

# Support Vector Machines - Dual formulation

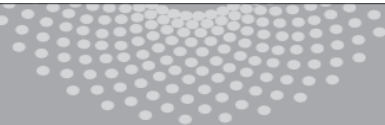
Aarti Singh

Machine Learning 10-701

Feb 6, 2023



**MACHINE LEARNING** DEPARTMENT



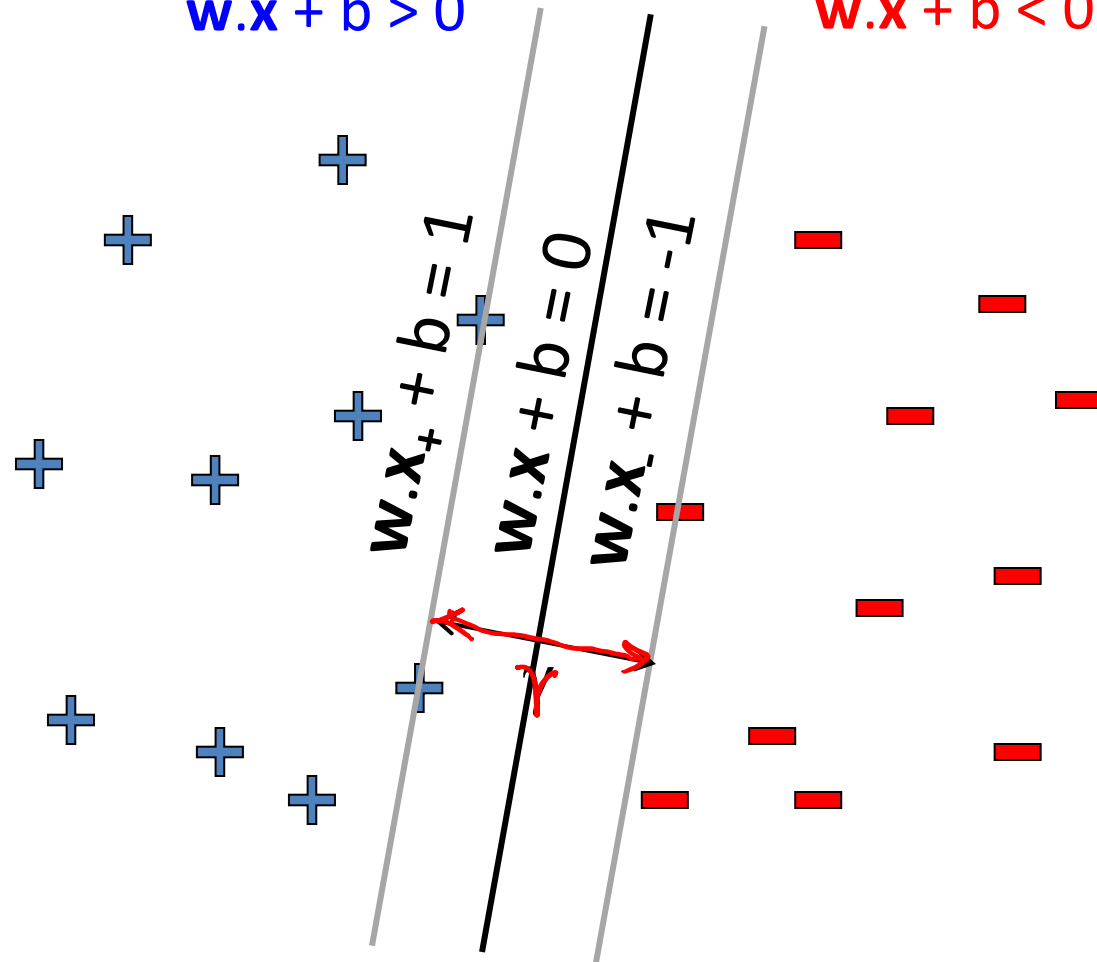
**Carnegie Mellon.**  
School of Computer Science

