alpha}^*$, \bar{w}^* , and b^* be a solution to the Lagrangian optimization problem such that

$$\max_{\bar{\alpha}} \min_{\bar{w}, b} L(\bar{\alpha}, \bar{w}, b) = L(\bar{\alpha}^*, \bar{w}^*, b^*)$$

- since ϕ is convex and the constraints g_i are linear, the solution $\bar{\alpha}^*$, \bar{w}^* , and b^* satisfy the KKT conditions

$$\frac{\partial L}{\partial \bar{w}}(\bar{\alpha}^*, \bar{w}^*, b^*) = \bar{0},$$

$$\frac{\partial L}{\partial b}(\bar{\alpha}^*, \bar{w}^*, b^*) = 0,$$

$$\alpha_i^*(y_i(\bar{w}^* \bullet \bar{x}_i - b^*) - 1) = 0,$$

$$y_i(\bar{w}^* \bullet \bar{x}_i - b^*) - 1 \geq 0,$$

$$\alpha_i^* \geq 0$$

assure that \bar{w}^* and b^* lie on the saddle point of the Lagrangian

complementarity condition

- **Complementarity condition implies that \bar{w}^* and b^* are also solutions for the primal optimization problem**

$$\max_{\bar{\alpha}} \min_{\bar{w}, b} L(\bar{\alpha}, \bar{w}, b) = L(\bar{\alpha}^*, \bar{w}^*, b^*) = \phi(\bar{w}^*, b^*)$$

- **To solve the Lagrangian optimization, we construct the Lagrangian dual**

- applying the KKT condition; taking the partial derivative of Lagrangian L with respect to the primal variable \bar{w} ; evaluating it at the saddle point \bar{w}^* , and setting it to zero:

$$\frac{\partial L}{\partial \bar{w}}(\bar{\alpha}, \bar{w}^*, b) = \bar{w}^* - \sum_{i=1}^l \alpha_i y_i \bar{x}_i = \bar{0}$$

$$\frac{\partial L}{\partial \bar{w}}(\bar{\alpha}^*, \bar{w}^*, b^*) = \bar{0}$$

$$\bar{w}^* = \sum_{i=1}^l \alpha_i y_i \bar{x}_i$$

- for b

$$\frac{\partial L}{\partial b}(\bar{\alpha}, \bar{w}, b^*) = \sum_{i=1}^l \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial b}(\bar{\alpha}^*, \bar{w}^*, b^*) = 0$$

■ In a maximum-margin classifier

- the supporting hyperplane for the +1 class has to go through some point \bar{x}_p +1 closest to the class boundary
- this allows us to compute the offset b^+ of the supporting hyperplane for class +1 as

$$b^+ = \bar{w}^* \bullet \bar{x}_p$$

- for a given rotation \bar{w}^* , the point \bar{x}_p closest to the class boundary will produce the smallest offset
 - that is, we can compute b^+ as an optimization as follows:

$$b^+ = \min\{\bar{w}^* \bullet \bar{x} \mid (\bar{x}, y) \in D \text{ with } y = +1\}$$

- applying similar reasoning to the class -1

$$b^- = \max\{\bar{w}^* \bullet \bar{x} \mid (\bar{x}, y) \in D \text{ with } y = -1\}$$

- decision surface is located right between the two supporting hyperplanes

$$b^* = \frac{b^+ + b^-}{2}$$

■ Both \bar{w}^* and b^* can be expressed in terms of the dual variable by repeated use of

$$\bar{w}^* = \sum_{i=1}^l \alpha_i y_i \bar{x}_i$$

- optimal decision surface $\bar{w}^* \bullet \bar{x} = b^*$ is completely determined by the value of $\bar{\alpha}$
- finding a solution $\bar{\alpha}^*$ will give us our decision surface
 - we can find a solution $\bar{\alpha}^*$ by solving the Lagrangian dual

