

SVM Review

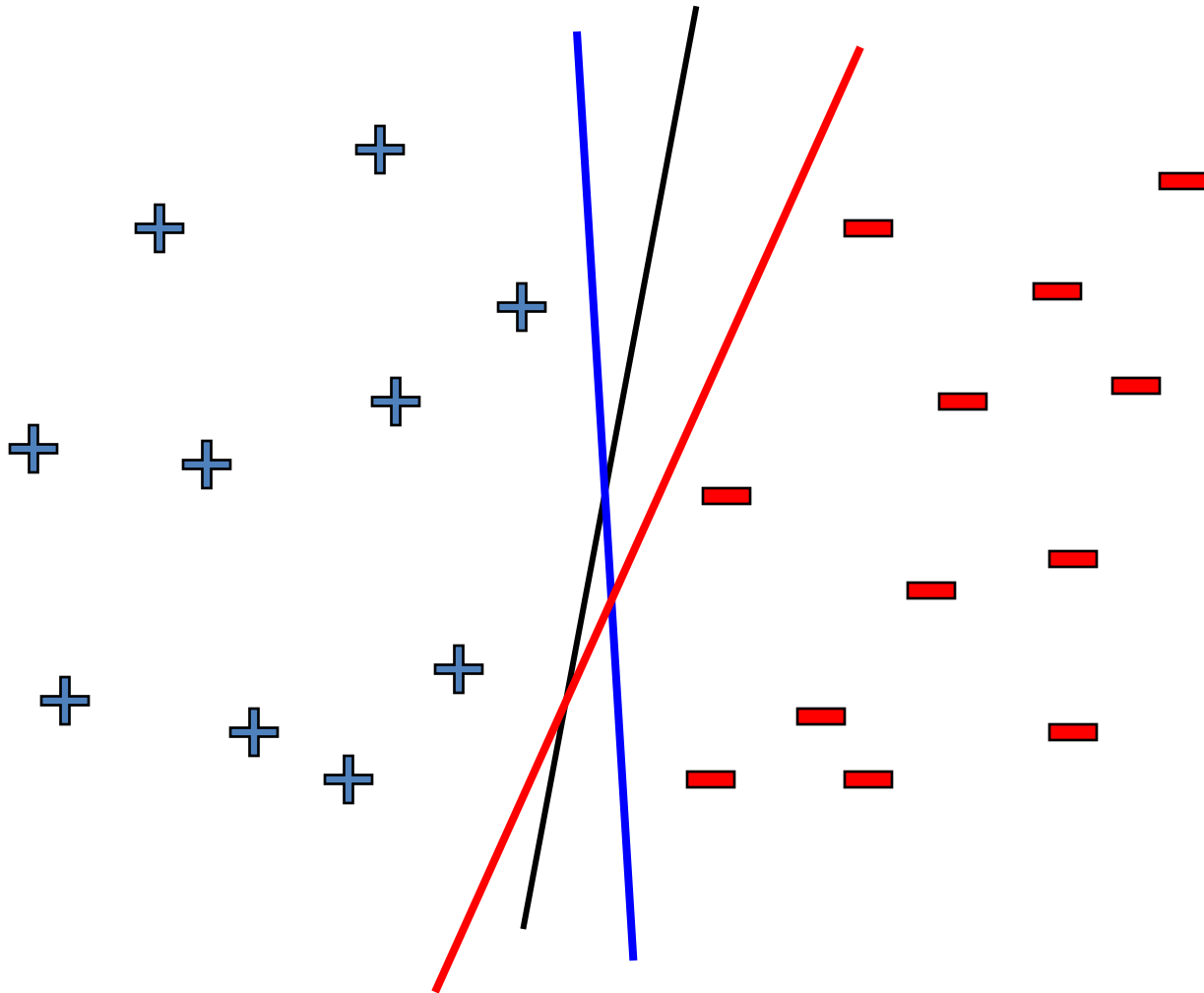
Siddharth Ancha

Slides from Aarti's Lectures

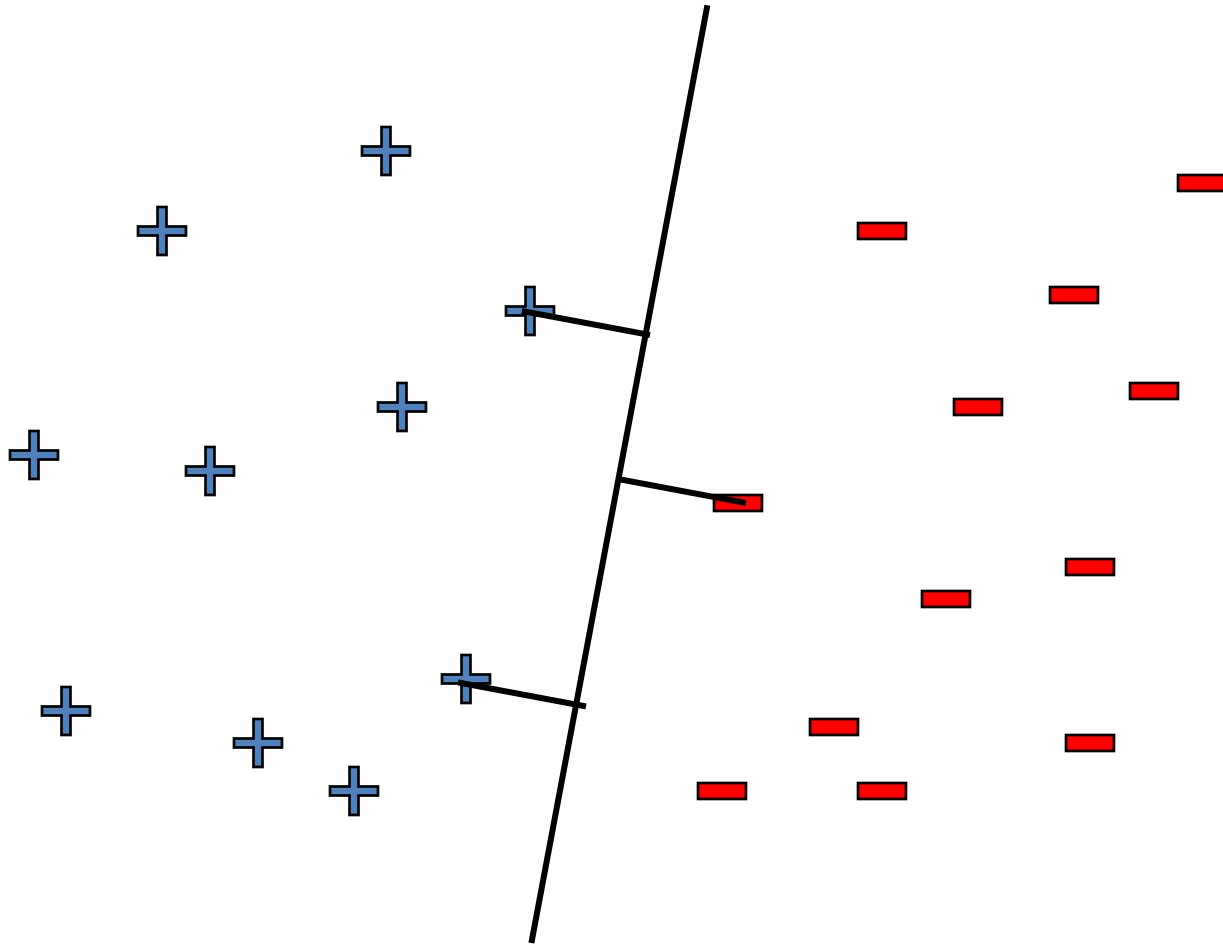
At Pittsburgh G-20 summit ...