# Hard margin SVM

$$\gamma \propto \frac{1}{\|w\|}$$

$$w \cdot x + b > 0$$

$$w \cdot x + b < 0$$



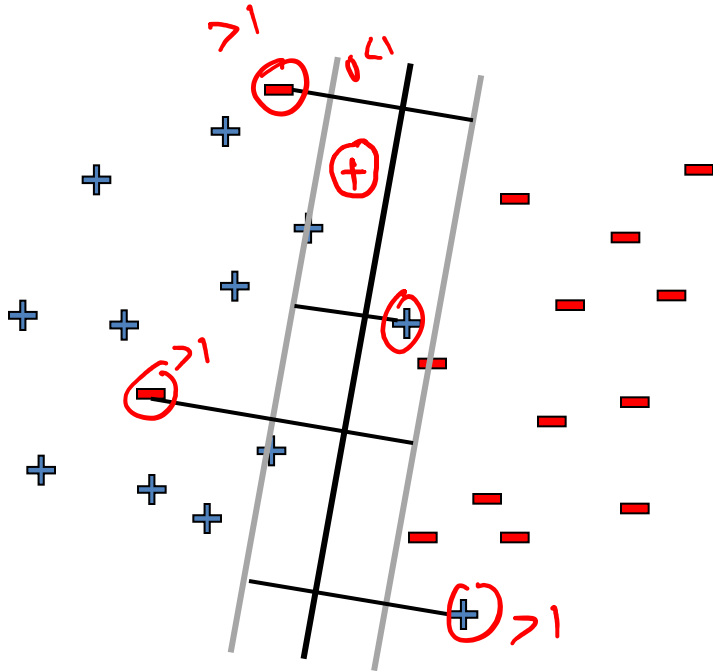
$$\min_{w,b} w \cdot w = \|w\|^2$$
$$\text{s.t. } (w \cdot x_j + b) y_j \geq 1 \quad \forall j$$

Solve efficiently by quadratic programming (QP)

Assumes data is linearly separable

# Soft margin SVM

Allow “error” in classification



$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_j\}} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ \text{s.t.} \quad & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j \\ & \xi_j \geq 0 \quad \forall j \end{aligned}$$

$\xi_j$  - “slack” variables  
= (>1 if  $x_j$  misclassified)

pay linear penalty if mistake

$C$  - tradeoff parameter ( $C = \infty$   
recovers hard margin SVM)

Still QP 😊

# SVM – linearly separable case

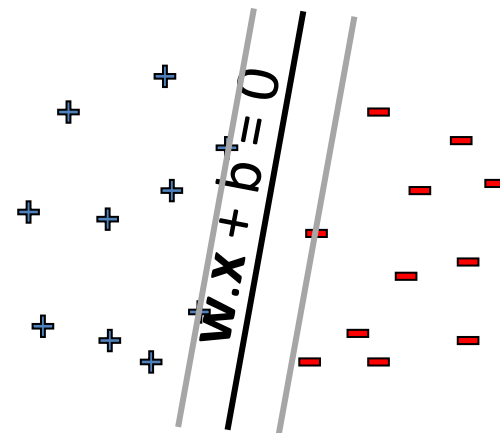
n training points

$(\mathbf{x}_1, \dots, \mathbf{x}_n)$

d features

$\mathbf{x}_j$  is a d-dimensional vector

- Primal problem: minimize  $\frac{1}{2} \mathbf{w} \cdot \mathbf{w}$  }  
 $\rightarrow (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \forall j$  }



**w - weights on features (d-dim problem)**

- Convex quadratic program – quadratic objective, linear constraints
- But expensive to solve if d is very large
- Often solved in dual form (n-dim problem)