■ Constructing Lagrangian dual

- substituting

$$\bar{w}^* = \sum_{i=1}^l \alpha_i y_i \bar{x}_i$$

into

$$\begin{aligned} L(\bar{\alpha}, \bar{w}, b) &= \phi(\bar{w}, b) - \sum_{i=1}^l \alpha_i g_i(\bar{w}, b) \\ &= \frac{1}{2} \bar{w} \bullet \bar{w} - \sum_{i=1}^l \alpha_i (y_i (\bar{w} \bullet \bar{x}_i - b) - 1) \\ &= \frac{1}{2} \bar{w} \bullet \bar{w} - \sum_{i=1}^l \alpha_i y_i \bar{w} \bullet \bar{x}_i + b \sum_{i=1}^l \alpha_i y_i + \sum_{i=1}^l \alpha_i \end{aligned}$$

- applying the constraint

$$\frac{\partial L}{\partial b}(\bar{\alpha}, \bar{w}, b^*) = \sum_{i=1}^l \alpha_i y_i = 0.$$

$$\phi'(\bar{\alpha}) = L(\bar{\alpha}, \bar{w}^*, b^*) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \bar{x}_i \bullet \bar{x}_j$$

■ Proposition (Maximum-Margin Lagrangian Dual)

- given the maximum-margin optimization as in the previous Proposition, the Lagrangian dual optimization for maximum-margin classifiers is

$$\max_{\bar{\alpha}} \phi'(\bar{\alpha}) = \max_{\bar{\alpha}} \left(\sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \bar{x}_i \bullet \bar{x}_j \right)$$

- subject to the constraints

$$\sum_{i=1}^l \alpha_i y_i = 0.$$

$$\alpha_i \geq 0$$

■ **Case1: $\alpha_j^* > 0$ for some point $(\bar{x}_j, y_j) \in D$**

➤ to satisfy the complementarity condition we have $y_j(\bar{w}^* \bullet \bar{x}_j - b^*) - 1 = 0$

➤ or
$$\begin{aligned} \bar{w}^* \bullet \bar{x}_j &= b^* + 1 & \text{if } y_j &= +1 \\ \bar{w}^* \bullet \bar{x}_j &= b^* - 1 & \text{if } y_j &= -1 \end{aligned}$$

■ **Training set point (\bar{x}_j, y_j) with a nonzero Lagrangian multiplier $\alpha_j^* > 0$ lies on one of the two supporting hyperplanes**

- it represents a constraint on the margin in that the supporting hyperplanes cannot be moved beyond it
- we call points with nonzero Lagrangian multipliers **support vectors**
- **only support vectors contribute to the solution of the dual maximum-margin optimization from**

$$\bar{w}^* = \sum_{i=1}^l \alpha_i y_i \bar{x}_i$$

and

$$b^* = \frac{b^+ + b^-}{2}$$

■ Case2: $\alpha_j^* = 0$ for some point $(\bar{x}_j, y_j) \in \mathbf{D}$

- the point \bar{x}_j is a point that does not lie in the vicinity of the class boundary because we have $y_j(\bar{w}^* \bullet \bar{x}_j - b^*) - 1 > 0$

or

$$\bar{w}^* \bullet \bar{x}_j > b^* + 1 \quad \text{if } y_j = +1.$$

$$\bar{w}^* \bullet \bar{x}_j < b^* - 1 \quad \text{if } y_j = -1.$$

- this implies that points with zero-valued Lagrangian multipliers do not constrain the size of the margin

■ Recall

- the primal maximum-margin optimization problem finds the respective supporting hyperplanes that are farthest apart
 - i.e., that create the maximum margin between them
- the points in the training set that limit the size of the margin were called **support vectors**

■ Statement

- the primal maximum-margin optimization **computes the supporting hyperplanes** whose margin is limited by support vectors
- the dual maximum-margin optimization **computes the support vectors** that limit the size of the margin of the supporting hyperplanes

■ Insight

- only support vectors contribute to our dual solution allows us to express the value for b^* in a more elegant way
- we already know which training set points constitute the constraints on the supporting hyper-planes
 - the **points with non-zero Lagrangian multipliers**

