

# Support Vector Machines (SVMs)

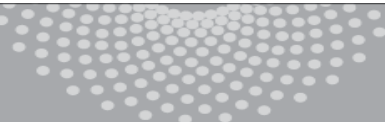
Aarti Singh & Geoff Gordon

Machine Learning 10-701

Mar 22, 2021



**MACHINE LEARNING** DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science

# Discriminative Classifiers

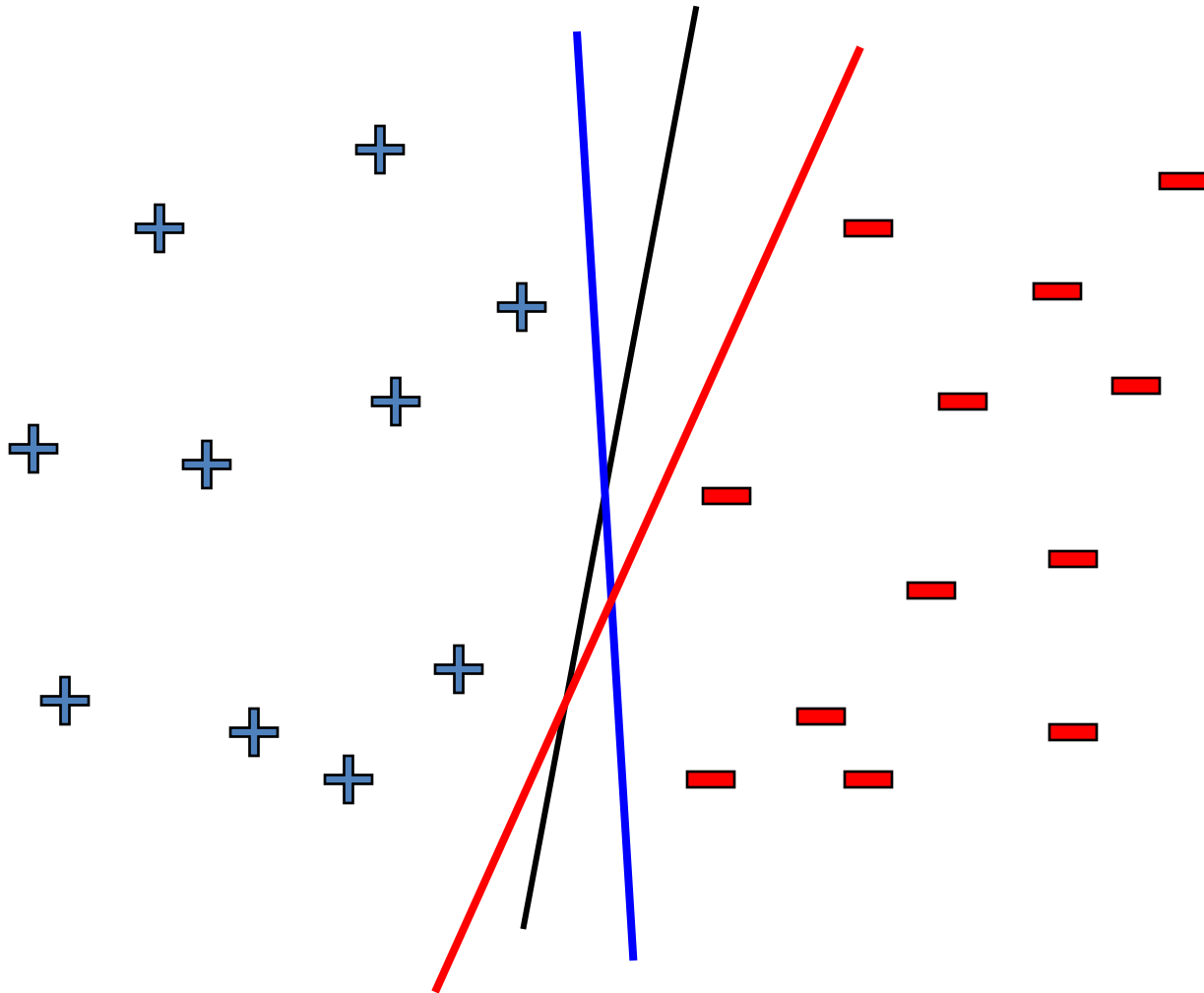
Optimal Classifier:

$$\begin{aligned} f^*(x) &= \arg \max_{Y=y} P(Y = y | X = x) \\ &= \arg \max_{Y=y} P(X = x | Y = y) P(Y = y) \end{aligned}$$