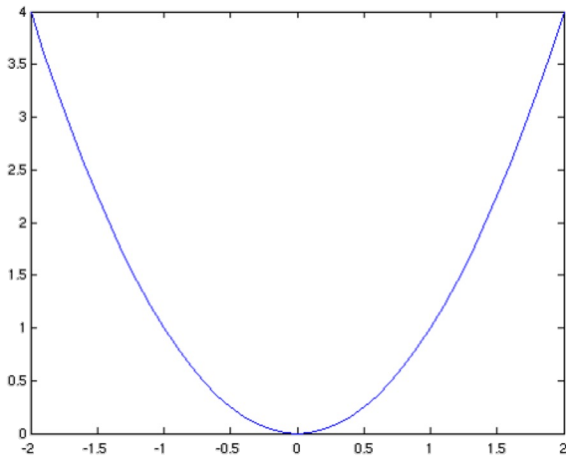
# Detour - Constrained Optimization

$x \in (w, b)$

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq b \end{aligned}$$

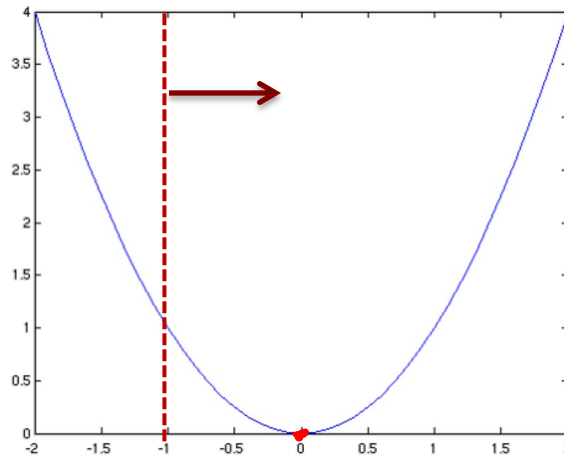
$$x^* = \max(b, 0)$$

$$\min_x x^2$$



$$x^* = 0$$

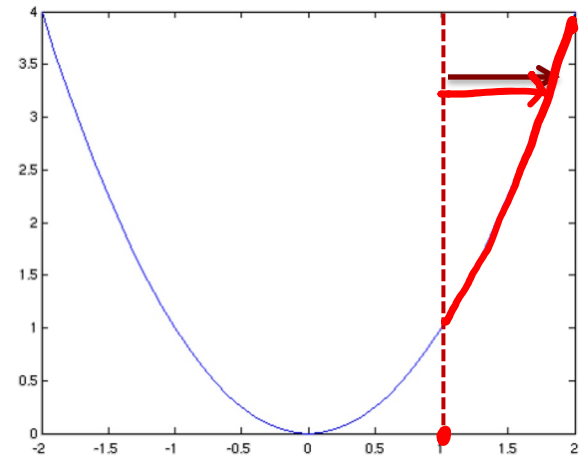
$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq -1 \end{aligned}$$



$$x^* = 0$$

Constraint inactive

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq 1 \end{aligned}$$

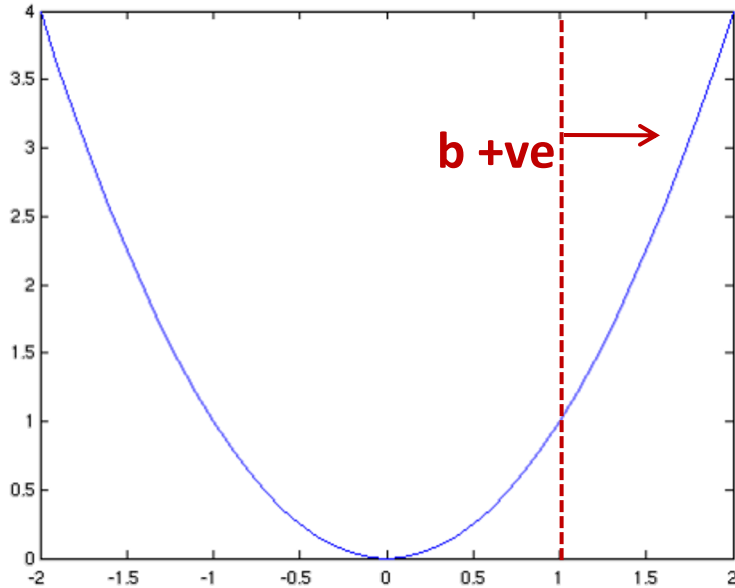


$$x^* = 1$$

Constraint active

(tight)

# Constrained Optimization



$$x^* = b$$

$$I(x) = I(x-b) = \begin{cases} \infty & x < b \\ 0 & x \geq b \end{cases}$$

$$x^2 - \alpha(x-b)$$



$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & \underline{x \geq b} \end{aligned}$$

Equivalent unconstrained optimization:

$$\min_x \quad \underline{x^2 + I(x-b)}$$

Replace with lower bound ( $\alpha \geq 0$ )

$$\underline{x^2 + I(x-b)} \geq \underline{x^2 - \alpha(x-b)} \quad \text{Lagrangian}$$

# Primal and Dual Problems

Notice that

$$\text{Primal problem: } p^* = \min_x x^2 \quad = \quad \min_x \max_{\alpha \geq 0} L(x, \alpha)$$

s.t.  $x \geq b$ 
 $\alpha \geq 0$

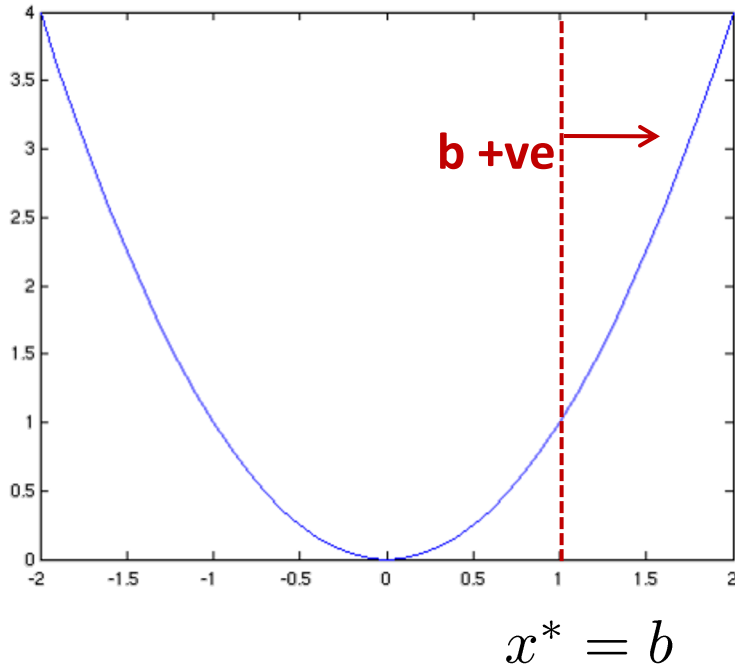
Why?  $L(x, \alpha) = x^2 - \alpha(x - b)$  ✓