- we can compute b^* as

$$b^* = \bar{w}^* \bullet \bar{x}_{sv^+} - 1 = \sum_{i=1}^l \alpha_i^* y_i \bar{x}_i \bullet \bar{x}_{sv^+} - 1$$

Dual decision function

■ Recall

- decision functions in linear classifiers are based on linear decision surfaces
- decision function itself returns a +1 (or -1) label for a point that lies above (or below) the decision surface

■ Optimal decision surface

$$\bar{w}^* \bullet \bar{x} = b^*$$

- maximum-margin decision function

$$\hat{f}(\bar{x}) = \text{sgn}(\bar{w}^* \bullet \bar{x} - b^*)$$

$$\hat{f}(\bar{x}) = \text{sgn}\left(\sum_{i=1}^l \alpha_i^* y_i \bar{x}_i \bullet \bar{x} - \sum_{i=1}^l \alpha_i^* y_i \bar{x}_i \bullet \bar{x}_{sv} + 1\right)$$

- **support vector machine**: dual maximum-margin classifier is determined completely by the support vectors

Linear support vector machine

■ Given

- A dot product space \mathbb{R}^n as our data universe with points $\bar{x} \in \mathbb{R}^n$
- Some target function $f : \mathbb{R}^n \rightarrow \{+1, -1\}$
- A labeled, linearly separable training set

$$D = \{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_l, y_l)\} \subseteq \mathbb{R}^n \times \{+1, -1\},$$

where $y_i = f(\bar{x}_i)$

- ## ■ Compute a model $\hat{f} : \mathbb{R}^n \rightarrow \{+1, -1\}$ using D such that $\hat{f}(\bar{x}) \cong f(\bar{x})$

$$\hat{f}(\bar{x}) = \text{sgn} \left(\sum_{i=1}^l \alpha_i^* y_i \bar{x}_i \cdot \bar{x} - \sum_{i=1}^l \alpha_i^* y_i \bar{x}_i \cdot \bar{x}_{sv^+} + 1 \right)$$

- pick one support vector from the set of available support vectors

$$(\bar{x}_{sv^+}, +1) \in \{(\bar{x}_i, +1) \mid (\bar{x}_i, +1) \in D \text{ and } \alpha_i^* > 0\}$$

- Train our support vector models with the Lagrangian dual optimization for maximum-margin classifiers

$$\bar{\alpha}^* = \operatorname{argmax}_{\bar{\alpha}} \left(\sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \bar{x}_i \bullet \bar{x}_j \right)$$

- subject to the constraints

$$\sum_{i=1}^l \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

- We can solve our classification problem using linear support vector machines as long as the training data are linearly separable

Nonlinear support vector machines

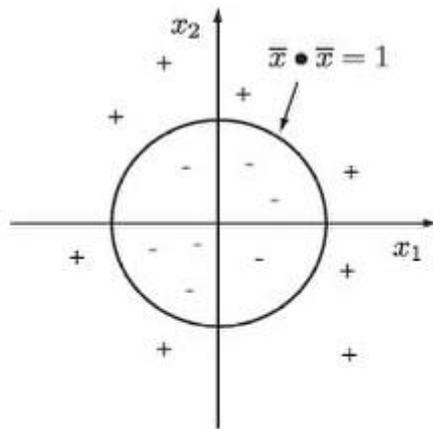
- Only very few real data sets are linearly separable
- Support vector machines
 - basic linear framework is easily extended to the case where the data set is not linearly separable
 - fundamental idea is to transform the input space into a higher-dimensional space called a feature space
 - input space: the data set is not linearly separable
 - feature space: the data are linearly separable
 - all the computations associated with the feature space can be performed in the input space
 - **kernel function**: the functions associated with transformations
 - **kernel trick**: the process of using kernel functions to move from a linear to a nonlinear support vector machine