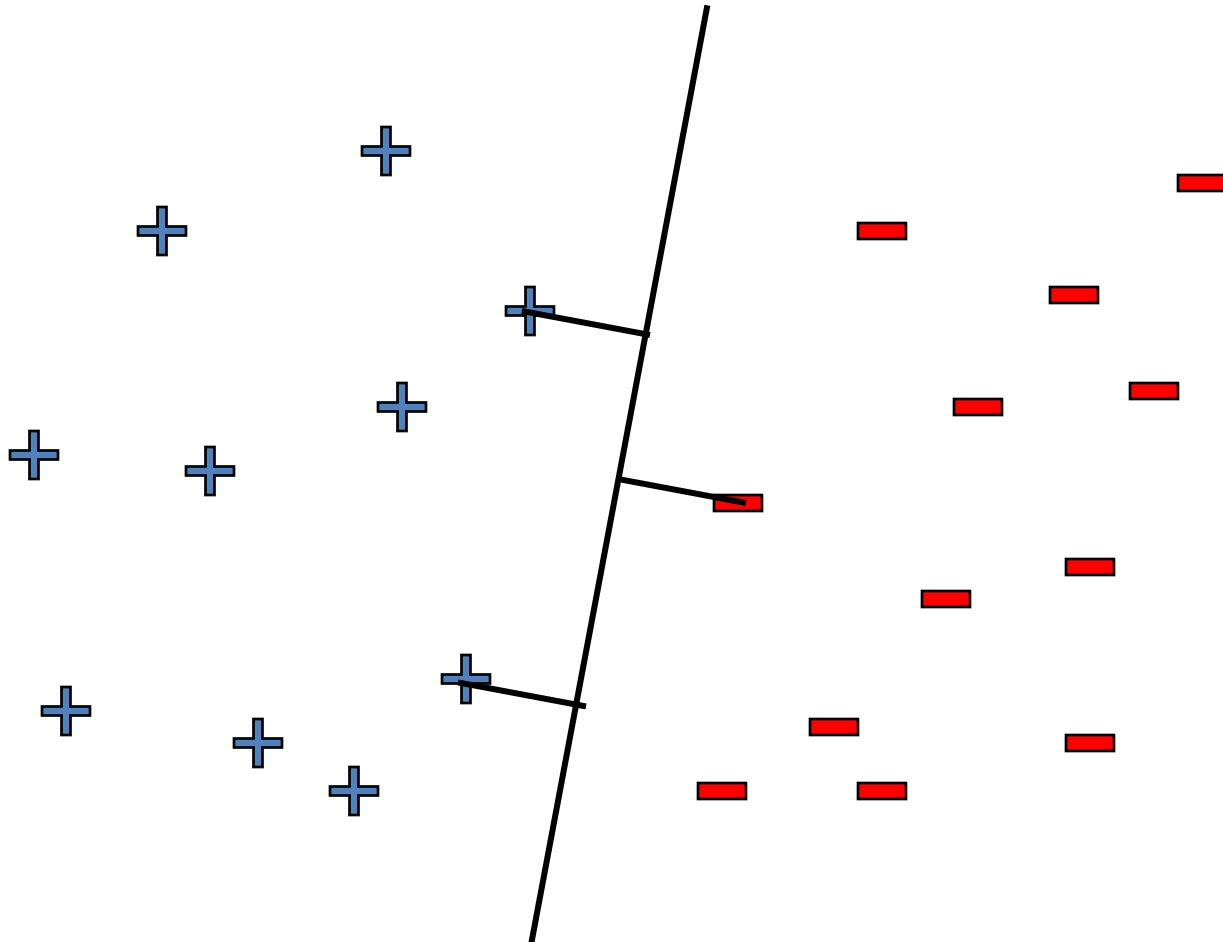
Why not learn  $P(Y|X)$  directly? Or better yet, why not learn the decision boundary directly?

- Assume some functional form for  $P(Y|X)$  (e.g. Logistic Regression) or for the decision boundary (e.g. Neural nets, SVMs - today)
- Estimate parameters of functional form directly from training data

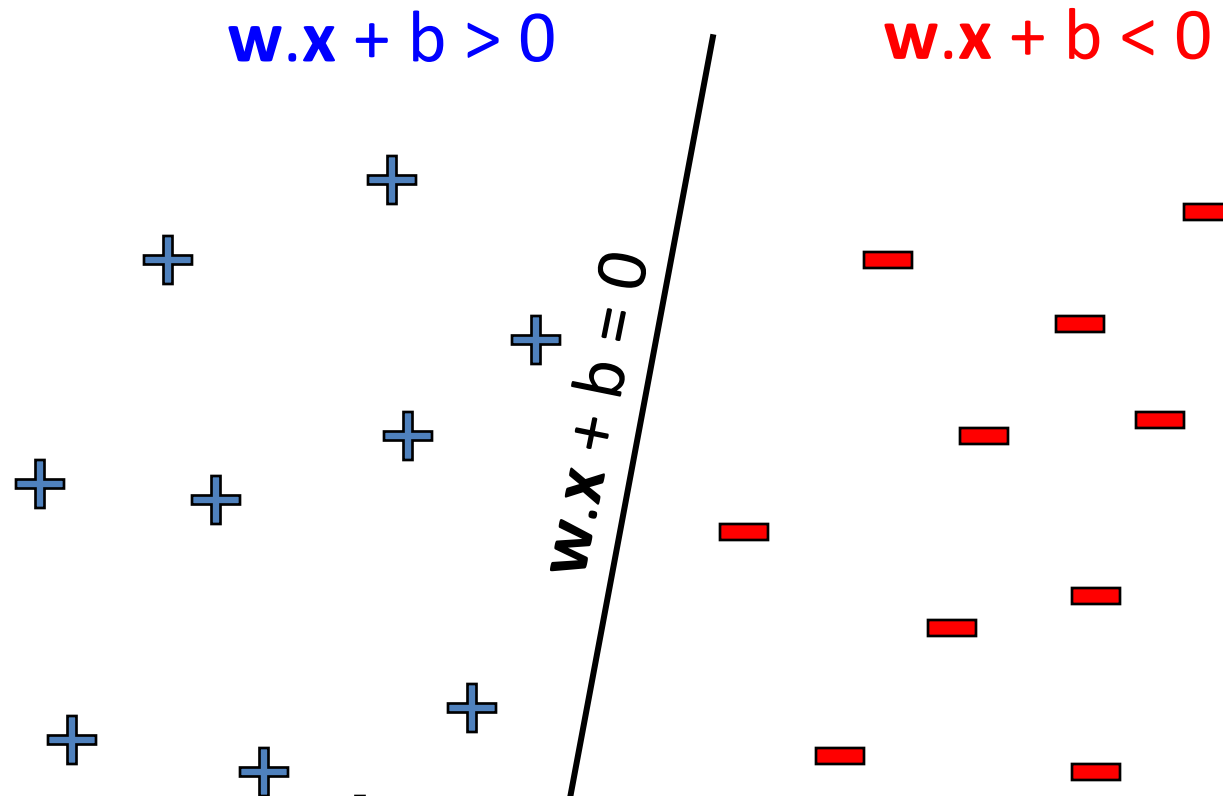
# Linear classifiers – which line is better?



