

# Support Vector Machines (SVMs) contd...

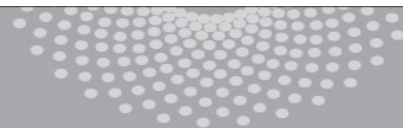
Aarti Singh

Machine Learning 10-315

Mar 23, 2022



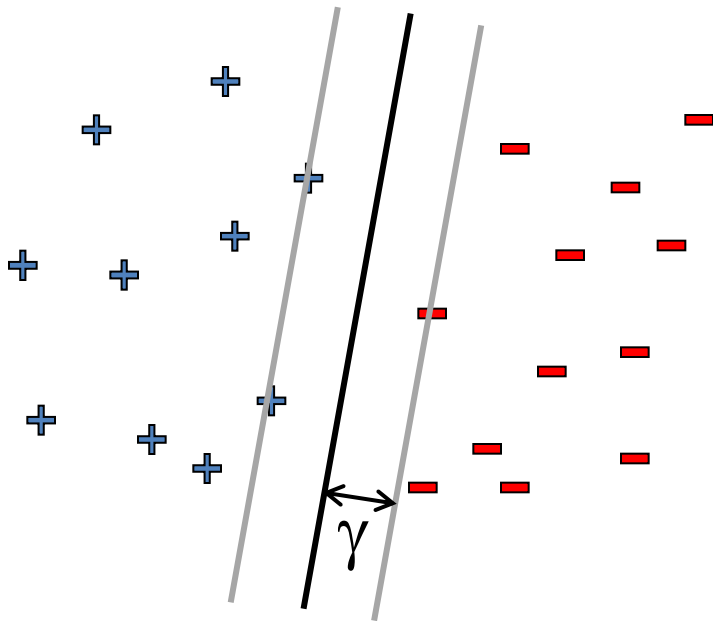
**MACHINE LEARNING** DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science

# Hard-margin SVM

Data perfectly separable by a linear decision boundary



Hard margin approach

$$\min_{\mathbf{w}, b} \mathbf{w} \cdot \mathbf{w}$$

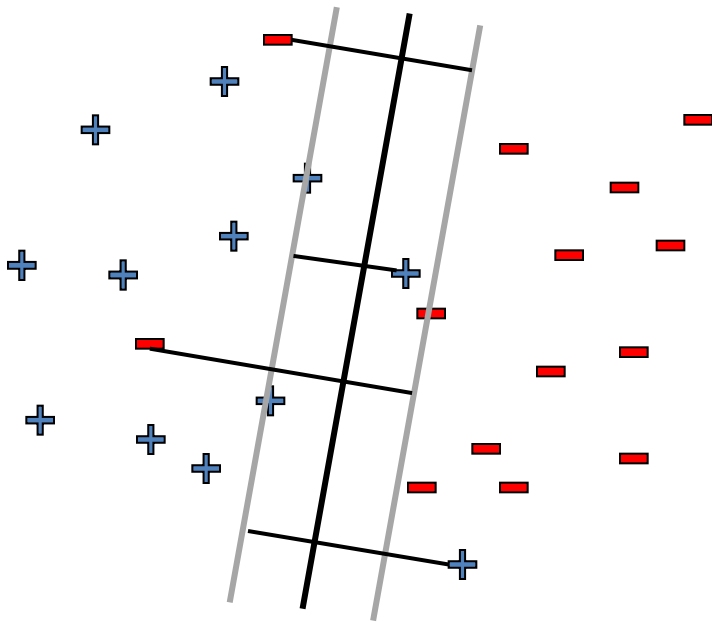
$$\text{s.t. } (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 \quad \forall j$$

Solve using Quadratic Programming (QP)

$$\text{Margin, } \gamma \propto 1/\|\mathbf{w}\|$$

# Soft-margin SVM

Allow “error” in classification



**Soft margin approach**

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

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

$C$  - tradeoff parameter (chosen by cross-validation)