$$\max_{\alpha \geq 0} L(x, \alpha) = x^2 - \min_{\alpha \geq 0} \alpha(x - b) = \begin{cases} +\infty & x < b \\ x^2 & x \geq b \end{cases}$$

$$\text{Dual problem: } d^* = \max_{\alpha} d(\alpha) \quad = \quad \max_{\alpha} \min_x L(x, \alpha) = d(\alpha)$$

s.t.  $\alpha \geq 0$ 
 $\alpha \geq 0$

# Recipe for deriving Dual Problem



Primal problem:

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & \underline{x \geq b} \end{aligned} \quad \begin{array}{l} g(x) \geq 0 \\ \underline{x - b \geq 0} \end{array}$$

Moving the constraint to objective function  
Lagrangian:

$$\begin{aligned} L(x, \alpha) &= x^2 - \alpha(x - b) \\ \text{s.t.} \quad & \underline{\alpha \geq 0} \end{aligned}$$

Dual problem:

$$\begin{aligned} \max_{\alpha} \quad & d(\alpha) \longrightarrow \min_x L(x, \alpha) \\ \text{s.t.} \quad & \alpha \geq 0 \end{aligned}$$



# Why solve the Dual?

$$\underline{d(\lambda)} = \min_x L(x, \alpha)$$

**Primal problem:**  $p^* = \min_x \overset{f(x)}{x^2}$   
s.t.  $x \geq b$

**Dual problem:**  $d^* = \max_{\alpha} \underline{d(\alpha)}$   
s.t.  $\alpha \geq 0$

$$= \min_x \max_{\alpha \geq 0} L(x, \alpha)$$

$$= \max_{\alpha} \min_x L(x, \alpha)$$

s.t.  $\alpha \geq 0$

- **Dual problem (maximization) is always concave even if primal is not convex**

Why? Pointwise infimum of concave functions is concave.  
[Pointwise supremum of convex functions is convex.]

$$L(x, \alpha) = x^2 - \alpha(x - b) \quad \underline{\quad \quad \quad} \quad \neq \quad \neq$$

- **As many dual variables  $\alpha$  as constraints, helpful if fewer constraints than dimension of primal variable  $x$**

# Connection between Primal and Dual

**Primal problem:**  $p^* = \min_x x^2$   
s.t.  $x \geq b$

**Dual problem:**  $d^* = \max_{\alpha} d(\alpha)$   
s.t.  $\alpha \geq 0$

➤ **Weak duality:** The dual solution  $d^*$  lower bounds the primal solution  $p^*$  i.e.  $d^* \leq p^*$

To see this, recall  $L(x, \alpha) = x^2 - \alpha(x - b)$

For every feasible  $x'$  (i.e.  $x' \geq b$ ) and feasible  $\alpha'$  (i.e.  $\alpha' \geq 0$ ), notice that

$$d(\alpha) = \min_x L(x, \alpha) \leq x'^2 - \alpha'(x' - b) \leq x'^2$$

Since above holds true for every feasible  $x'$ , we have  $d(\alpha) \leq x^{*2} = p^*$

# Connection between Primal and Dual

$$x \equiv (w, b)$$

$$\text{Primal problem: } p^* = \min_x x^2 \quad \checkmark$$

$x^*$       =      s.t.  $x \geq b$

$g(x) \geq 0$

$$\text{Dual problem: } d^* = \max_{\alpha} d(\alpha)$$

$d^*$       =      s.t.  $\alpha \geq 0$

- **Weak duality:** The dual solution  $d^*$  lower bounds the primal solution  $p^*$  i.e.  $d^* \leq p^*$

$$\text{Duality gap} = p^* - d^*$$

- **Strong duality:**  $d^* = p^*$  holds often for many problems of interest e.g. if the primal is a feasible convex objective with linear constraints (Slater's condition)

# Connection between Primal and Dual

What does strong duality say about  $\alpha^*$  (the  $\alpha$  that achieved optimal value of dual) and  $x^*$  (the  $x$  that achieves optimal value of primal problem)?