# Pick the one with the largest margin!



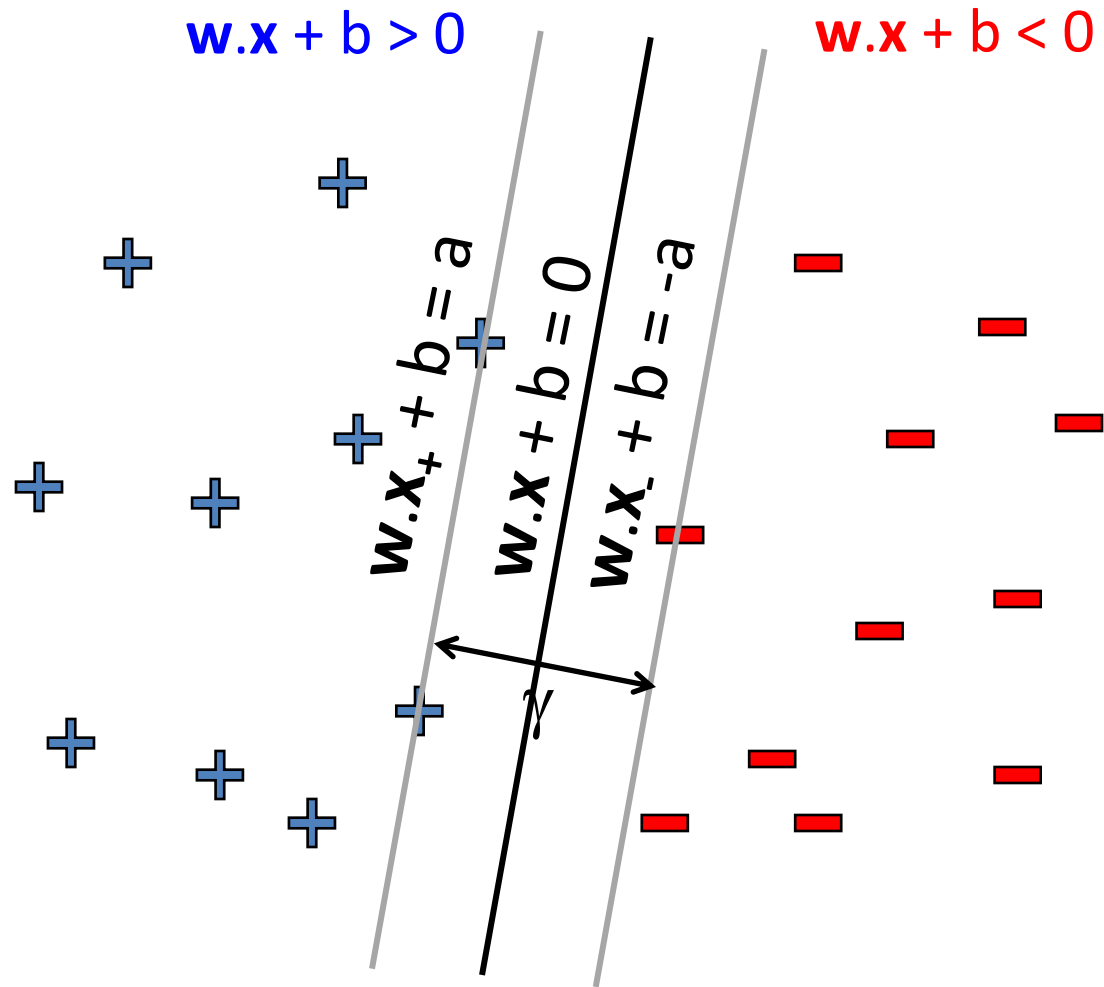
# Parameterizing the decision boundary



$y_j \in \{-1, +1\}$  — class

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