Still QP 😊

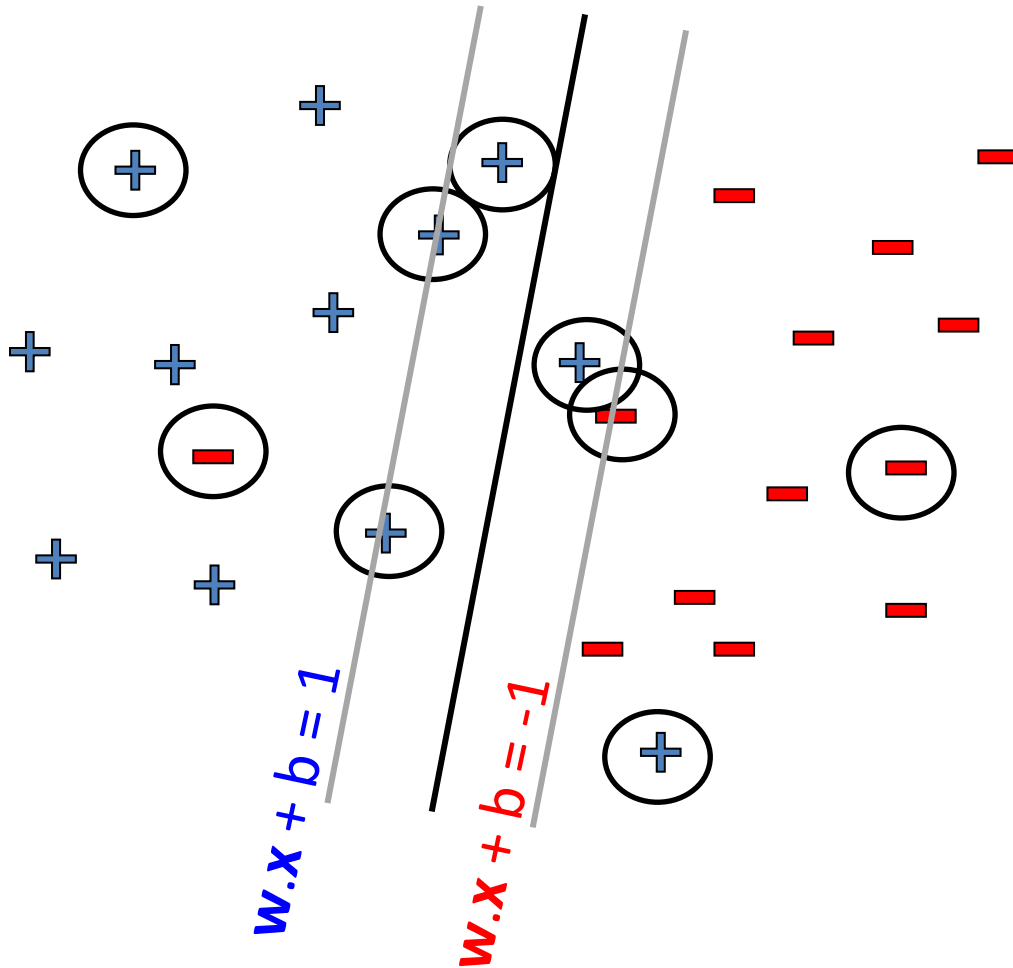
$$\begin{aligned}
 \min_{\mathbf{w}, b, \{\xi_j\}} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum \xi_j \\
 \text{s.t.} \quad & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j \\
 & \xi_j \geq 0 \quad \forall j
 \end{aligned}$$