## KKT (Karush-Kuhn-Tucker conditions)

Whenever strong duality holds, the following conditions (known as KKT conditions) are true for  $\alpha^*$  and  $x^*$ :

- 1.  $\nabla L(\underline{x^*}, \underline{\alpha^*}) = \underline{0}$  i.e. Gradient of Lagrangian at  $\underline{x^*}$  and  $\underline{\alpha^*}$  is zero. ✓
- 2.  $\underline{x^*} \geq b$  i.e.  $x^*$  is primal feasible ✓
- 3.  $\alpha^* \geq 0$  i.e.  $\alpha^*$  is dual feasible ✓✓
- 4.  $\alpha^*(x^* - b) = 0$  (called as complementary slackness)

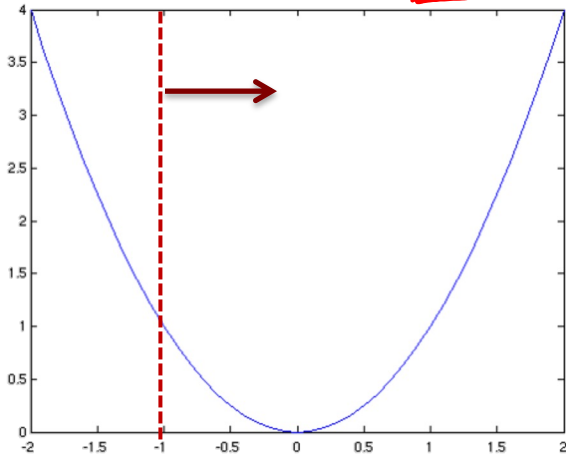
$$\underbrace{b = -1} \quad \underbrace{\alpha^*(x^* + 1) = 0} \quad \Rightarrow \quad \text{either } \underbrace{\alpha^* = 0} \quad \text{or} \quad \underbrace{x^* = -1}$$

$$\mathcal{L}^*(x^*)$$

# Constrained Optimization – Dual Problem

$$\mathcal{L}^*(x^*-1) = 0$$

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq -1 \end{aligned}$$



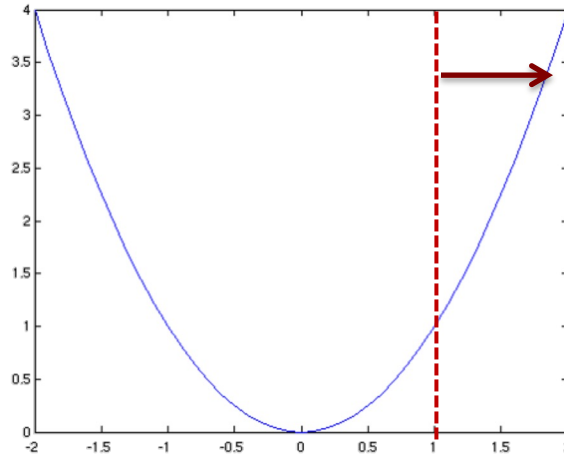
$$x^* = 0$$

$$\alpha^* = 0$$

constraint is inactive

$$x^* > -1$$

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq 1 \end{aligned}$$



$$x^* = 1$$

$$\alpha^* \geq 0$$

constraint is active

$$x^* = 1$$

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq -1 \\ & x \geq 1 \end{aligned}$$

$$g(x) \geq 0$$

$$\alpha_1(x+1)$$

$$\alpha_2(x-1)$$

► Poll:

Lagrangian

Dual variables

$$x^2 - \alpha_1(x+1) - \alpha_2(x-1)$$

$$x^2 - \alpha g(x)$$

# Dual SVM – linearly separable case

n training points, d features  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  where  $\mathbf{x}_i$  is a d-dimensional vector

- Primal problem: minimize<sub>w,b</sub>  $\frac{1}{2} \mathbf{w} \cdot \mathbf{w}$   
 $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \forall j$   
 $\alpha_j \geq 0$

$\mathbf{w}$  - weights on features (d-dim problem)

- Dual problem (derivation):

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j [(\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1]$$

$\alpha_j \geq 0, \forall j$

$\alpha$  - weights on training pts (n-dim problem)

# Dual SVM – linearly separable case

$$\alpha = (\alpha_1, \dots, \alpha_n)$$

- Dual problem (derivation):

$$\max_{\alpha} \min_{\mathbf{w}, b} \overset{d(\alpha)}{L(\mathbf{w}, b, \alpha)} = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j \left[ (\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1 \right]$$

$\alpha_j \geq 0, \forall j$

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_j \alpha_j \mathbf{x}_j y_j = 0$$

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

$$\frac{\partial L}{\partial b} = 0 \quad \Rightarrow \quad \sum_j \alpha_j y_j = 0$$

$$\frac{\partial L}{\partial b} = -\sum_j \alpha_j y_j = 0$$

# Dual SVM – linearly separable case

- Dual problem:

$$\max_{\alpha} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j [(\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1]$$

$$\alpha_j \geq 0, \forall j$$

$$\Rightarrow \mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

$$\Rightarrow \sum_j \alpha_j y_j = 0$$

$$\min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha) = \frac{1}{2} \sum_j \alpha_j y_j \mathbf{x}_j \cdot \sum_i \alpha_i y_i \mathbf{x}_i - \sum_j \alpha_j [(\sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}_j + b) y_j - 1]$$

$$= \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j - \sum_j \sum_i \alpha_j \alpha_i y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j + \sum_j \alpha_j$$

$$= -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j + \sum_j \alpha_j = d(\alpha)$$