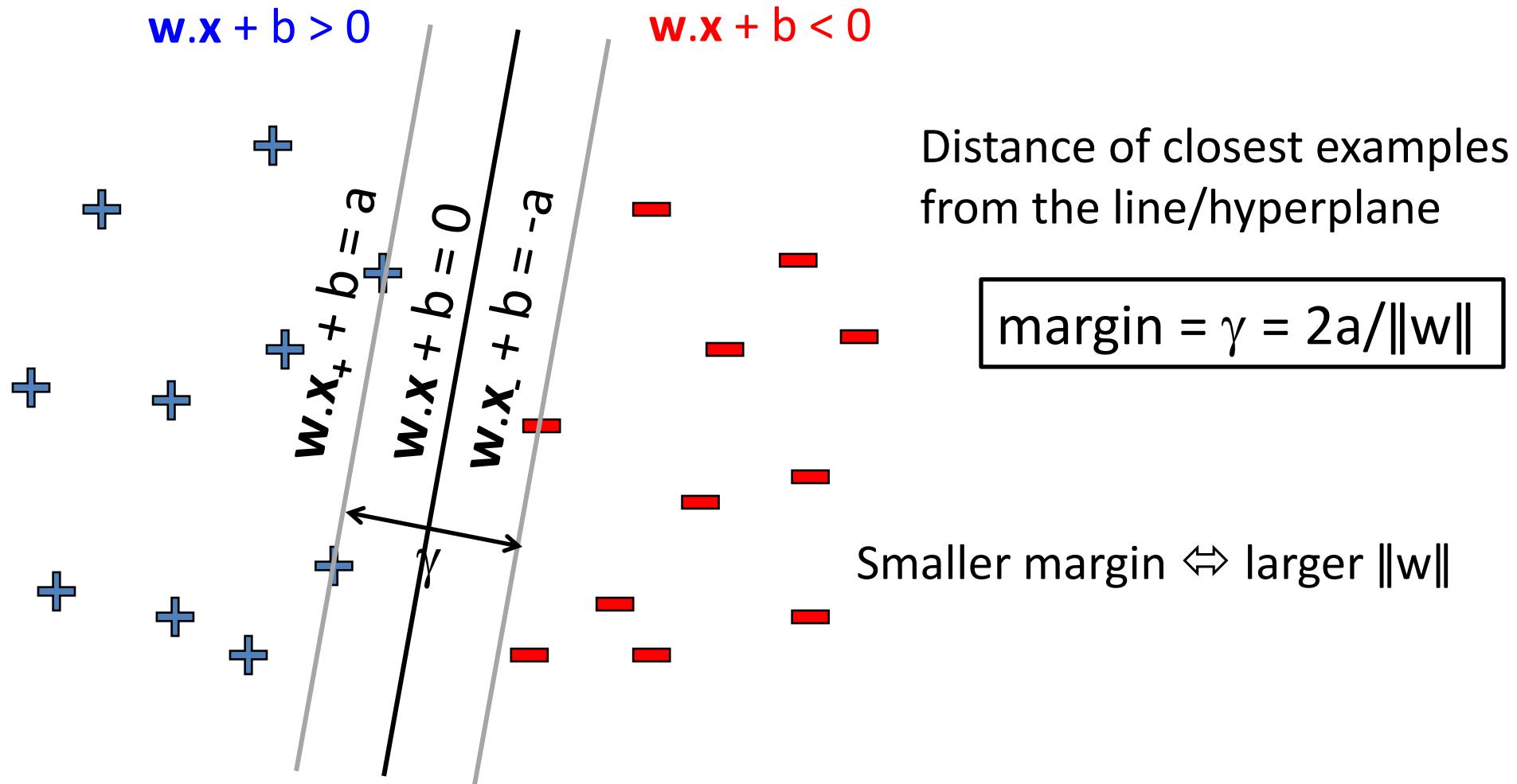
# Maximizing the margin



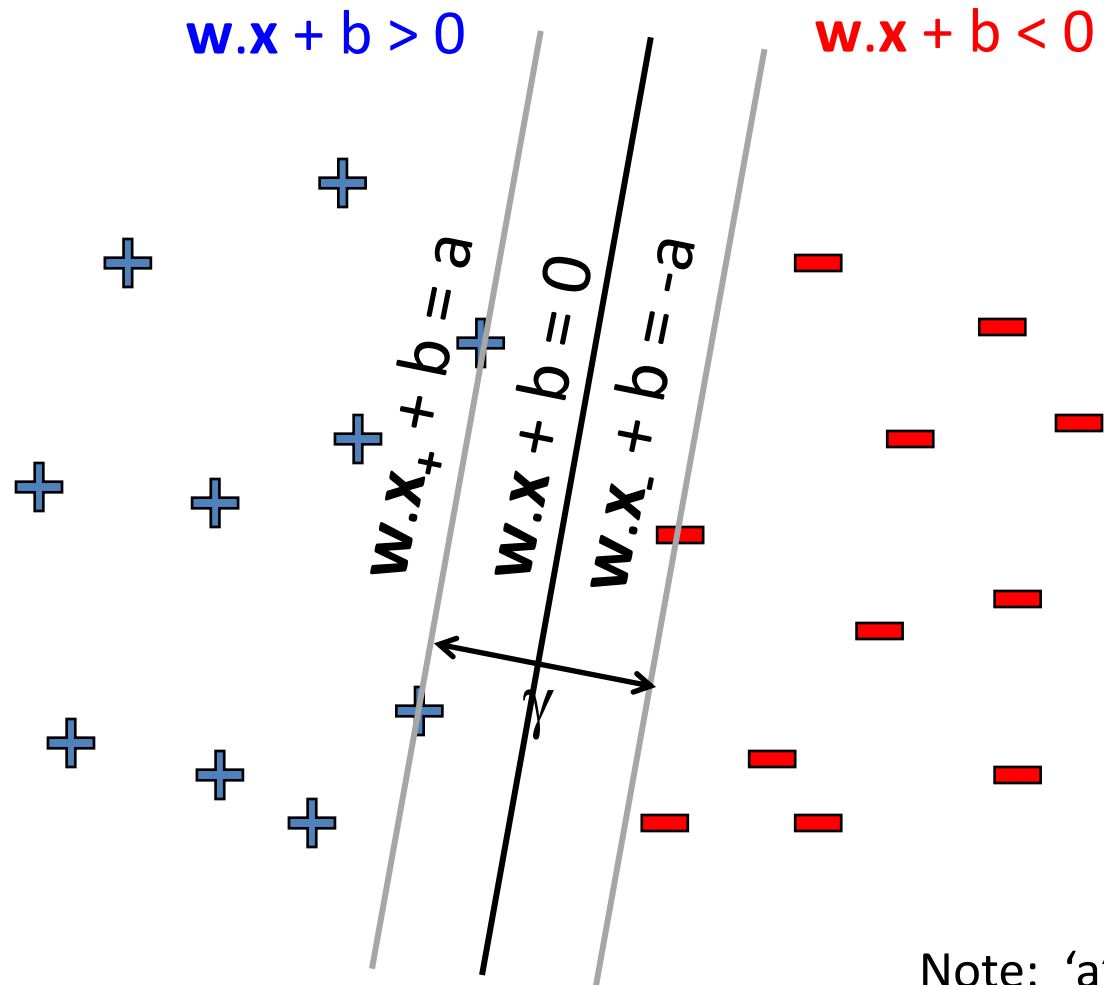
Distance of closest examples from the line/hyperplane