# Slack variables

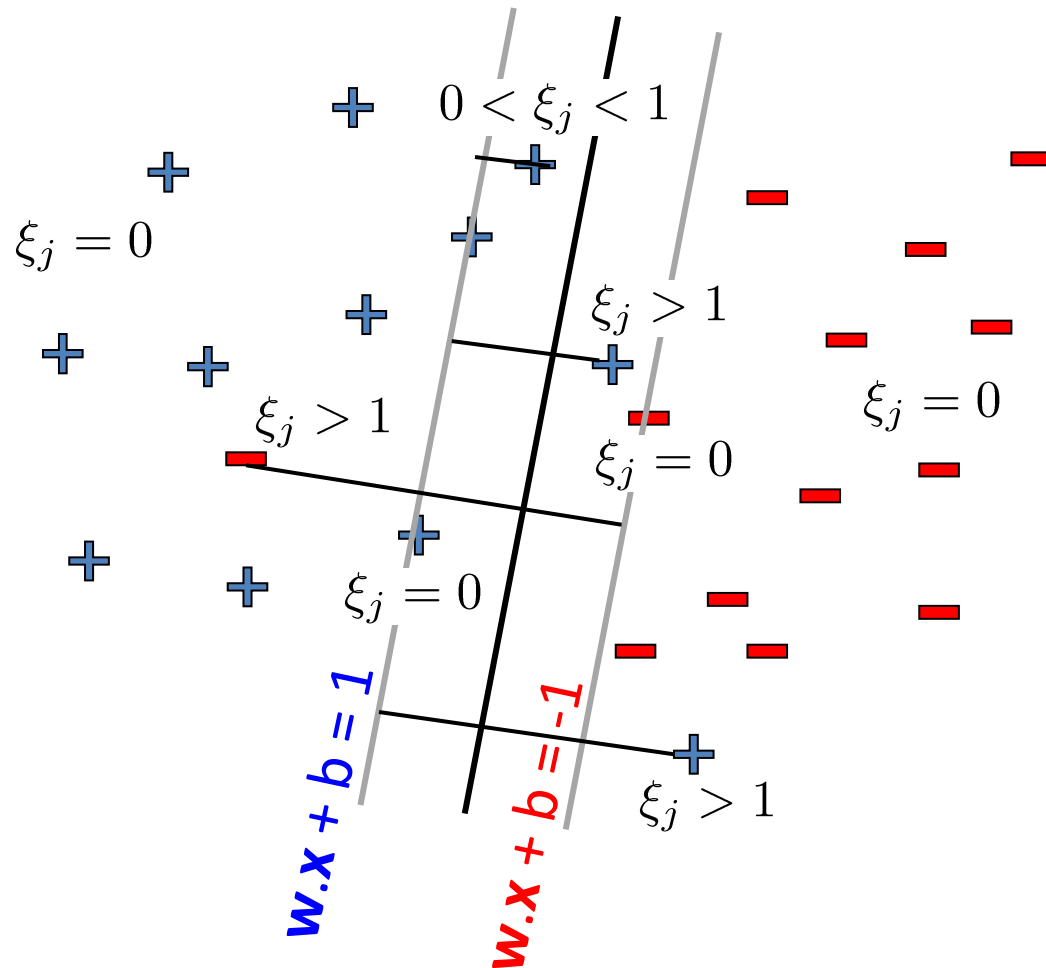
$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j$$

What is the slack  $\xi_j$  for the following points?