Example

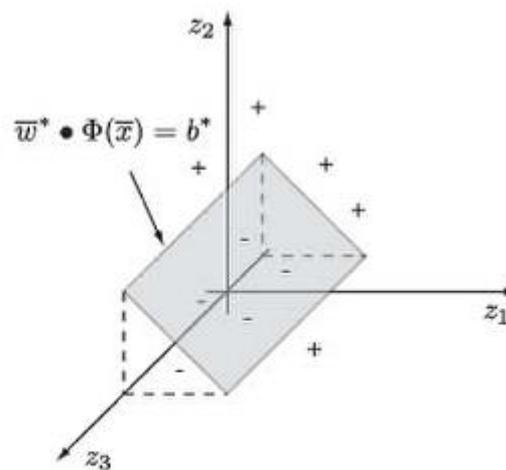
- There is no linear decision surface of the form $\bar{w} \bullet \bar{x} = b$
- In contrast, the nonlinear decision surface $\bar{x} \bullet \bar{x} = 1$
- Decision surface in a higher-dimensional space

$$\hat{f}(\bar{x}) = \text{sgn}(\bar{w} \bullet \Phi(\bar{x}) - b)$$

- Kernel function: $\Phi(\bar{x}) = \Phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) = (z_1, z_2, z_3) = \bar{z}$



(a)



(b)

$$\bar{w}^* \bullet \Phi(\bar{x}) = b^* \quad \text{plane}$$

$$\bar{w}^* = (w_1^*, w_2^*, w_3^*) = (1, 1, 0)$$

$$b^* = 1$$

■ Point $\bar{q} = (1, 0)$ in input space

- clearly lies on the nonlinear decision surface $\bar{q} \bullet \bar{q} = (1, 0) \bullet (1, 0) = 1^2 + 0^2 = 1$
- also lies on the plane in feature space:

$$\begin{aligned}\bar{w}^* \bullet \Phi(\bar{q}) &= (1, 1, 0) \bullet (1^2, 0^2, \sqrt{2} \times 1 \times 0) \\ &= (1, 1, 0) \bullet (1, 0, 0) \\ &= 1^2 + 1 \times 0 + 0 \times 0 \\ &= 1 \\ &= b^*.\end{aligned}$$

■ Decision surface $\hat{f}(\bar{x}) = \text{sgn}(\bar{w}^* \bullet \Phi(\bar{x}) - b^*)$

$$\begin{aligned}\hat{f}(\bar{x}) &= \text{sgn}(\bar{w}^* \bullet \Phi(\bar{x}) - b^*) \\ &= \text{sgn}(w_1^* x_1^2 + w_2^* x_2^2 + w_3^* \sqrt{2} x_1 x_2 - b^*) \\ &= \text{sgn}(\bar{w}^* \bullet \bar{z} - b^*) \\ &= \text{sgn}\left(\sum_{i=1}^3 w_i^* z_i - b^*\right).\end{aligned}$$

$$\Phi(\bar{x}) = \Phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2} x_1 x_2) = (z_1, z_2, z_3) = \bar{z}$$

- this shows that the complexity of the decision function is directly related to the number of dimensions in the feature space

$$\hat{f}(\bar{x}) = \text{sgn}\left(\sum_{i=1}^d w_i z_i - b\right)$$

Kernel trick

- **Dual representation of the normal vector of our decision function** $\hat{f}(\bar{x}) = \text{sgn}(\bar{w}^* \bullet \Phi(\bar{x}) - b^*)$

$$\bar{w}^* = \sum_{i=1}^l \alpha_i^* y_i \Phi(\bar{x}_i)$$

$$\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

- assume that the values α_i^* represent the appropriate Lagrangian multipliers for this dual representation
- notice that the transformation of the training points $\Phi(\bar{x}_i)$ is necessary since \bar{w}^* is a normal vector in the feature space