# Dual SVM – linearly separable case

$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

$\text{sign}(w \cdot x + b)$

Dual problem is also QP

Solution gives  $\alpha_j$ s



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

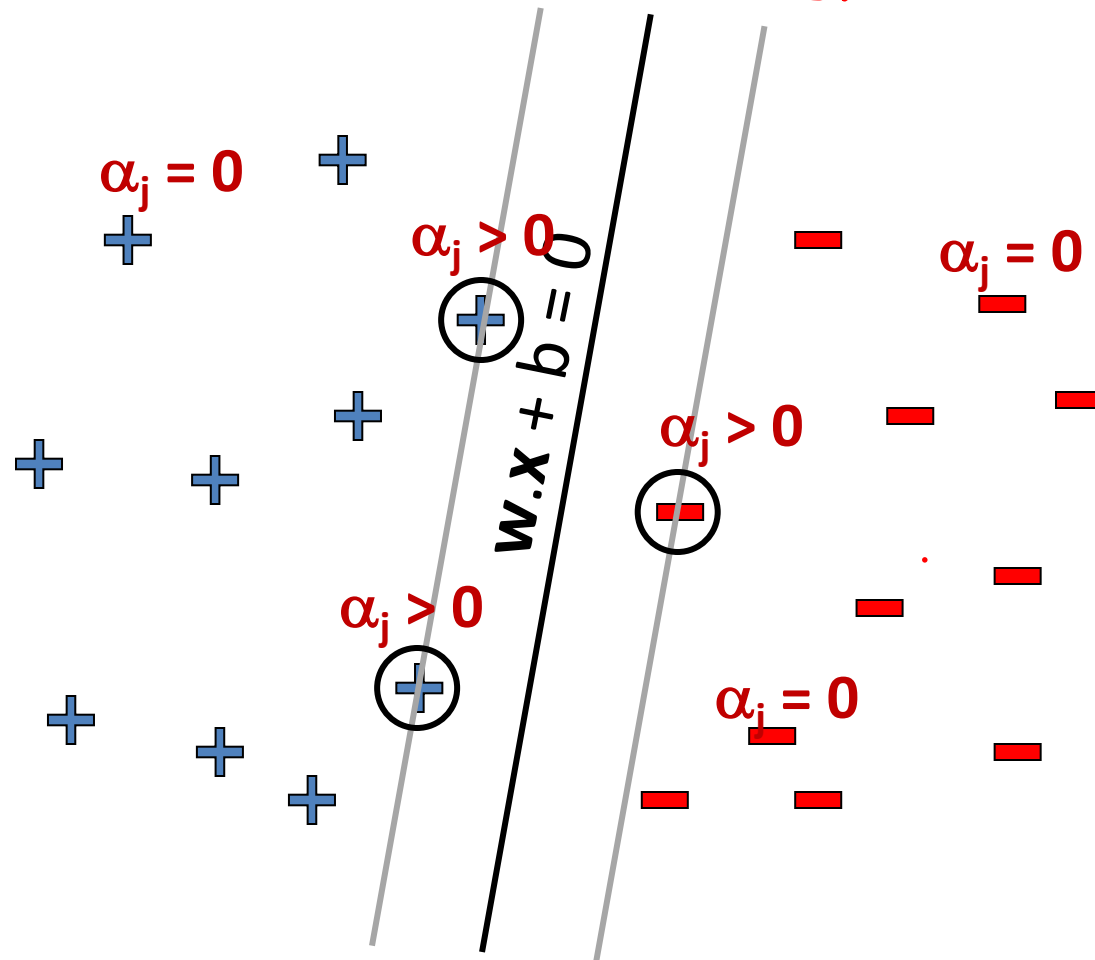
What about  $b$ ?