Confidence | Slack

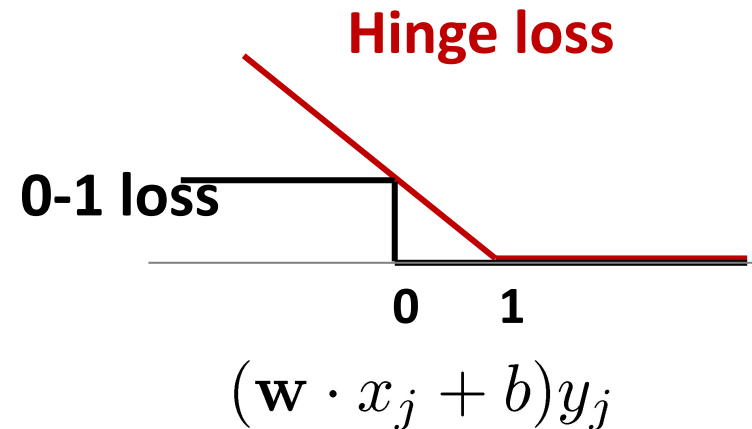


# Slack variables – Hinge loss



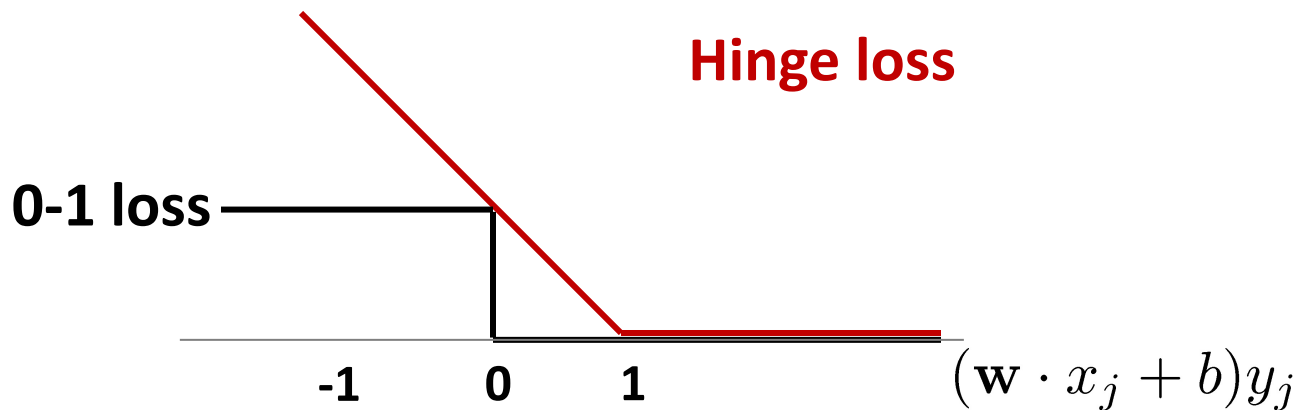
Notice that

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



# Slack variables – Hinge loss

$$\xi_j = (1 - (\mathbf{w} \cdot x_j + b)y_j)_+$$



$$\min_{\mathbf{w}, b, \{\xi_j\}} \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j$$

$$\text{s.t. } (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j$$

$$\xi_j \geq 0 \quad \forall j$$

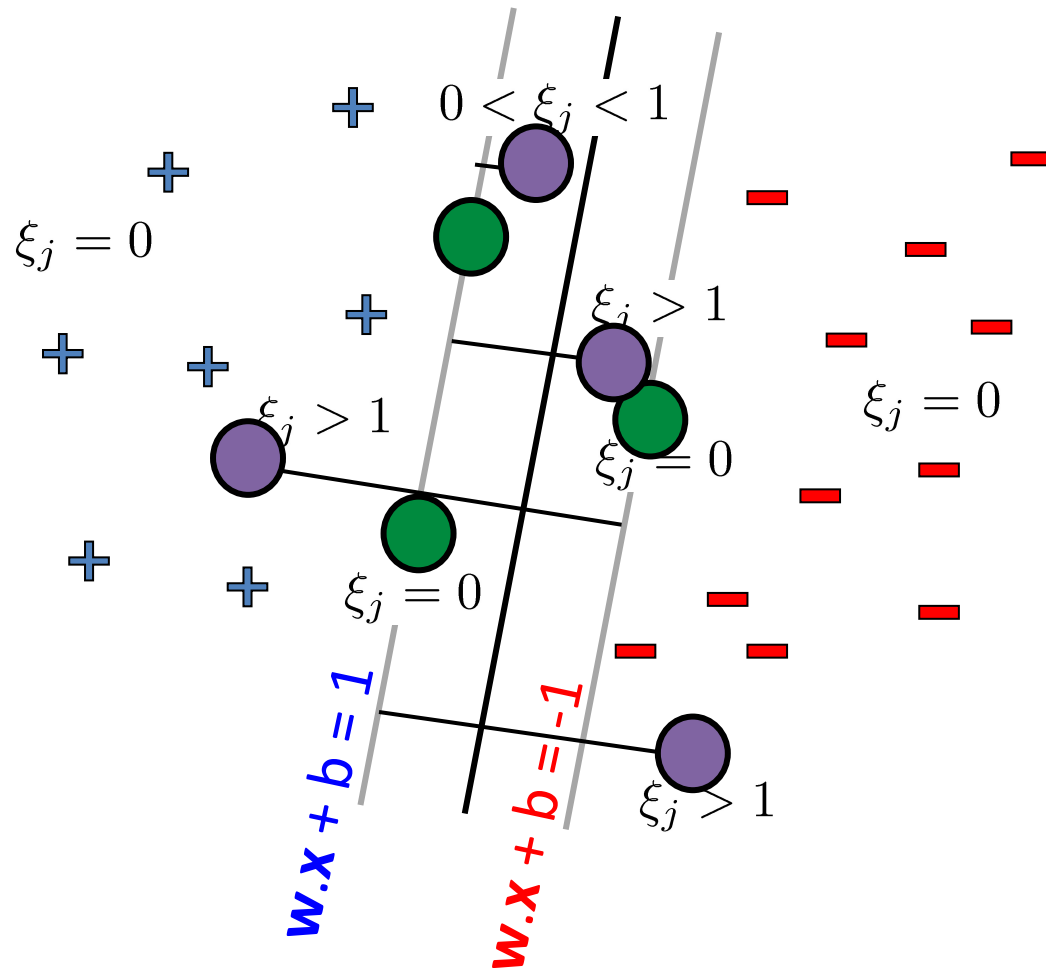


Regularized hinge loss

$$\min_{\mathbf{w}, b} \mathbf{w} \cdot \mathbf{w} + C \sum_j (1 - (\mathbf{w} \cdot \mathbf{x}_j + b)y_j)_+$$