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

# 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

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

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

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

$\gamma$

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

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

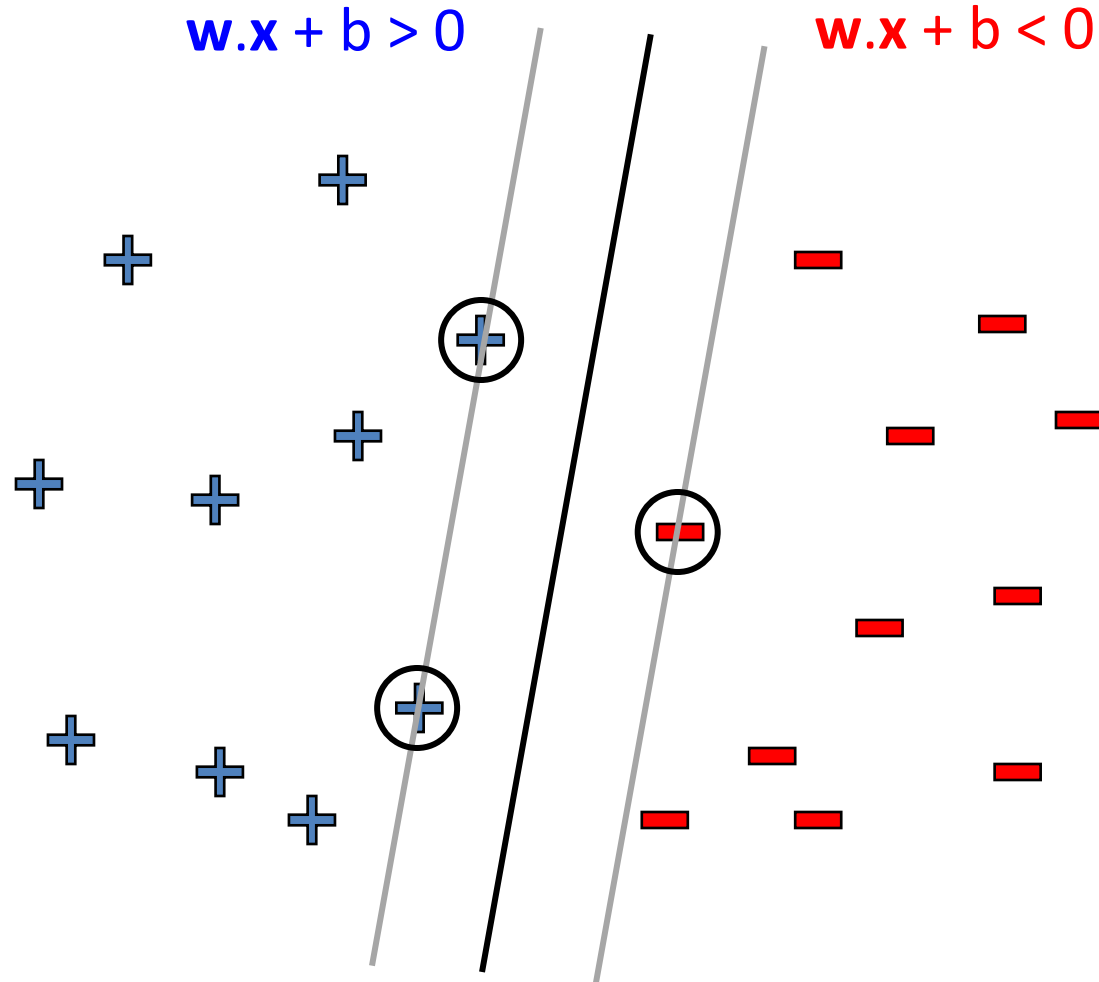
Solve efficiently by quadratic programming (QP)

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

# Support Vectors

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

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



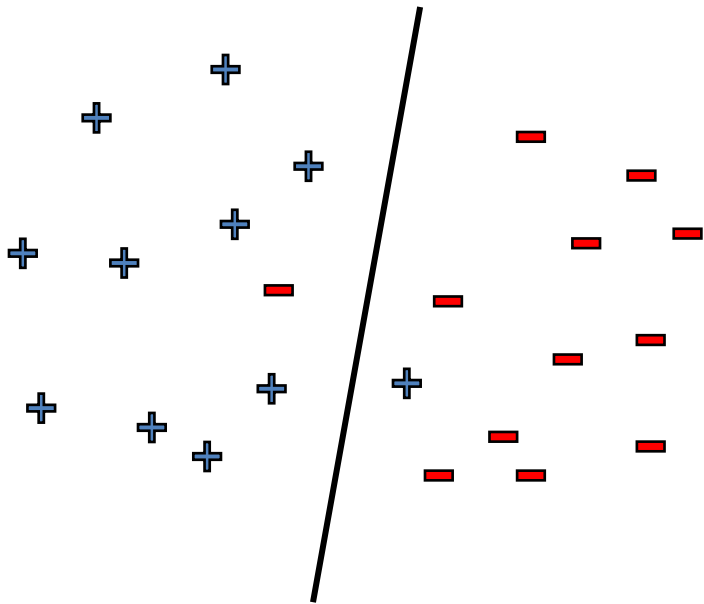
Linear hyperplane defined by  
“support vectors”

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 not linearly separable?