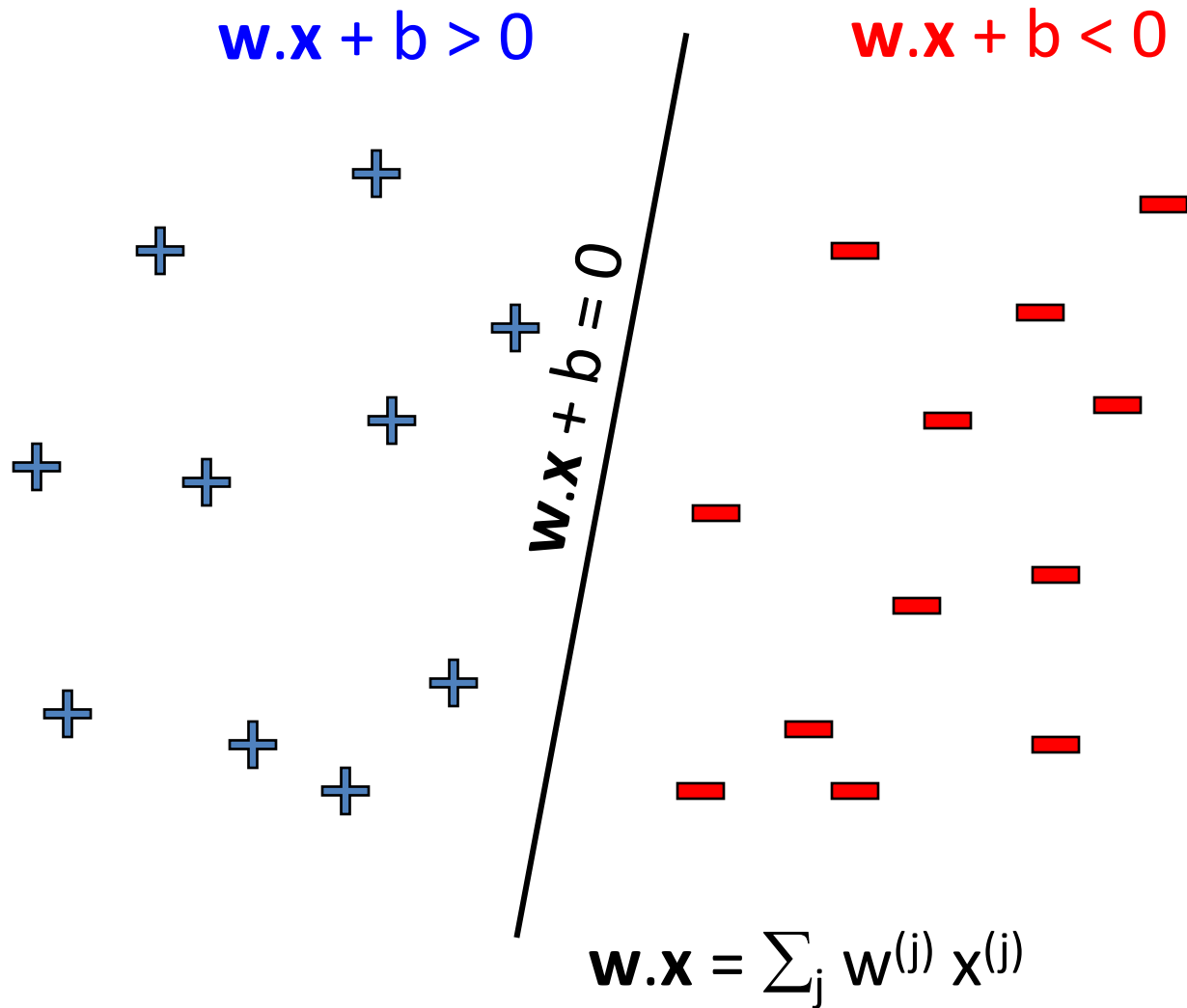
Linear classifiers – which line is better?



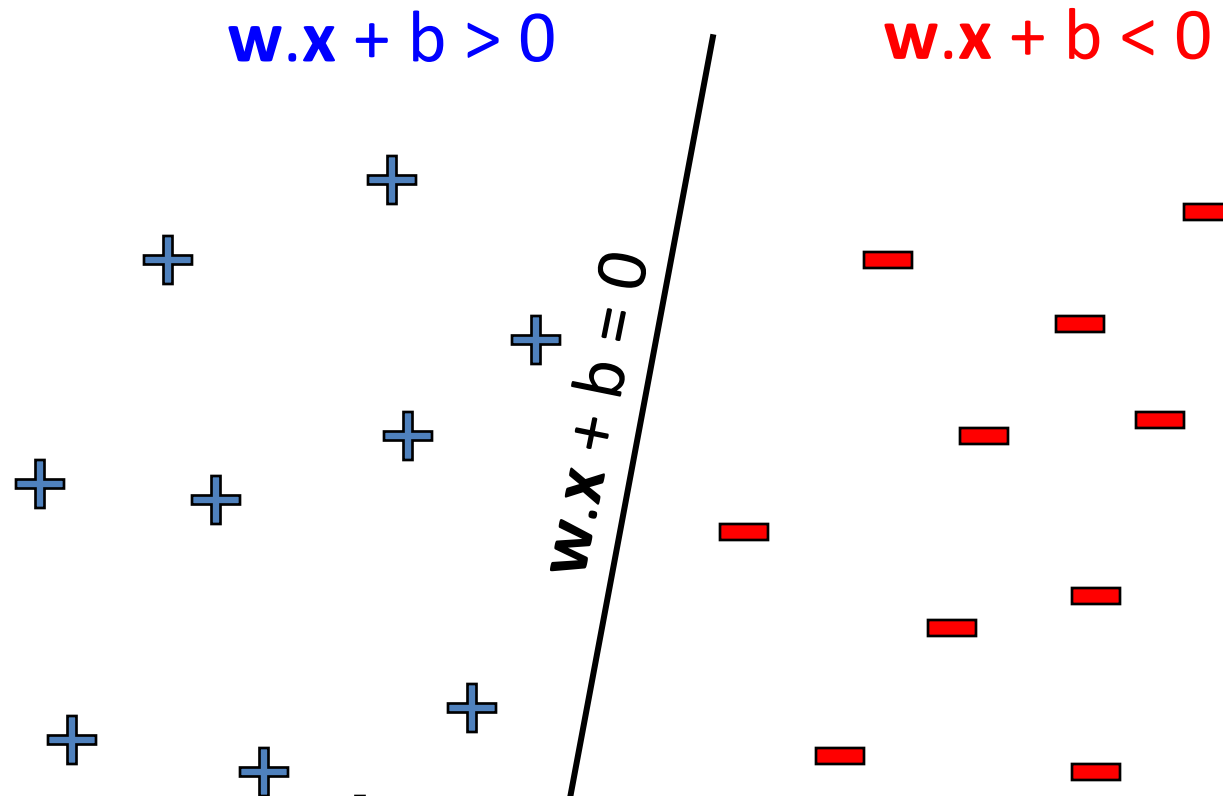
Pick the one with the largest margin!