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

# Support Vectors



## Margin support vectors

$\xi_j = 0$ ,  $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j = 1$   
(don't contribute to objective but enforce constraints on solution)

Correctly classified but on margin

## Non-margin support vectors

$\xi_j > 0$   
(contribute to both objective and constraints)

$1 > \xi_j > 0$  Correctly classified but inside margin

$\xi_j > 1$  Incorrectly classified

# SVM – linearly separable case

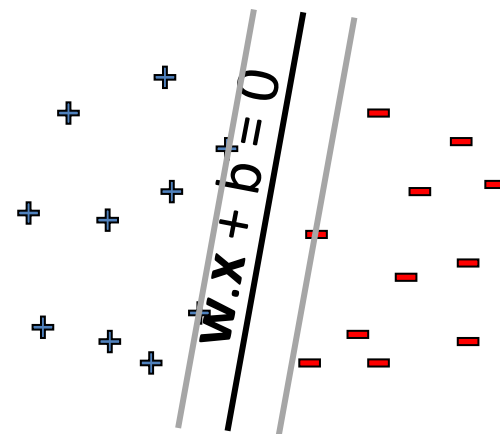
n training points

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

d features

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

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



**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)

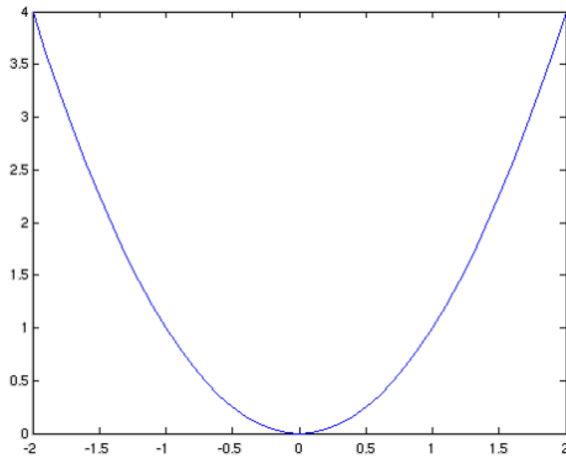


# Detour - Constrained Optimization

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

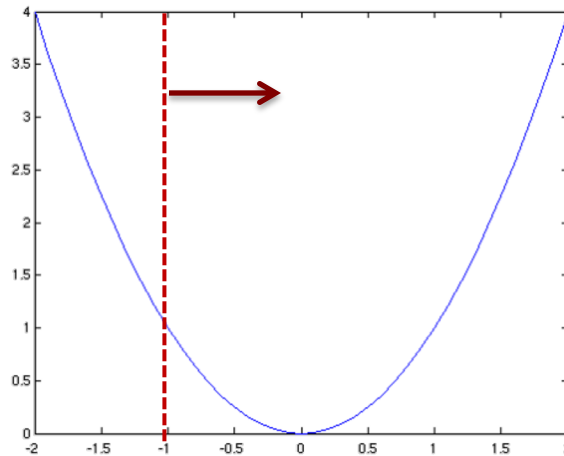
$$x^* = \max(b, 0)$$

$$\min_x x^2$$



$$x^* = 0$$

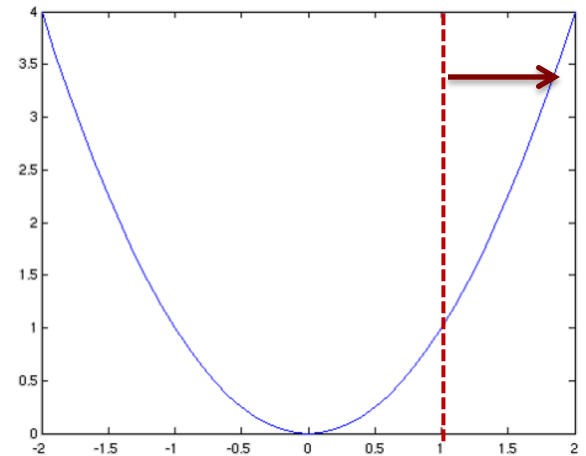
$$\begin{array}{ll}\min_x & x^2 \\ \text{s.t.} & x \geq -1\end{array}$$