■ Plugging dual representation into decision function

$$\begin{aligned}\hat{f}(\bar{x}) &= \text{sgn}(\bar{w}^* \bullet \Phi(\bar{x}) - b^*) \\ &= \text{sgn}\left(\sum_{i=1}^l \alpha_i^* y_i \Phi(\bar{x}_i) \bullet \Phi(\bar{x}) - b^*\right)\end{aligned}$$

■ Simplifying it using the calculations in the table

$$\begin{aligned}\hat{f}(\bar{x}) &= \text{sgn}\left(\sum_{i=1}^l \alpha_i^* y_i \Phi(\bar{x}_i) \bullet \Phi(\bar{x}) - b^*\right) \\ &= \text{sgn}\left(\sum_{i=1}^l \alpha_i^* y_i (\bar{x}_i \bullet \bar{x})^2 - b^*\right).\end{aligned}$$

Given the mapping $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ defined as $\Phi(\bar{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$, the value of the dot product $\Phi(\bar{x}) \bullet \Phi(\bar{y})$ with $\bar{x}, \bar{y} \in \mathbb{R}^2$ can be computed in the input space \mathbb{R}^2 :

$$\begin{aligned}\Phi(\bar{x}) \bullet \Phi(\bar{y}) &= (x_1^2, x_2^2, \sqrt{2}x_1x_2) \bullet (y_1^2, y_2^2, \sqrt{2}y_1y_2) \\ &= x_1^2y_1^2 + x_2^2y_2^2 + 2x_1x_2y_1y_2 \\ &= (x_1y_1 + x_2y_2)(x_1y_1 + x_2y_2) \\ &= (\bar{x} \bullet \bar{y})(\bar{x} \bullet \bar{y}) \\ &= (\bar{x} \bullet \bar{y})^2.\end{aligned}$$

■ Obtaining an expression whose complexity is proportional to the number of support vectors

- instead of obtaining a function whose complexity is proportional to the dimensions of the feature space

■ Given an appropriate mapping $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ($m \geq n$)

- functions of the form $k(\bar{x}, \bar{y}) = \Phi(\bar{x}) \bullet \Phi(\bar{y})$ where $\bar{x}, \bar{y} \in \mathbb{R}^n$ are called **kernels** or **kernel functions**
- kernel functions evaluate a dot product in feature space
- characteristic of a kernel : the value of this dot product is actually computed in the input space

■ Rewriting decision function in terms of kernel functions

$$\begin{aligned} \hat{f}(\bar{x}) &= \text{sgn}(\bar{w}^* \bullet \Phi(\bar{x}) - b^*) \\ &= \text{sgn}\left(\sum_{i=1}^l \alpha_i^* y_i \Phi(\bar{x}_i) \bullet \Phi(\bar{x}) - b^*\right) \end{aligned} \quad \xrightarrow{\text{kernel trick}} \quad \hat{f}(\bar{x}) = \text{sgn}\left(\sum_{i=1}^n \alpha_i^* y_i k(\bar{x}_i, \bar{x}) - b^*\right)$$

■ Appropriate kernel function

- we can take advantages of mappings into feature spaces without having to pay the price of actually having to compute the explicit mappings
 - the computations in the feature space always simplify to computations in the input space
- we can control the complexity of this model by selecting the kernel function judiciously
- trick lies in finding the appropriate kernel in order to construct a model for a particular data set

■ Types of kernels

Kernel Name	Kernel Function ^a	Free Parameters
Linear kernel	$k(\bar{x}, \bar{y}) = \bar{x} \bullet \bar{y}$	none
Homogeneous polynomial kernel	$k(\bar{x}, \bar{y}) = (\bar{x} \bullet \bar{y})^d$	$d \geq 2$
Nonhomogeneous polynomial kernel	$k(\bar{x}, \bar{y}) = (\bar{x} \bullet \bar{y} + c)^d$	$d \geq 2, c > 0$
Gaussian kernel	$k(\bar{x}, \bar{y}) = e^{-\frac{ \bar{x} - \bar{y} ^2}{2\sigma^2}}$	$\sigma > 0$