Parameterizing the decision boundary



Parameterizing the decision boundary



$y_j \in \{-1, +1\}$ — class

“confidence” $= (w \cdot x_j + b) y_j$

Maximizing the margin

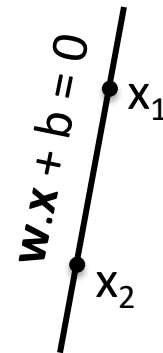
$$w \cdot x + b > 0$$

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

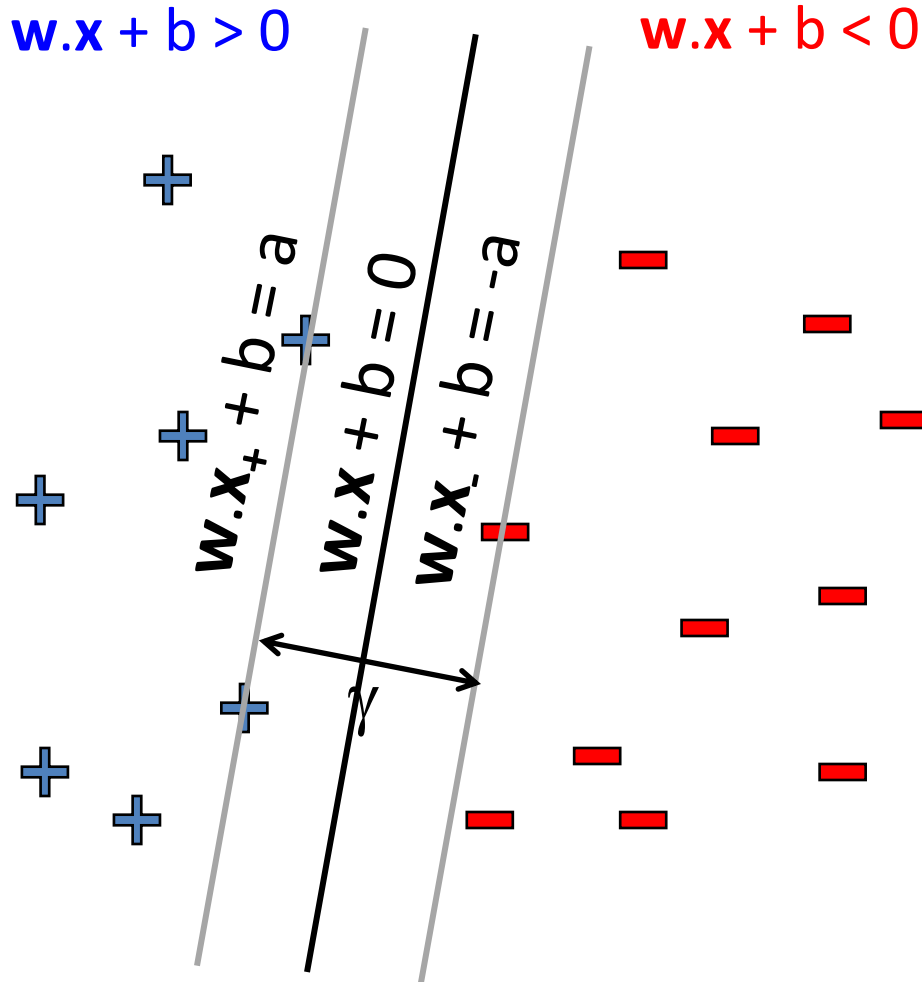
Distance of closest examples from the line/hyperplane

$$\text{margin} = \gamma = 2a / \|w\|$$

Step 1: w is perpendicular to lines since for any x_1, x_2 on line $w \cdot (x_1 - x_2) = 0$



Maximizing the margin



$$\text{margin} = \gamma = 2a / \|w\|$$

Step 1: w is perpendicular to lines

Step 2: Take a point x_- on $w \cdot x + b = -a$ and move to point x_+ that is γ away on line $w \cdot x + b = a$