$$x^* = 0$$

Constraint inactive

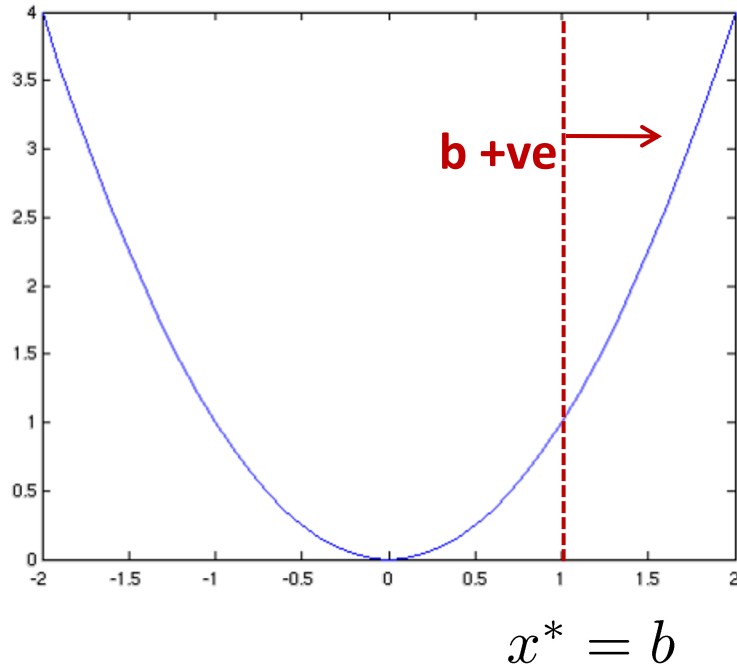
$$\begin{array}{ll}\min_x & x^2 \\ \text{s.t.} & x \geq 1\end{array}$$



$$x^* = 1$$

Constraint active  
(tight)

# Constrained Optimization



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

Equivalent unconstrained optimization:  
 $\min_x x^2 + l(x-b)$

Replace with lower bound ( $\alpha \geq 0$ )  
$$x^2 + l(x-b) \geq \underbrace{x^2 - \alpha(x-b)}_{L(x,\alpha)}$$

# Primal and Dual Problems

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

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

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

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

s.t.  $\alpha \geq 0$

where Lagrangian  $L(x, \alpha) = x^2 - \alpha(x - b)$

How to form the Lagrangian?

For each constraint, introduce a positive Lagrange multiplier  
Fold constraints into objective

# Why solve the Dual problem?

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

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

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

- **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

**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.  $d^* \leq p^*$
- **Strong duality:**  $d^* = p^*$  holds often for many problems of interest e.g. if the primal is a feasible convex objective with linear constraints

# 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)?

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

- 1.  $\nabla L(x^*, \alpha^*) = 0$  i.e. Gradient of Lagrangian at  $x^*$  and  $\alpha^*$  is zero.
- 2.  $x^* \geq b$  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)

We use the first one to relate  $x^*$  and  $\alpha^*$ . We use the last one (complimentary slackness) to argue that  $\alpha^* = 0$  if constraint is inactive and  $\alpha^* > 0$  if constraint is active and tight.

# Primal and Dual Problems

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

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

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

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

s.t.  $\alpha \geq 0$

where Lagrangian  $L(x, \alpha) = x^2 - \alpha(x - b)$

How to form the Lagrangian?

For each constraint, introduce a positive Lagrange multiplier  
Fold constraints into objective



# 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: 
$$\begin{aligned} &\text{minimize}_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w} \cdot \mathbf{w} \\ &\quad \left( \mathbf{w} \cdot \mathbf{x}_j + b \right) y_j \geq 1, \quad \forall j \end{aligned}$$

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

- Dual problem (derivation):

$$\begin{aligned} L(\mathbf{w}, b, \alpha) &= \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j \left[ \left( \mathbf{w} \cdot \mathbf{x}_j + b \right) y_j - 1 \right] \\ \alpha_j &\geq 0, \quad \forall j \end{aligned}$$