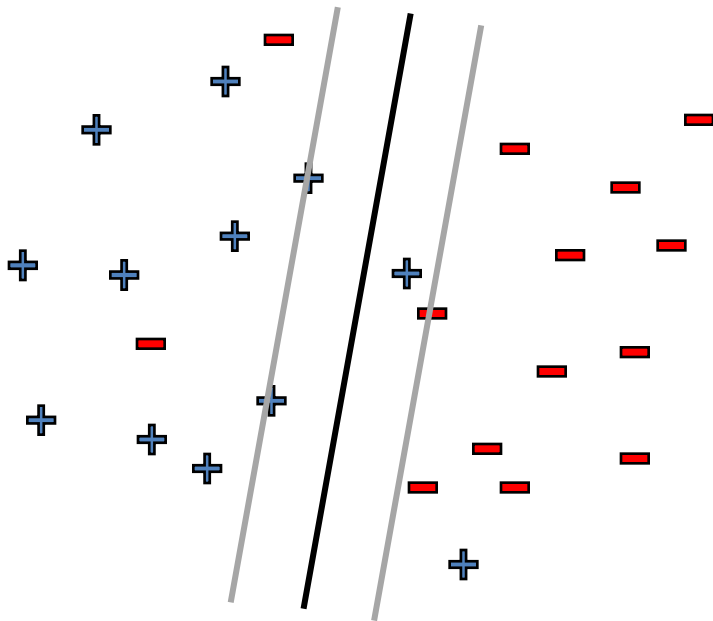
Use features of features  
of features of features....

$$x_1^2, x_2^2, x_1x_2, \dots, \exp(x_1)$$

But run risk of overfitting!

# What if data is still not linearly separable?

Allow “error” in classification



Smaller margin  $\Leftrightarrow$  larger  $\|w\|$

$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \cdot \text{\#mistakes} \\ \text{s.t.} \quad & (w \cdot x_j + b) y_j \geq 1 \quad \forall j \end{aligned}$$

Maximize margin and minimize  
# mistakes on training data

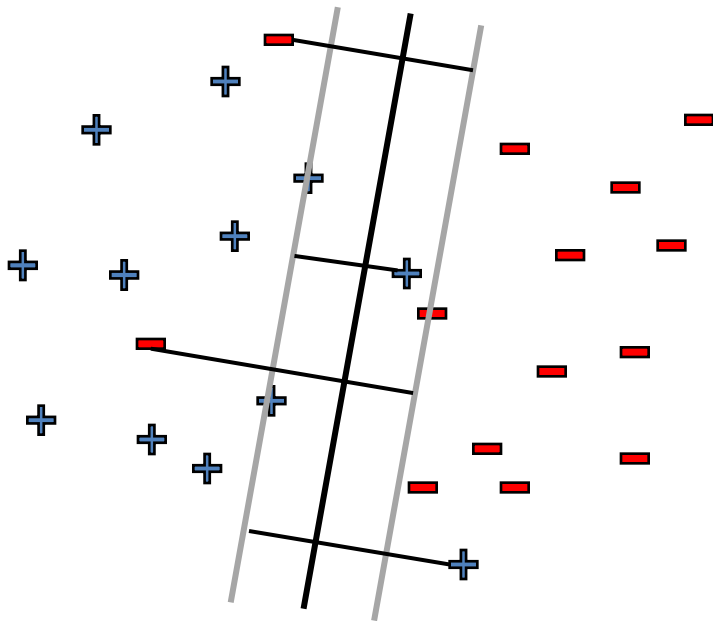
C - tradeoff parameter

Not QP ☹️

0/1 loss (doesn't distinguish between  
near miss and bad mistake)

# What if data is still not linearly separable?

Allow “error” in classification



**Soft margin approach**

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

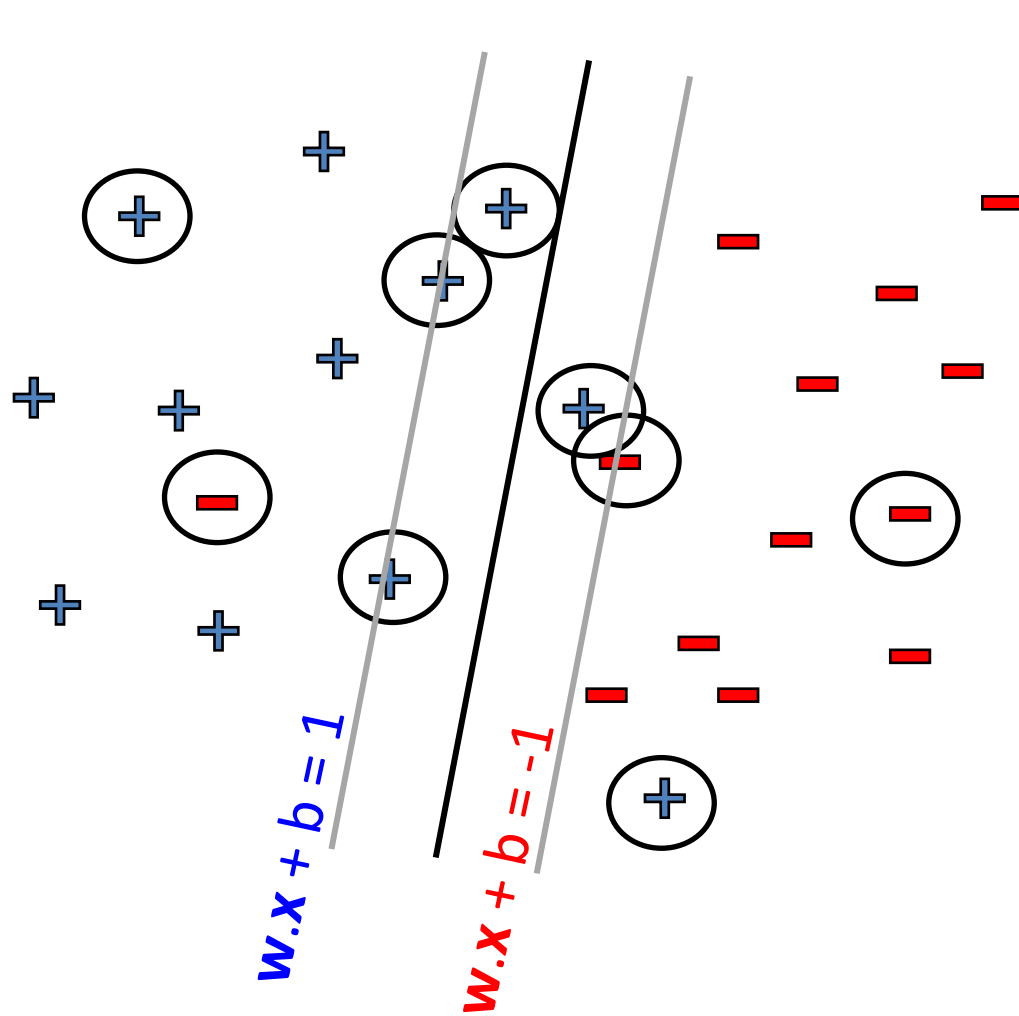
$\xi_j$  - “slack” variables  
= (>1 if  $x_j$  misclassified)