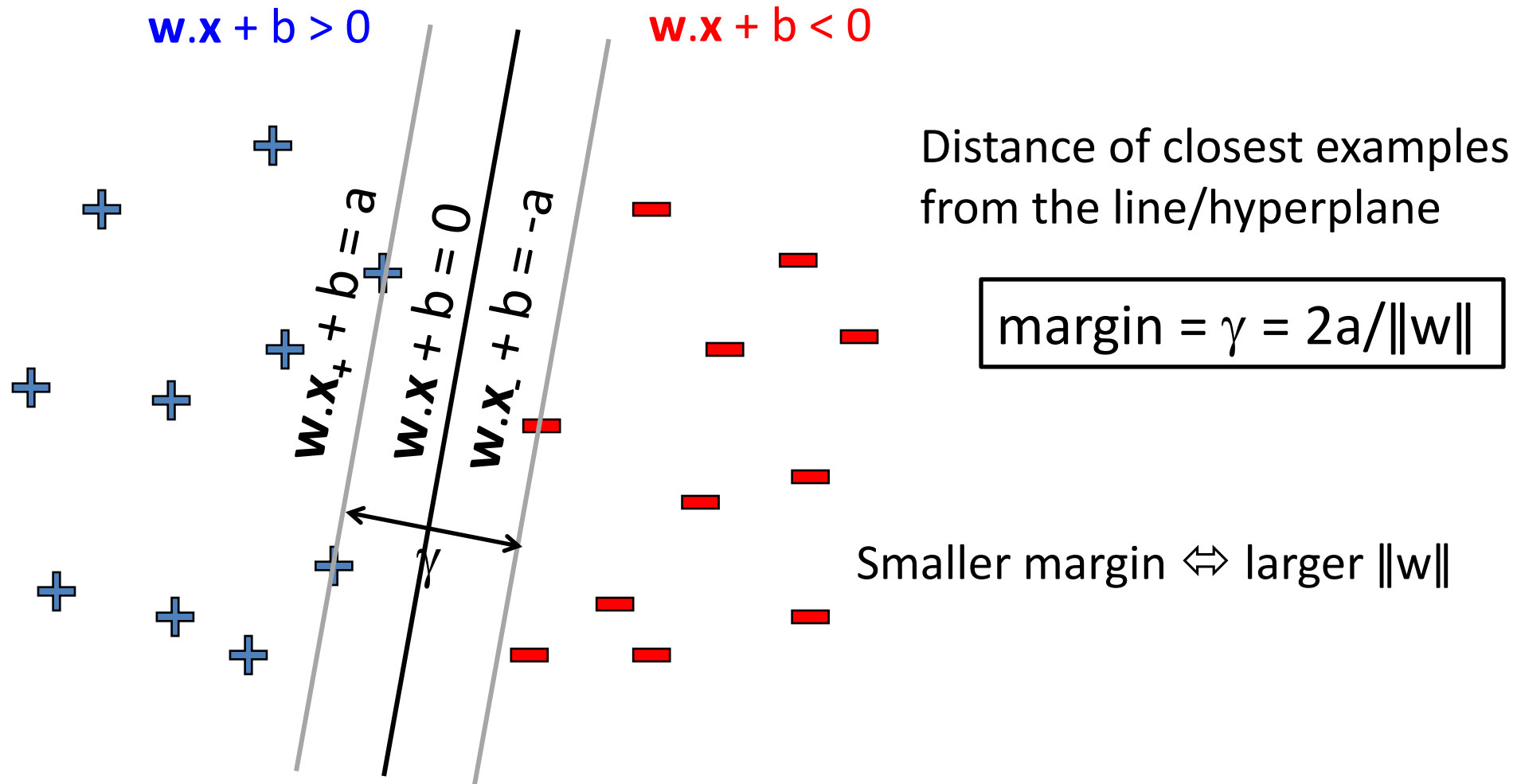
$$x_+ = x_- + \gamma w / \|w\|$$

$$w \cdot x_+ = w \cdot x_- + \gamma w \cdot w / \|w\|$$

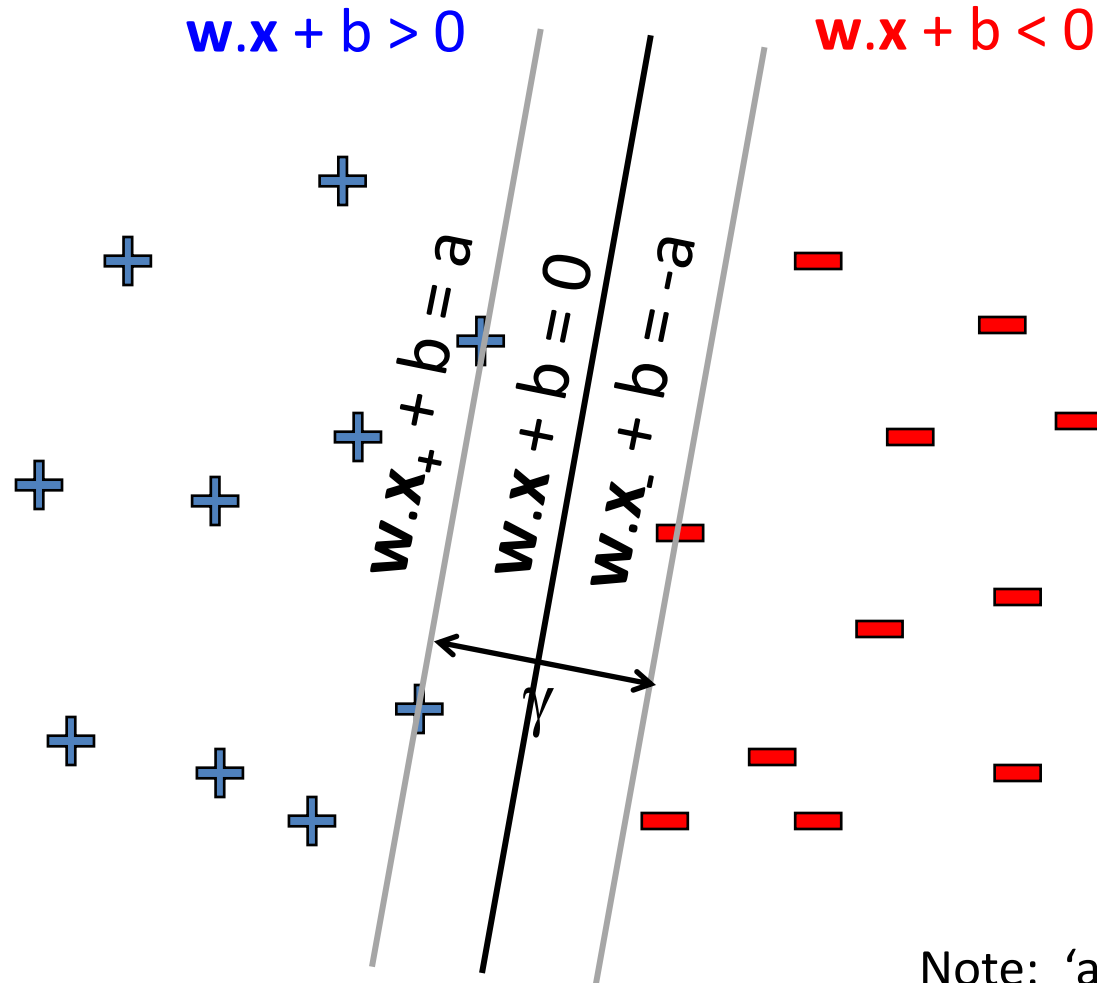
$$a - b = -a - b + \gamma \|w\|$$

$$2a = \gamma \|w\|$$

Maximizing the margin



Maximizing the margin



Distance of closest examples from the line/hyperplane

$$\text{margin} = \gamma = 2a / \|w\|$$

$$\begin{aligned} \max_{w, b} \quad & \gamma = 2a / \|w\| \\ \text{s.t.} \quad & (w \cdot x_j + b) y_j \geq a \quad \forall j \end{aligned}$$

Note: 'a' is arbitrary (can normalize equations by a)