■ Linear kernel

$$k(\bar{x}, \bar{y}) = \Phi(\bar{x}) \bullet \Phi(\bar{y}) = \bar{x} \bullet \bar{y}$$

- feature space is simply the same as the input space
- this kernel is useful for high-dimensional data sets in conjunction with soft-margin classifiers

■ Homogeneous polynomial kernel (of degree 2)

$$k(\bar{x}, \bar{y}) = \Phi(\bar{x}) \bullet \Phi(\bar{y}) = (\bar{x} \bullet \bar{y})^2$$

$$\Phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

■ Dual structure of the offset term b^* can also be represented with a kernel

$$\begin{aligned} b^* &= \bar{w}^* \bullet \Phi(\bar{x}_{sv+}) - 1 \\ &= \sum_{i=1}^l \alpha_i^* y_i \Phi(\bar{x}_i) \bullet \Phi(\bar{x}_{sv+}) - 1 \\ &= \sum_{i=1}^l \alpha_i^* y_i k(\bar{x}_i, \bar{x}_{sv+}) - 1, \end{aligned}$$

■ To actually find the support vectors in feature space

- we must also apply the kernel trick to our training algorithm for support vector machine models:

$$\bar{\alpha}^* = \operatorname{argmax}_{\bar{\alpha}} \left(\sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j k(\bar{x}_i, \bar{x}_j) \right)$$

- subject to constraints $\sum_{i=1}^l \alpha_i y_i = 0$, $\alpha_i \geq 0, \quad i = 1, \dots, l$

■ To find our optimal Lagrangian multipliers $\bar{\alpha}^*$

- replace the expression $\kappa(\bar{x}_i, \bar{x}_j)$ with any appropriate kernel

■ Constraints $\sum_{i=1}^n y_i \alpha_i = 0$ and $\alpha_i \geq 0$ for $i = 1, \dots, l$ are not affected by the kernel trick

Feature search

- We can rewrite our support vector machine model in terms of a kernel function

$$\hat{f}(\bar{x}) = \text{sgn} \left(\sum_{i=1}^l \alpha_i^* y_i k(\bar{x}_i, \bar{x}) - b^* \right)$$

- Similarly, we can write our training algorithm in terms of a kernel function

$$\bar{\alpha}^* = \underset{\bar{\alpha}}{\text{argmax}} \left(\sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j k(\bar{x}_i, \bar{x}_j) \right)$$

■ We are free to change kernels according to the requirements of the classification problem

- if our classification problem involves quadratic decision surfaces in input space
 - choose a polynomial kernel of degree 2 $k(\bar{x}, \bar{y}) = (\bar{x} \bullet \bar{y})^2$
 - it maps quadratic decision surfaces in the input space to linear decision surfaces in the feature space
- for more complex decision surfaces in the input space
 - try polynomial kernels of higher degrees or even more complex kernels (e.g., Gaussian kernel)

■ Feature search

- process of selecting a kernel and the associated values of its free parameters, such as the degree d for the polynomial kernel
- in general, non-trivial and requiring some trade-offs in model complexity and model accuracy
- e.g., grid search (provided by many packages)

Soft-margin classifiers

■ Generalization of maximum-margin classifiers

- to deal with noisy training data by allowing the classifiers to make mistakes

■ Recall

- maximum-margin classifiers are models of the form

$$\hat{f}(\bar{x}) = \text{sgn}(\bar{w} \bullet \bar{x} - b)$$

- where the normal vector \bar{x} and the offset term b of the decision surface are computed via the primal optimization problem

$$\min \phi(\bar{w}, b) = \min \frac{1}{2} \bar{w} \bullet \bar{w}$$

- subject to the constraints $y_i(\bar{w} \bullet \bar{x}_i - b) - 1 \geq 0$

- given the training set $(\bar{x}_1, y_1), \dots, (\bar{x}_l, y_l) \in \mathbb{R}^n \times \{+1, -1\}$