pay linear penalty if mistake

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

Still QP 😊

# Slack variables – Hinge loss

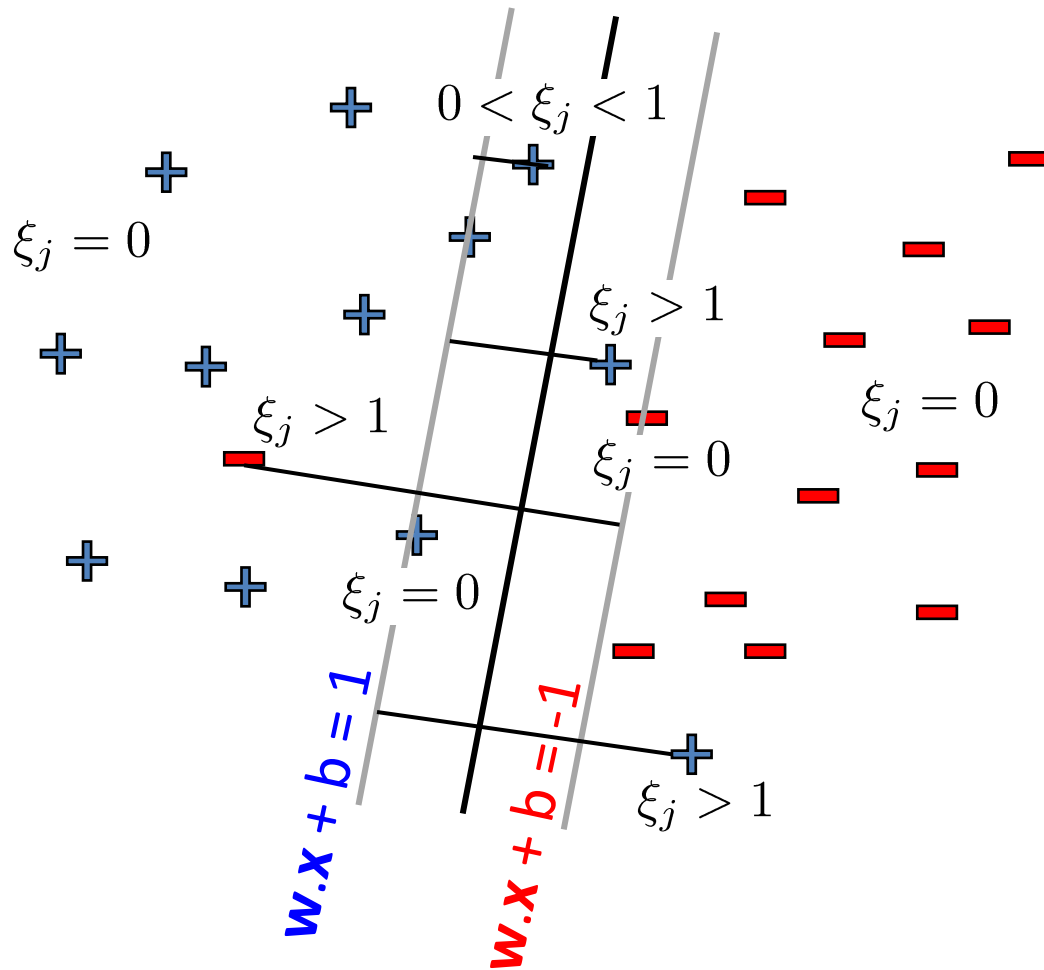


$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j$$

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

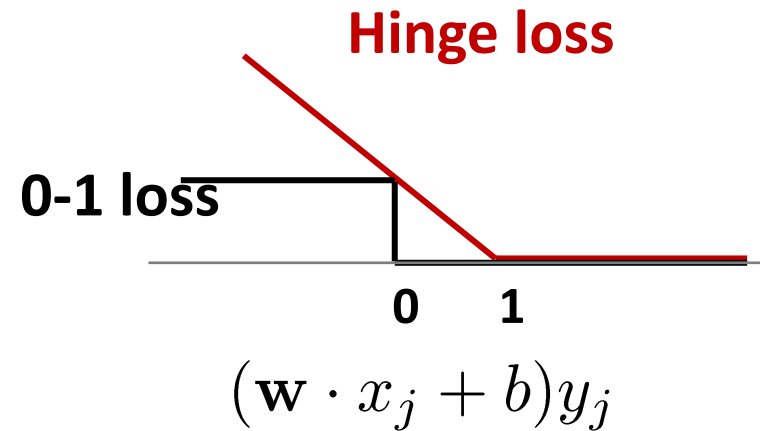
Confidence	Slack
------------	-------

# Slack variables – Hinge loss



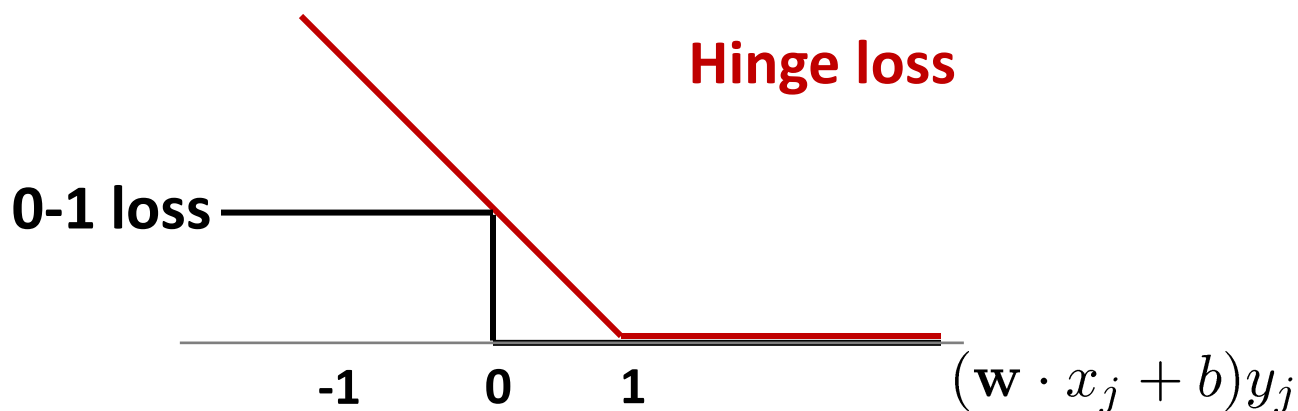
Notice that

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