Support Vector Machines

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

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

$$w \cdot x_+ + b = 1$$
$$w \cdot x + b = 0$$
$$w \cdot x_- + b = -1$$

γ

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

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

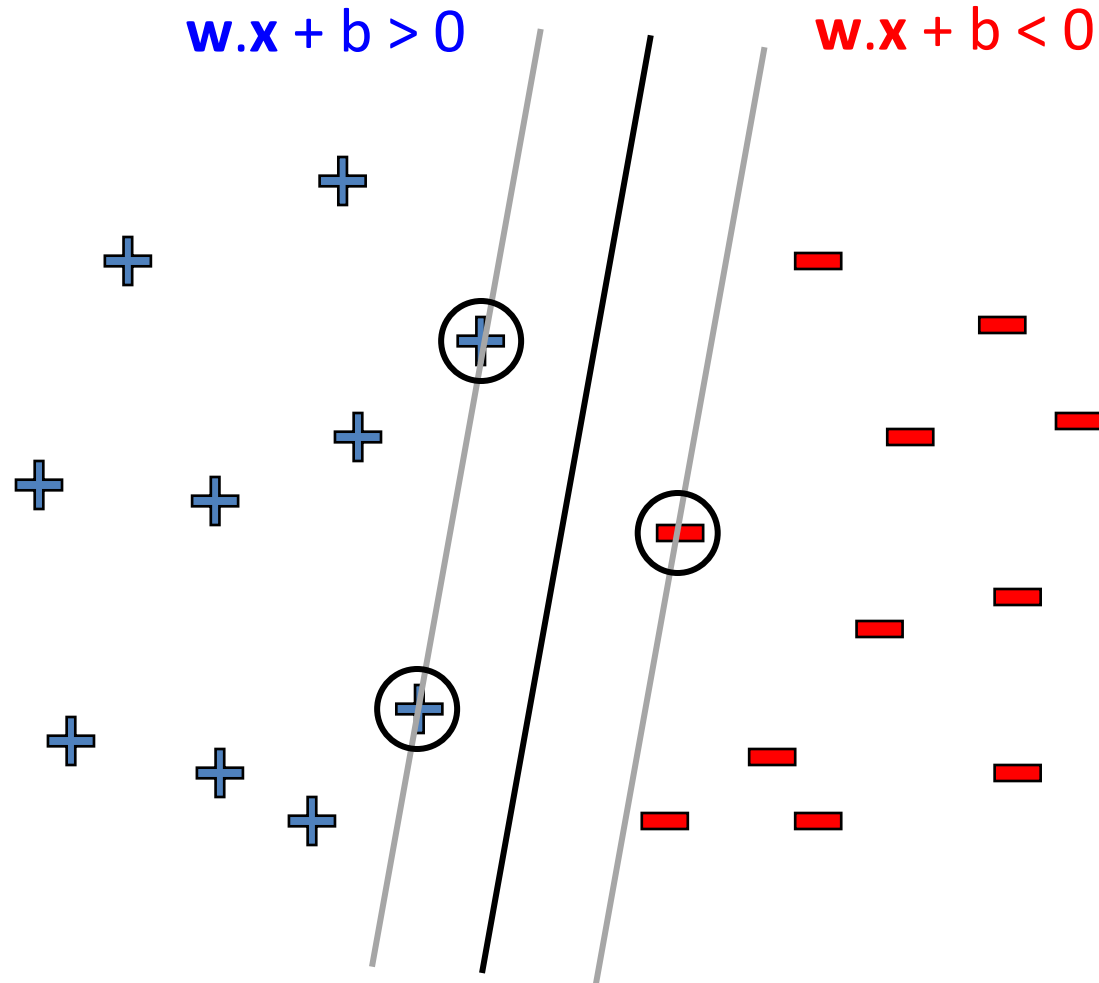
Solve efficiently by quadratic programming (QP)

- Quadratic objective, linear constraints
- Well-studied solution algorithms

Support Vectors

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

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



Linear hyperplane defined by
“support vectors”

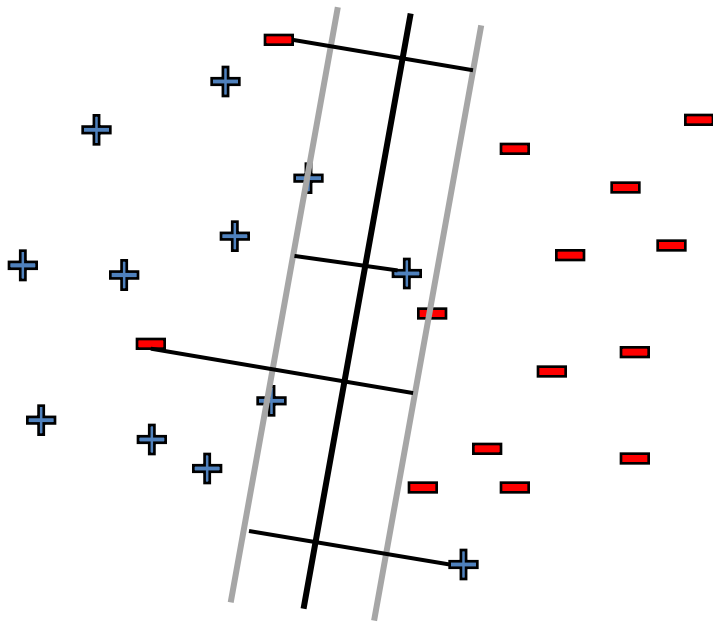
Moving other points a little
doesn't effect the decision
boundary

only need to store the
support vectors to predict
labels of new points

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

What if data is still not linearly separable?

Allow “error” in classification



Soft margin approach

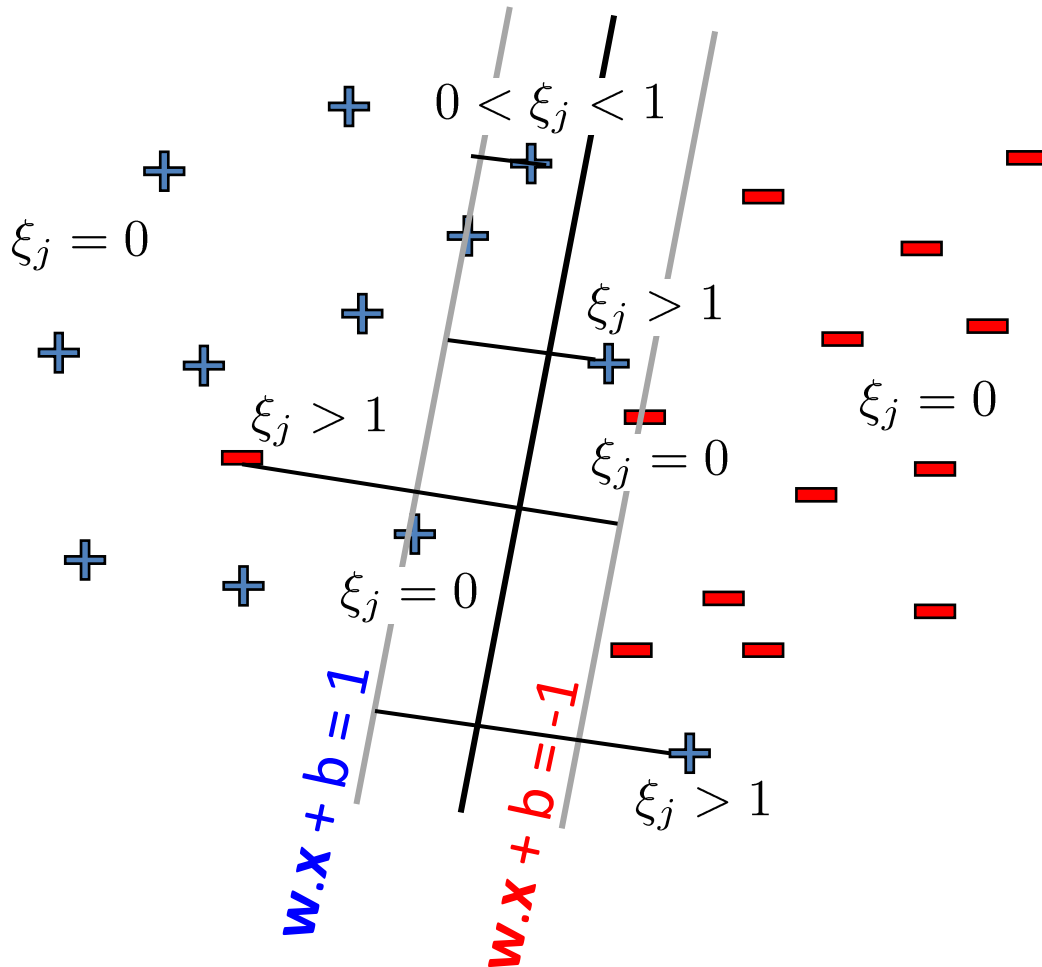
$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_j\}} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ \text{s.t.} \quad & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j \\ & \xi_j \geq 0 \quad \forall j \end{aligned}$$