# Dual SVM: Sparsity of dual solution

complementary slackness

$$\alpha_j ((w \cdot x_j + b) y_j - 1) = 0$$

$$w = \sum_j \alpha_j y_j x_j$$



Only few  $\alpha_j$ s can be non-zero : where constraint is active

$$(w \cdot x_j + b) y_j = 1$$

**Support vectors** – training points  $j$  whose  $\alpha_j$ s are non-zero

# Dual SVM – linearly separable case

$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

Dual problem is also QP

Solution gives  $\alpha_j$ s  $\longrightarrow$

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

$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

for any  $k$  where  $\alpha_k > 0$

Use any one of support vectors with  $\alpha_k > 0$  to compute  $b$  since constraint is tight  $(\mathbf{w} \cdot \mathbf{x}_k + b)y_k = 1$

$$\Rightarrow b = \frac{1}{y_k} - \mathbf{w} \cdot \mathbf{x}_k$$

# Dual SVM – non-separable case

- Primal problem:

$$\begin{aligned}
 & \text{minimize}_{\mathbf{w}, b, \{\xi_j\}} \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\
 & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j \quad \leftarrow \\
 & \xi_j \geq 0, \quad \forall j \quad \leftarrow
 \end{aligned}$$

$\alpha_j$   
 $\mu_j$

$\alpha_j$   
 $\mu_j$

Lagrange  
Multipliers

- Dual problem:

$$\begin{aligned}
 & \max_{\alpha, \mu} \min_{\mathbf{w}, b, \{\xi_j\}} \overbrace{L(\mathbf{w}, b, \xi, \alpha, \mu)}^{d(\alpha, \mu)} \\
 & s.t. \alpha_j \geq 0 \quad \forall j \\
 & \mu_j \geq 0 \quad \forall j
 \end{aligned}$$

# Dual SVM – non-separable case

$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$\underbrace{C}_{\geq} \alpha_i \geq 0 \quad \leftarrow$$

$$C > \alpha_k > 0$$

comes from  $\frac{\partial L}{\partial \xi} = 0$

Intuition:

If  $C \rightarrow \infty$ , recover hard-margin SVM

Dual problem is also QP

Solution gives  $\alpha_j$   $\longrightarrow$

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

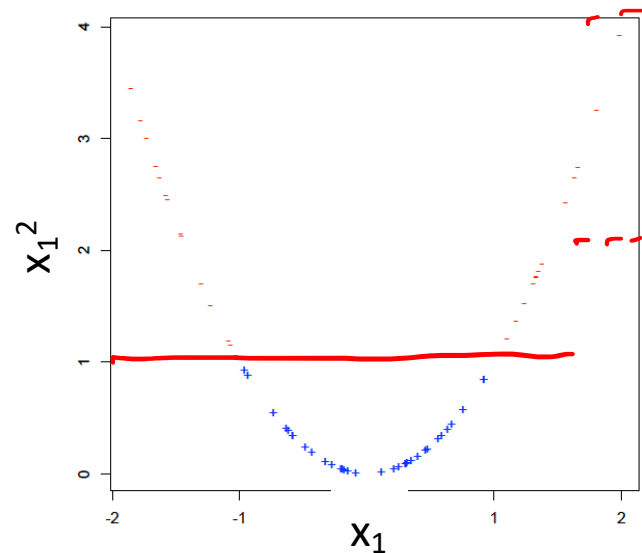
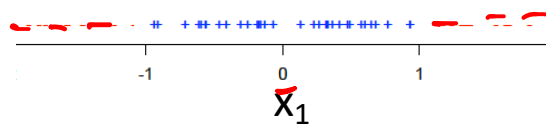
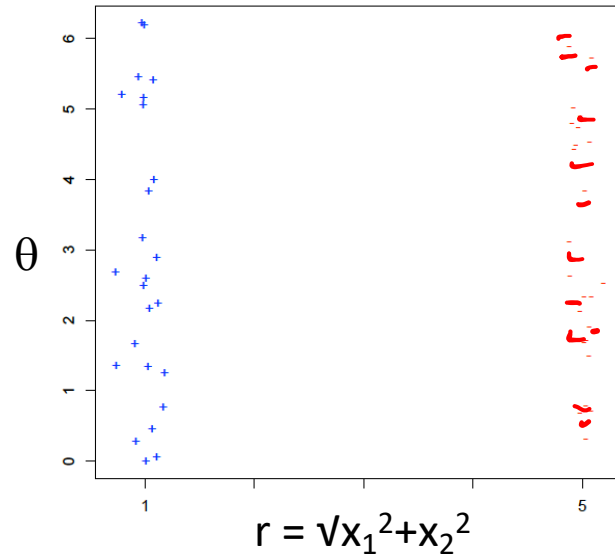
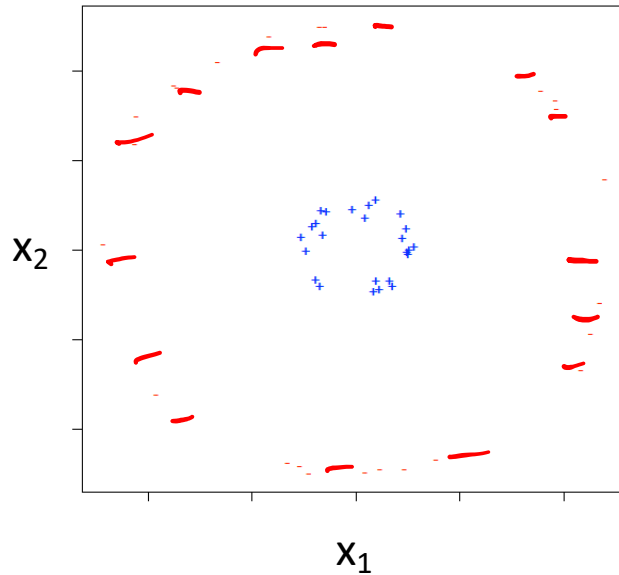
$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