# Slack variables – Hinge loss

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



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

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

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

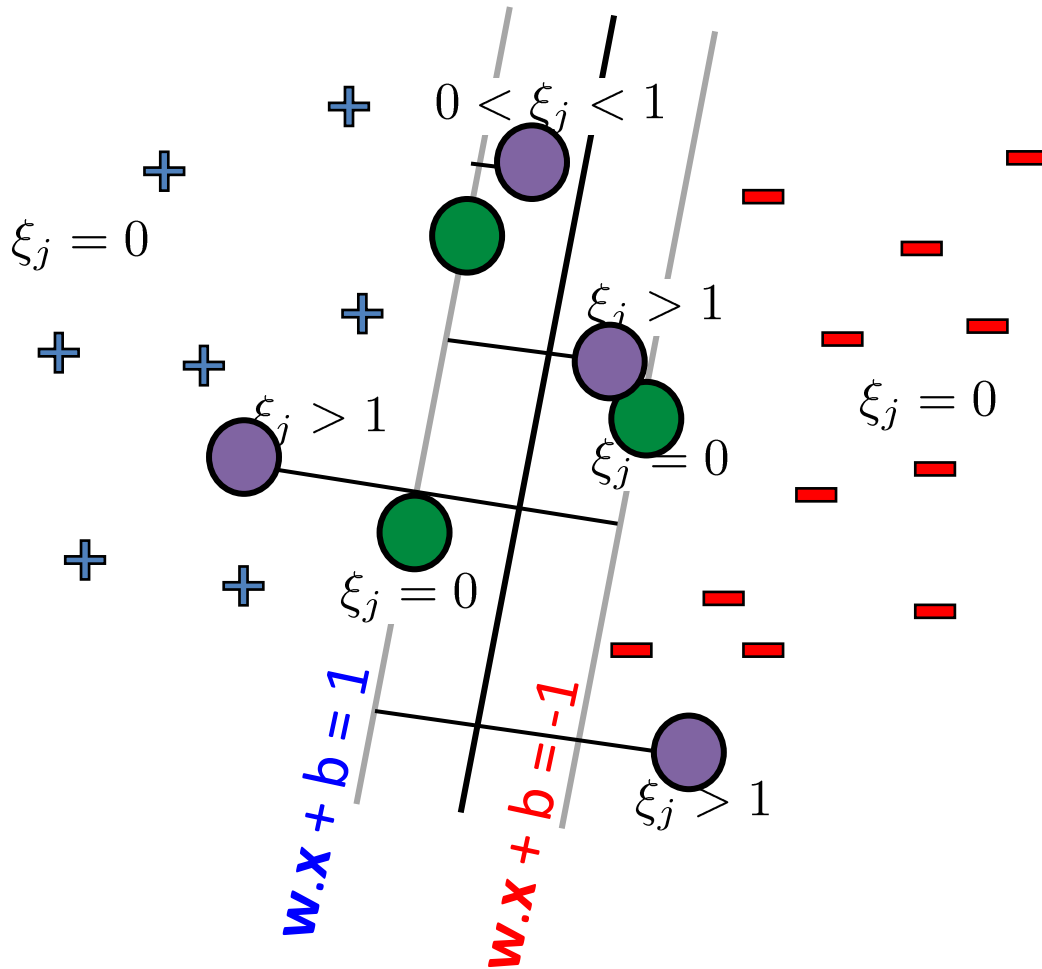


Regularized hinge loss

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



# Support Vectors



## Margin support vectors

$\xi_j = 0$ ,  $(w \cdot x_j + b) y_j = 1$   
 (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 17