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

C - tradeoff parameter (chosen by cross-validation)

Still QP 😊

Soft-margin SVM



Soften the constraints:

$$(w \cdot x_j + b) y_j \geq 1 - \xi_j \quad \forall j$$

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

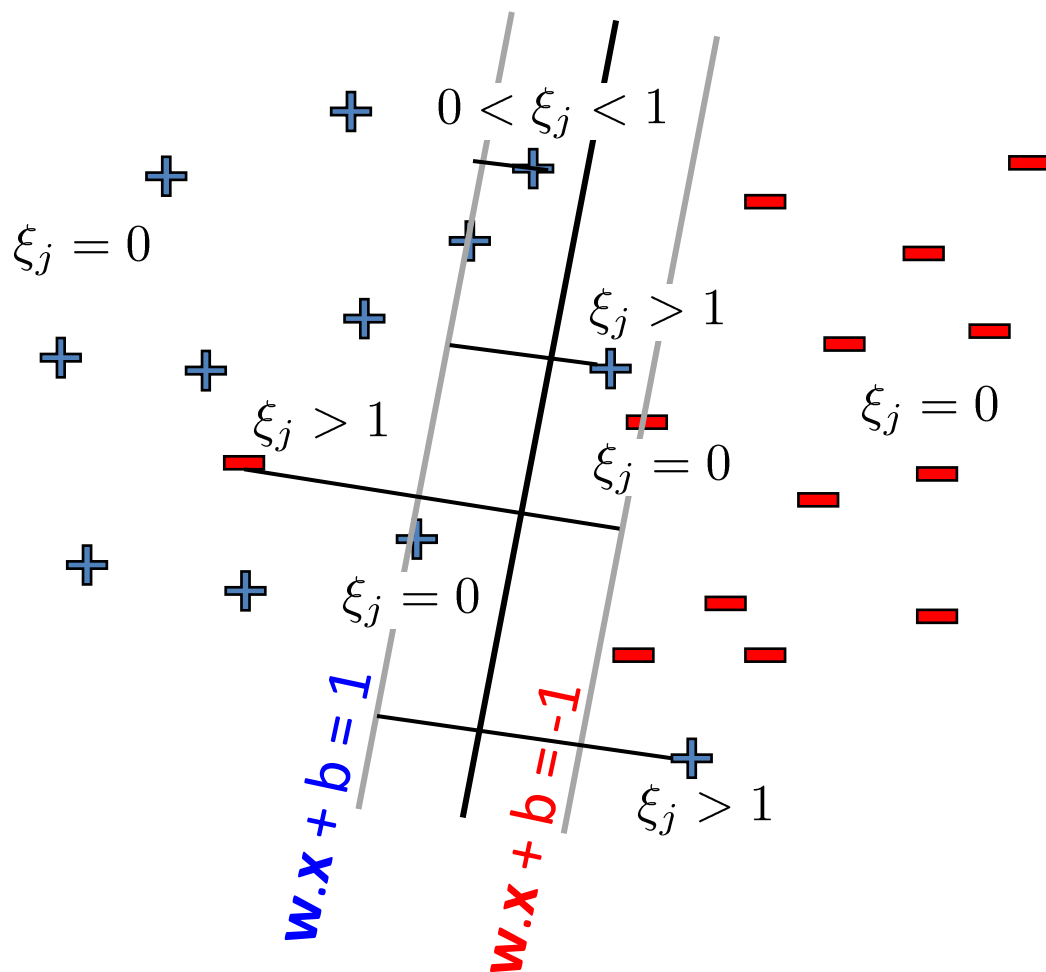
Penalty for misclassifying:

$$C \xi_j$$

How do we recover hard margin SVM?