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

# Dual SVM – linearly separable case

- Dual problem:

$$\max_{\alpha} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j \left[ (\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1 \right]$$
$$\alpha_j \geq 0, \quad \forall j$$

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

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

If we can solve for  $\alpha$ s (dual problem), then we have a solution for  $\mathbf{w}$  (primal problem)

# Dual SVM – linearly separable case

- Dual problem:

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

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

$$\Rightarrow \mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j \qquad \Rightarrow \sum_j \alpha_j y_j = 0$$

# 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

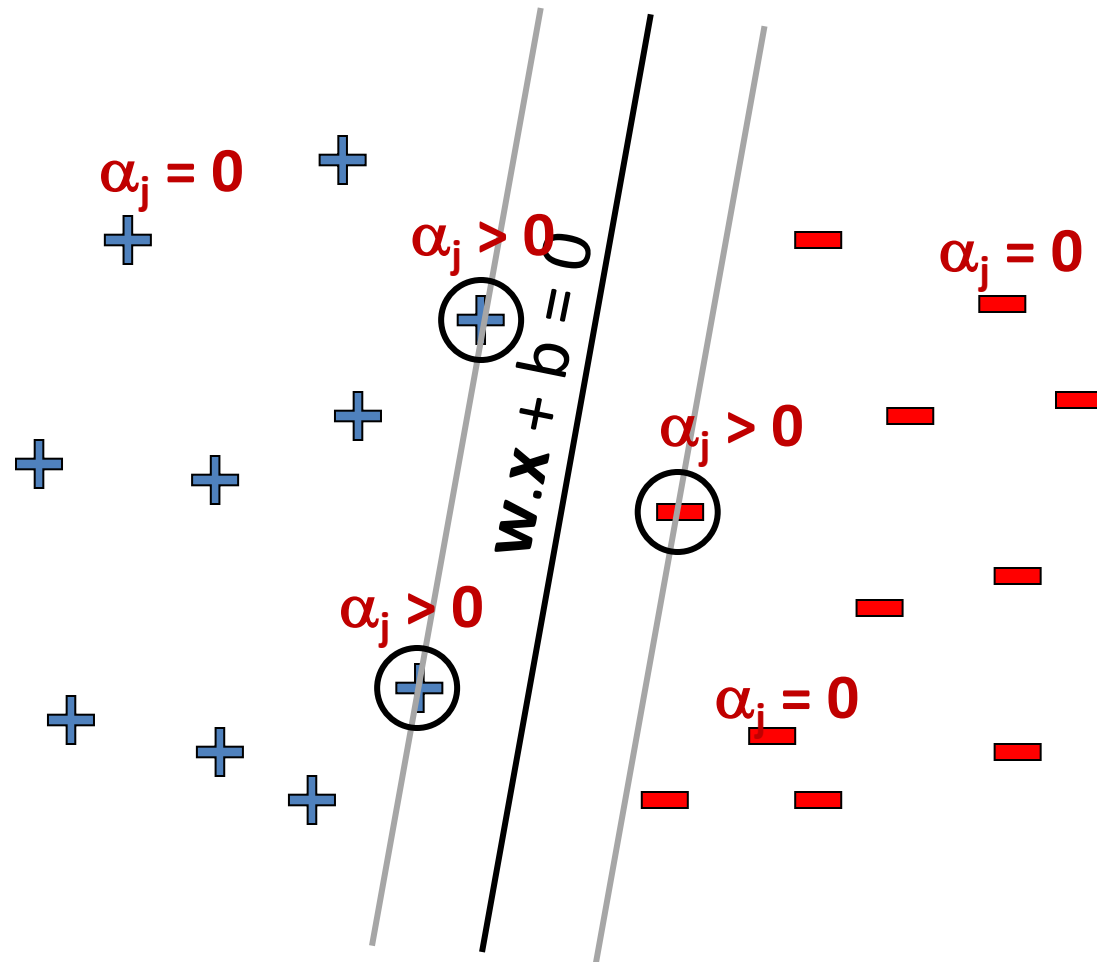
Solution gives  $\alpha_j$ s



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

What about b?

# Dual SVM: Sparsity of dual solution



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

**Complementary slackness implies**  
Only few  $\alpha_j$ s can be non-zero : where constraint is active and tight

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j = 1$$

**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  $\longrightarrow$

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

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

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

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

# 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**”!!!