■ Construction of optimized maximum-margin classifier

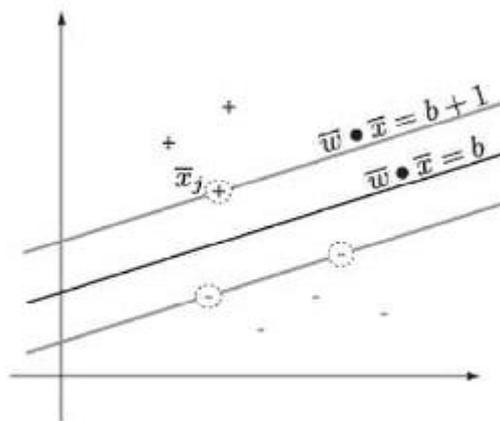
- positioning the supporting hyperplanes as far away from the decision surface as possible
 - so that they just touch their respective class boundaries
- only partially successful in the case of noisy training data where the size of the margin is limited by a few noisy training points
- reducing the impact that these points have on the size of the margin
 - by allowing them to lie on the “wrong” side of their respective supporting hyperplanes with the introduction of slack variables

■ Slack variable

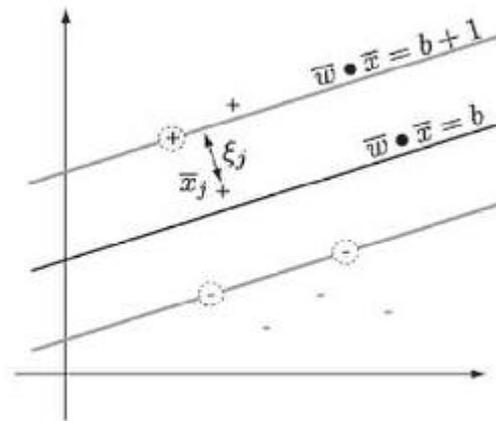
- error terms that measure how far a particular point lies on the wrong side of its respective supporting hyperplane
- measuring how much of an error is committed by allowing the supporting hyperplane to be unconstrained by that point

■ Example

- (a) maximum-margin classifier with its margin limited by the (perhaps noisy) point $(\bar{x}_j, +1)$
- (b) the supporting hyperplane $\bar{w} \cdot \bar{x} = b + 1$ is unconstrained by the training point $(\bar{x}_j, +1)$
- the training point is allowed to lie on the wrong side of the supporting hyperplane
- the amount of the error is measured by the corresponding slack variable ξ_j



(a)



(b)

- **Constraint $\bar{w} \cdot \bar{x}_j - b - 1 \geq 0$ of optimization problem is violated**

- **We can recover a sensible constraint by taking the slack variable into account $\bar{w} \cdot \bar{x}_j - b + \xi_j - 1 \geq 0$**
 - it creates the illusion that the point appears to be located right on the supporting hyperplane
 - it satisfies the original constraint
 - $\xi_j \geq 0$: error is always measured as a positive quantity

- **For each training point (\bar{x}_i, y_i)**
 - constraints are $y_i(\bar{w} \cdot \bar{x}_i - b) + \xi_i - 1 \geq 0$ with $\xi_j \geq 0$

■ Proposition (Soft-Margin Optimization)

- given a training set

$$D = \{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_l, y_l)\} \subseteq \mathbb{R}^n \times \{+1, -1\}$$

- we can compute a soft-margin decision surface

$$\bar{w}^* \bullet \bar{x} = b^*$$

- with an optimization

$$\min_{\bar{w}, \bar{\xi}, b} \phi(\bar{w}, \bar{\xi}, b) = \min_{\bar{w}, \bar{\xi}, b} \left(\frac{1}{2} \bar{w} \bullet \bar{w} + C \sum_{i=1}^l \xi_i \right)$$

- subject to the constraints

$$y_i (\bar{w} \bullet \bar{x}_i - b) + \xi_i - 1 \geq 0$$

$$\xi_i \geq 0,$$

$$\text{with } i = 1, \dots, l, \bar{\xi} = (\xi_1, \dots, \xi_l), \text{ and } C > 0$$