Set $C = \infty$

Slack variables – Hinge loss

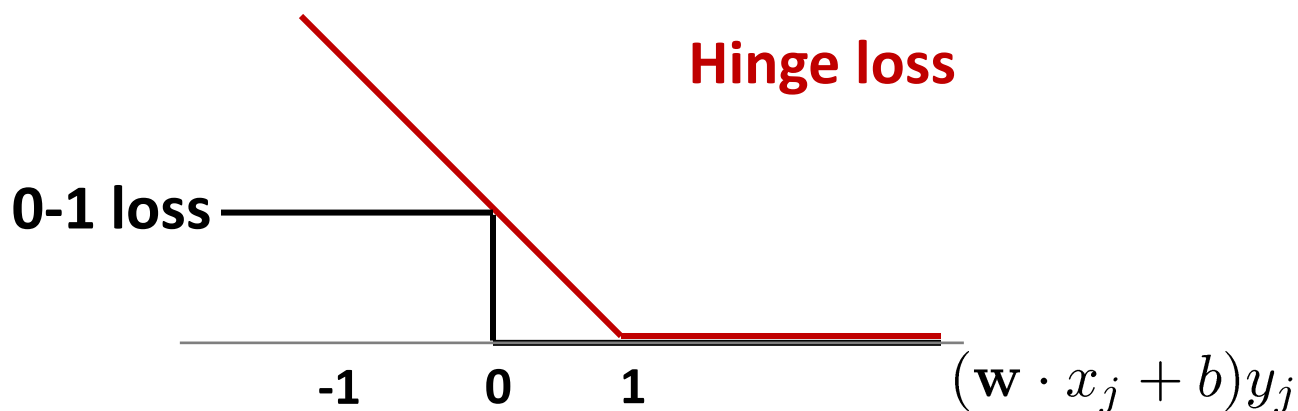


Notice that

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

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

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

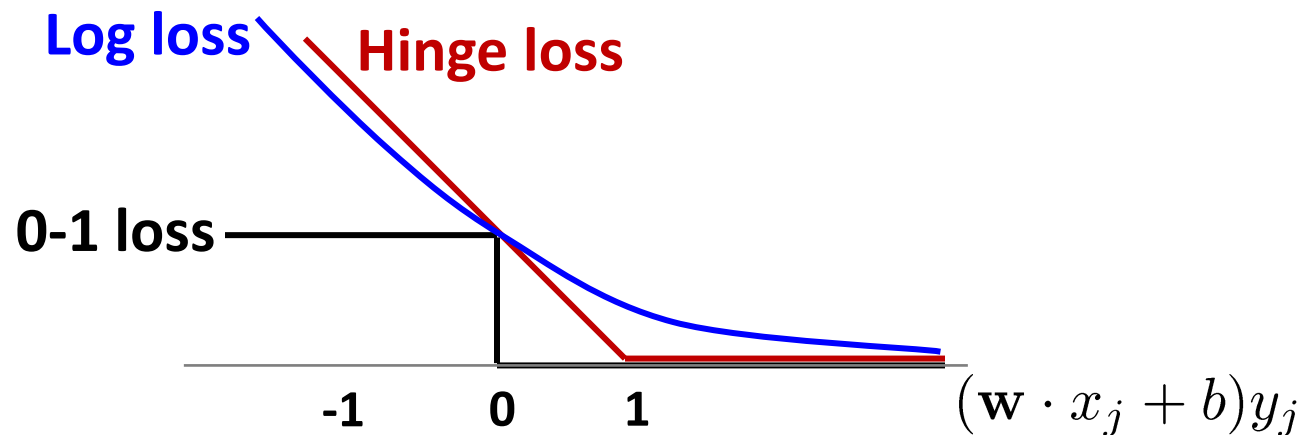
SVM vs. Logistic Regression

SVM : **Hinge loss**

$$\text{loss}(f(x_j), y_j) = (1 - (\mathbf{w} \cdot x_j + b)y_j)_+$$

Logistic Regression : **Log loss** (-ve log conditional likelihood)

$$\text{loss}(f(x_j), y_j) = -\log P(y_j | x_j, \mathbf{w}, b) = \log(1 + e^{-(\mathbf{w} \cdot x_j + b)y_j})$$



SVM – linearly separable case

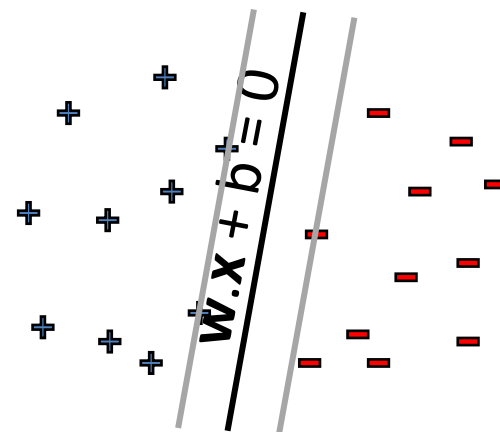
n training points

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

d features

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

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



w - weights on features (d-dim problem)

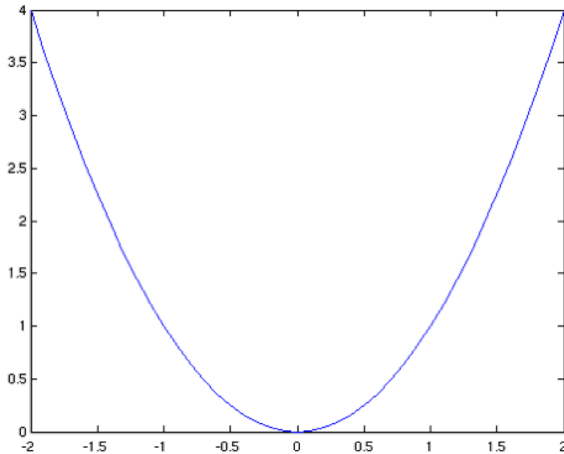
- Convex quadratic program – quadratic objective, linear constraints
- But expensive to solve if d is very large
- Often solved in dual form (n-dim problem)

Constrained Optimization

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

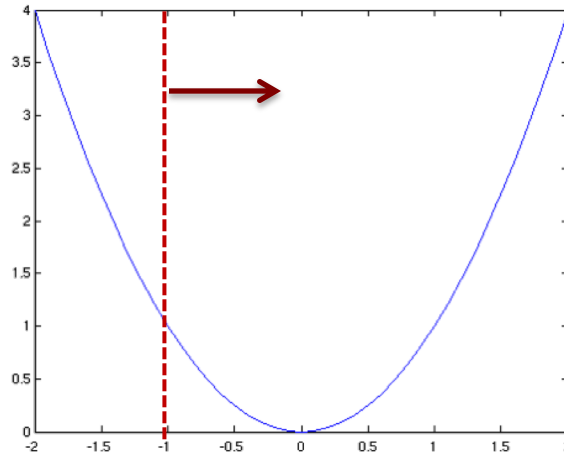
$$x^* = \max(b, 0)$$

$$\min_x x^2$$



$$x^* = 0$$

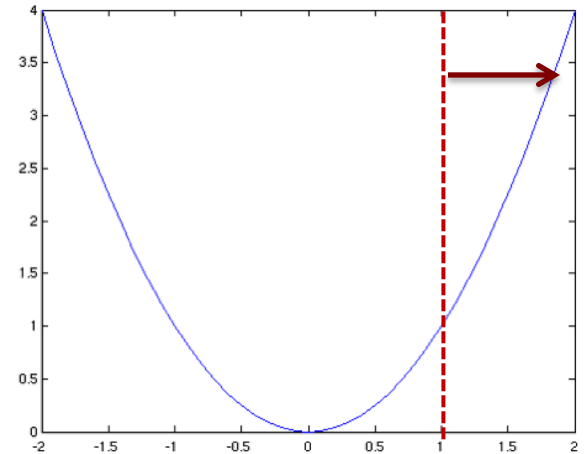
$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq -1 \end{aligned}$$