for any  $k$  where  $C > \alpha_k > 0$

# So why solve the dual SVM?

- There are some quadratic programming algorithms that can solve the dual faster than the primal, (specially in high dimensions  $d \gg n$ )
- But, more importantly, the “kernel trick”!!!

# Separable using higher-order features



# Dual formulation only depends on dot-products, not on $w$ !

$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{\mathbf{x}_i \cdot \mathbf{x}_j} \\ & \sum_i \alpha_i y_i = 0 \\ & C \geq \alpha_i \geq 0 \end{aligned}$$

$$\begin{aligned} & \overbrace{\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)} \\ & \phi(\mathbf{x}_i) = \begin{bmatrix} x_i \\ x_i^2 \\ \log x_i \\ e^{x_i} \end{bmatrix} \end{aligned}$$



$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{K(\mathbf{x}_i, \mathbf{x}_j)} \\ & K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) \\ & \sum_i \alpha_i y_i = 0 \\ & C \geq \alpha_i \geq 0 \end{aligned}$$

$\Phi(\mathbf{x})$  – High-dimensional feature space, but never need it explicitly as long as we can compute the dot product fast using some Kernel  $K$



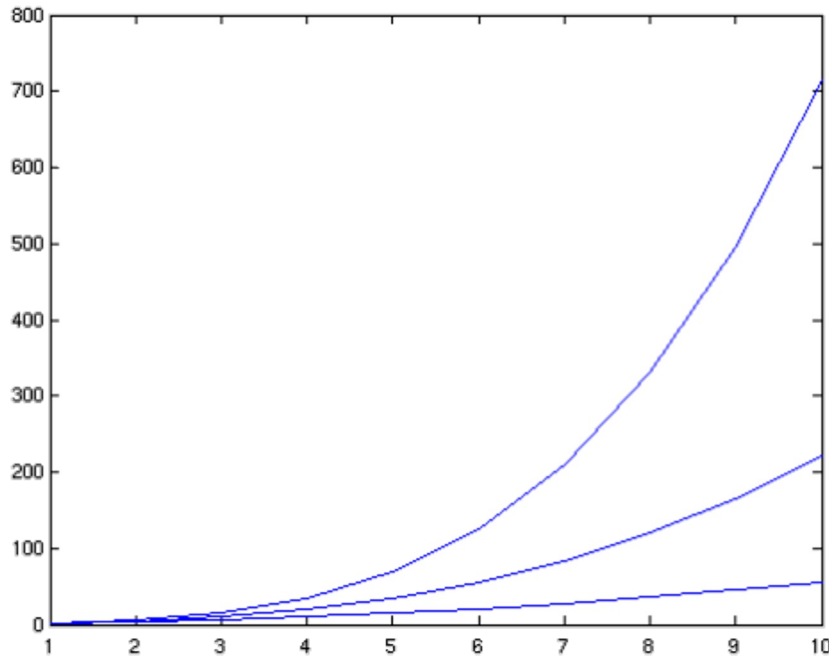
# Polynomial features $\phi(x)$

$$x \quad \phi(x) = \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^d \end{bmatrix}$$

$m$  – input features

$d$  – degree of polynomial

$$\text{num. terms} = \binom{d + m - 1}{d} = \frac{(d + m - 1)!}{d!(m - 1)!} \sim m^d$$



grows fast!

$d = 6, m = 100$

about 1.6 billion terms

$$x = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix}$$

$$\phi(x) = \begin{bmatrix} x^{(1)3} \\ x^{(2)3} \\ \vdots \\ x^{(1)2} x^{(2)} \\ \vdots \\ x^{(1)} x^{(2)} x^{(3)} \end{bmatrix} \quad d=3$$

# Dot Product of Polynomial features

$\Phi(\mathbf{x}) =$  polynomials of degree exactly  $d$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$d=1 \quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \underline{x_1 z_1} + \underline{x_2 z_2} = \underline{\mathbf{x} \cdot \mathbf{z}}$$

$$\begin{aligned} d=2 \quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) &= \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix} = x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 \\ &= (x_1 z_1 + x_2 z_2)^2 \\ &= (\mathbf{x} \cdot \mathbf{z})^2 \end{aligned}$$

$d \geq 6$

$$d \quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = K(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z})^d$$