■ Important fact

- **still a convex optimization** because all the error terms are constrained to be positive values

■ Optimizing the function $\phi(\bar{w}, \bar{\xi}, b)$

- it is a trade-off between the size of the margin and the size of the error
- the error is the sum of the values of the slack variables
- the larger we make the margin, the more training points will be on the wrong side of their respective supporting hyperplanes
 - therefore, the larger the error
- if we make the margin large, this will probably introduce a large number of nonzero slack variables
- If we make the margin small, we can reduce the number of nonzero slack variables
 - but, we are also back to where we started in the sense that noisy points will dictate the position of the decision surface

■ Controlling the trade-off between margin size and error

- constant C (called the cost) allows us to do
- we cannot simply ignore the slack variables by setting $C = 0$
- with a large value for C , the optimization will try to find a solution with as small a number of nonzero slack variables as possible
 - because errors are costly, due to the large C
- with a small value for C , the introduction of nonzero slack variables is much more forgiving and we can find solutions with a larger margin
 - ignoring some of the noisier points near the decision surface

large $C \sim$ small margin,
small $C \sim$ large margin.

- a solution \bar{w}^* , $\bar{\xi}^*$, and b^* to the optimization problem

$$\min_{\bar{w}, \bar{\xi}, b} \phi(\bar{w}, \bar{\xi}, b) = \frac{1}{2} \bar{w}^* \bullet \bar{w}^* + C \sum_{i=1}^l \xi_i^* = m^*$$

is then a trade-off between the size of the margin m^* and the size of the error $\sum_{i=1}^l \xi_i^*$.

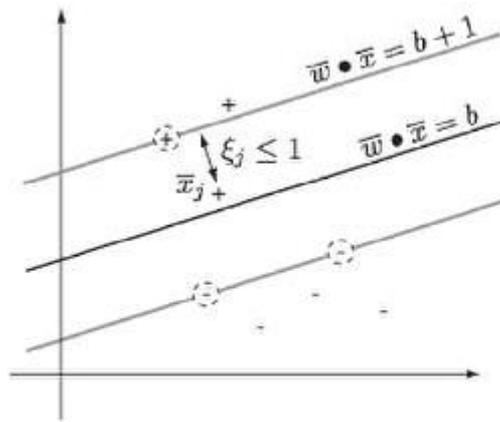
■ Slack variables only appear as part of the training algorithm

- maximum-margin classifier model itself remains unchanged

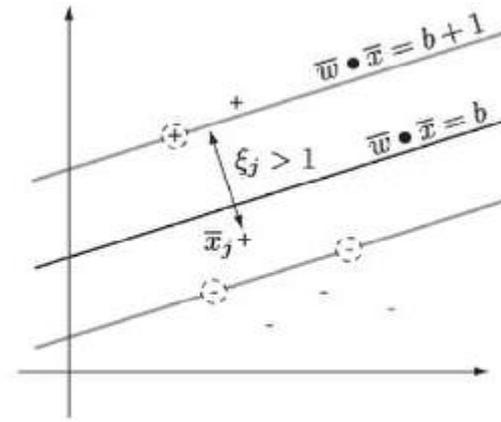
$$\hat{f}(\bar{x}) = \text{sgn}(\bar{w}^* \bullet \bar{x} - b^*)$$

- only difference between the model in a hard-margin setting and in a soft-margin setting : allowing our model to make a certain number of classification errors governed by the cost constant C

- If a particular point (\bar{x}_j, y_j) has a slack variable with value less than 1, $\xi_j \leq 1$
 - the point will be classified correctly by the decision function even though the point lies in the margin
- If the point has a slack variable such that $\xi_j > 1$
 - the point will be misclassified by the decision function



(a)



(b)

We skip

- **Dual Setting for Soft-Margin Classifiers**
- **Implementation**
 - gradient ascent
 - quadratic programming
 - **sequential minimal optimization (SMO)**