$$x^* = 0$$

Constraint inactive

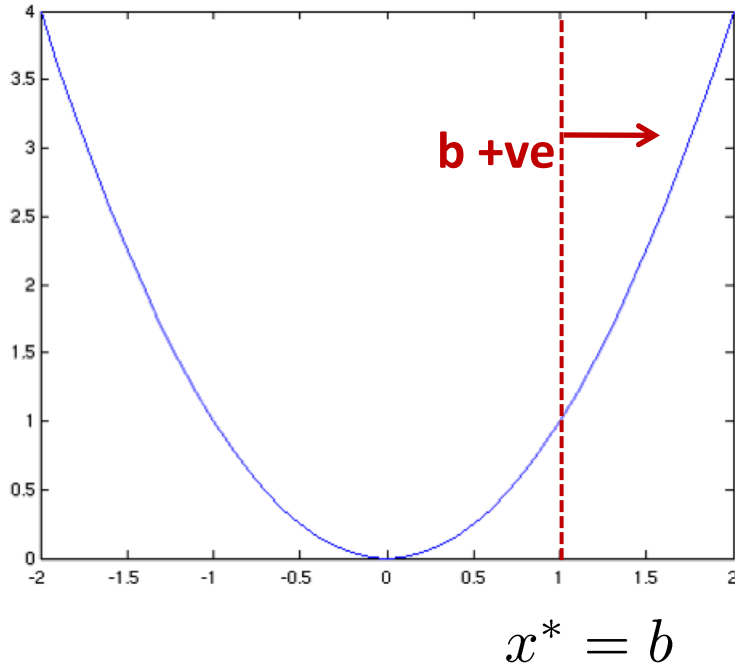
$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq 1 \end{aligned}$$



$$x^* = 1$$

Constraint active
and tight

Constrained Optimization – Dual Problem



$\alpha = 0$ constraint is inactive
 $\alpha > 0$ constraint is active

Primal problem:

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

Moving the constraint to objective function
Lagrangian:

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

Dual problem:

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

Connection between Primal and Dual

$$\text{Dual problem: } d^* = \max_{\alpha} d(\alpha) = \max_{\alpha} \min_x L(x, \alpha) \\ \text{s.t. } \alpha \geq 0 \quad \text{s.t. } \alpha \geq 0$$

Notice that

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

$$\text{Why? } L(x, \alpha) = x^2 - \alpha(x - b)$$

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

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 α (i.e. $\alpha \geq 0$), notice that

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

Since this holds for all feasible x , in particular it holds for x^* achieving the min of p^* , hence $d(\alpha) \leq p^*$ for all feasible $\alpha \geq 0$.

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 α^* (the α 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 α^* and x^* :

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

We use the first one to relate x^* and α^* . 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.

Solving the dual

Solving:

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

Find the dual: Optimization over x is unconstrained.

$$\begin{aligned} \frac{\partial L}{\partial x} = 2x - \alpha = 0 & \Rightarrow x^* = \frac{\alpha}{2} & L(x^*, \alpha) &= \frac{\alpha^2}{4} - \alpha \left(\frac{\alpha}{2} - b \right) \\ & & &= -\frac{\alpha^2}{4} + b\alpha \end{aligned}$$

Solve: Now need to maximize $L(x^*, \alpha)$ over $\alpha \geq 0$

Solve unconstrained problem to get α' and then take $\max(\alpha', 0)$

$$\frac{\partial}{\partial \alpha} L(x^*, \alpha) = -\frac{\alpha}{2} + b \Rightarrow \alpha' = 2b$$

$$\Rightarrow \alpha^* = \max(2b, 0) \quad \Rightarrow x^* = \frac{\alpha^*}{2} = \max(b, 0)$$

$\alpha = 0$ constraint is inactive, $\alpha > 0$ constraint is active and tight 10

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 _{\mathbf{w}, b} $\frac{1}{2}\mathbf{w} \cdot \mathbf{w}$
 $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \forall j$

w - weights on features (d-dim problem)

- Dual problem (derivation):

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

α - weights on training pts (n-dim problem)

Dual SVM – linearly separable case

- Dual problem:

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

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

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

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

Dual SVM – linearly separable case

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

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

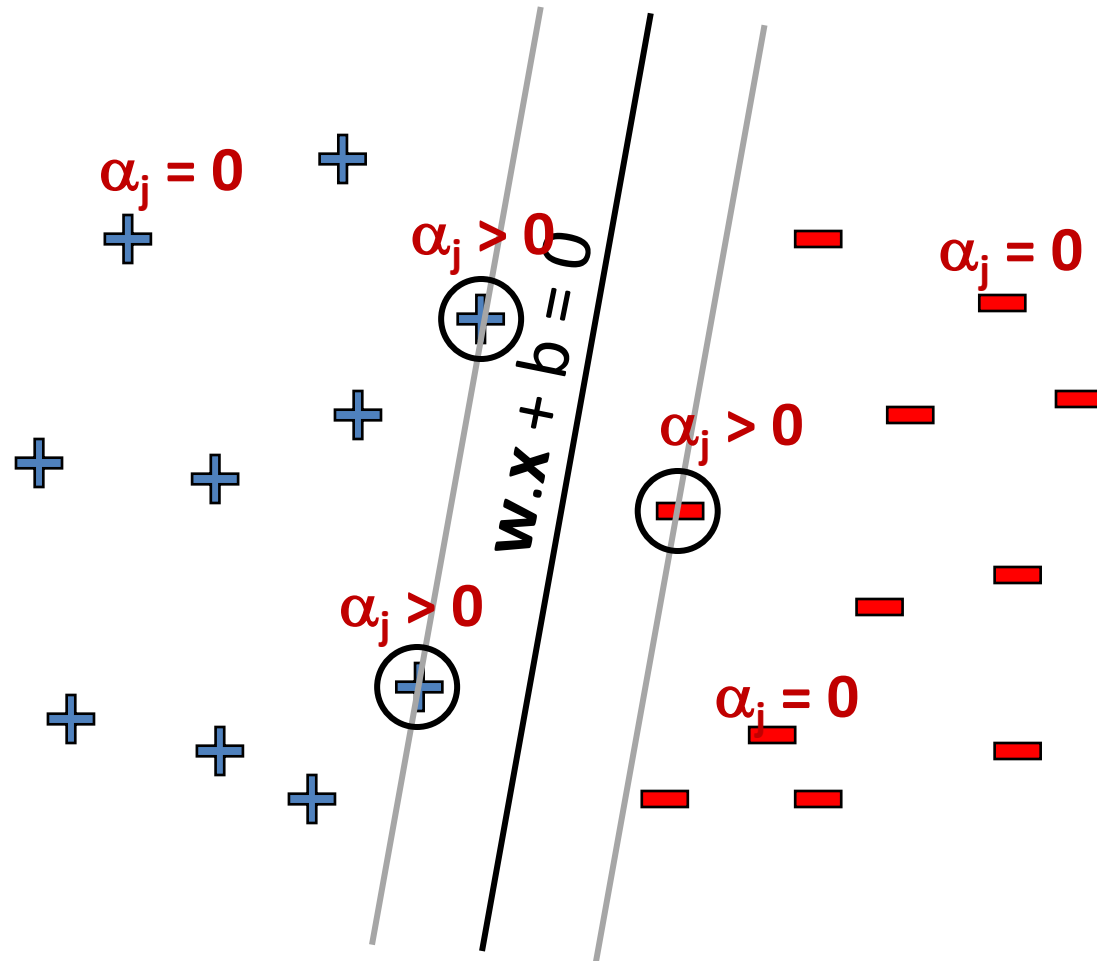
Dual problem is also QP

Solution gives α_j s \longrightarrow

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

What about b?

Dual SVM: Sparsity of dual solution



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

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