

Support Vector Machines (SVMs)

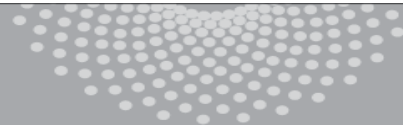
Aarti Singh

Machine Learning 10-701

Feb 1, 2023



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Discriminative Classifiers

Optimal Classifier:

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

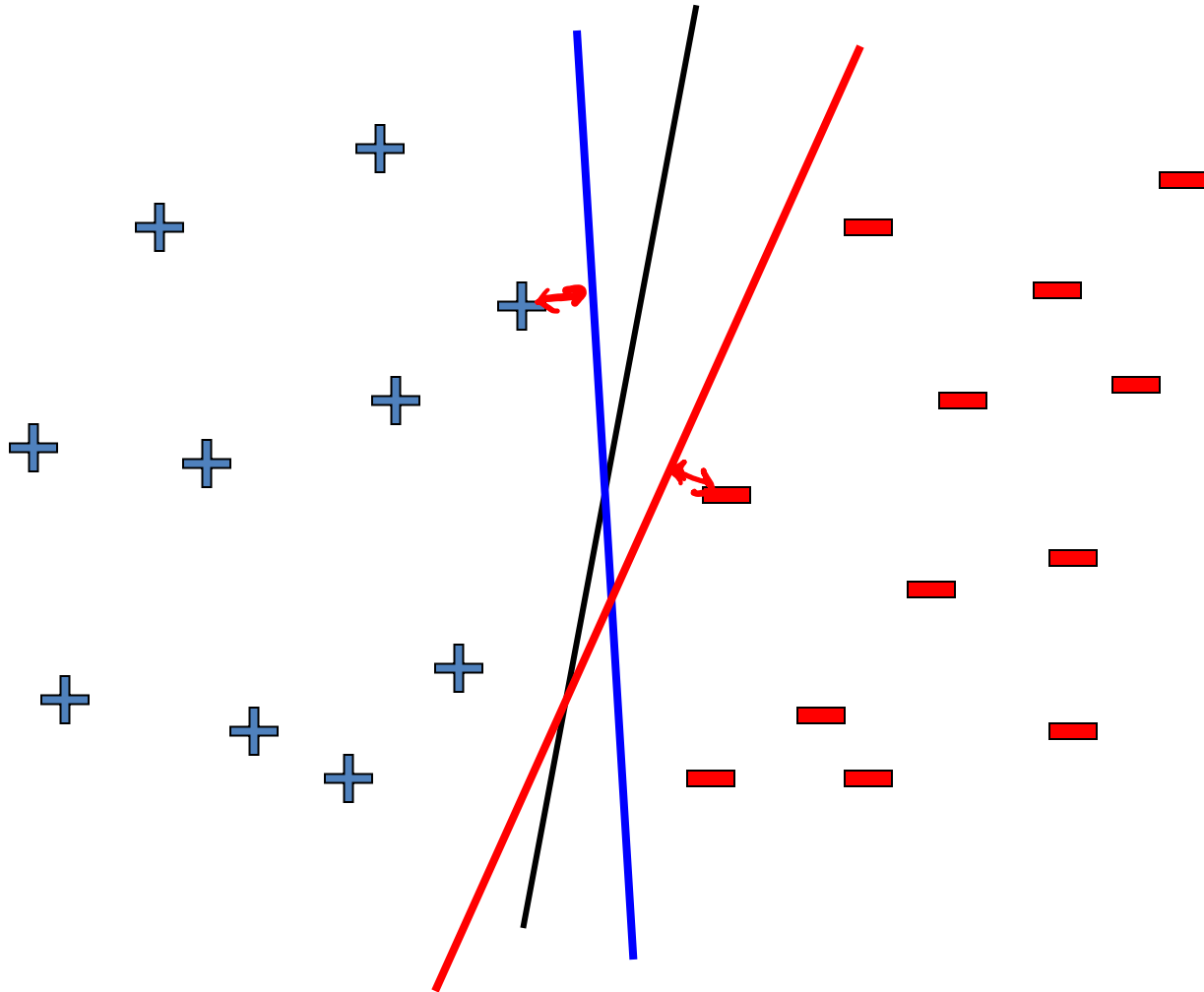
Why not learn $P(Y|X)$ directly? Or better yet, why not learn the decision boundary directly?

$$P(Y=1|X) = \frac{1}{1 + \exp(-\sum_j w_j^{(1)} x_j^{(1)})}$$

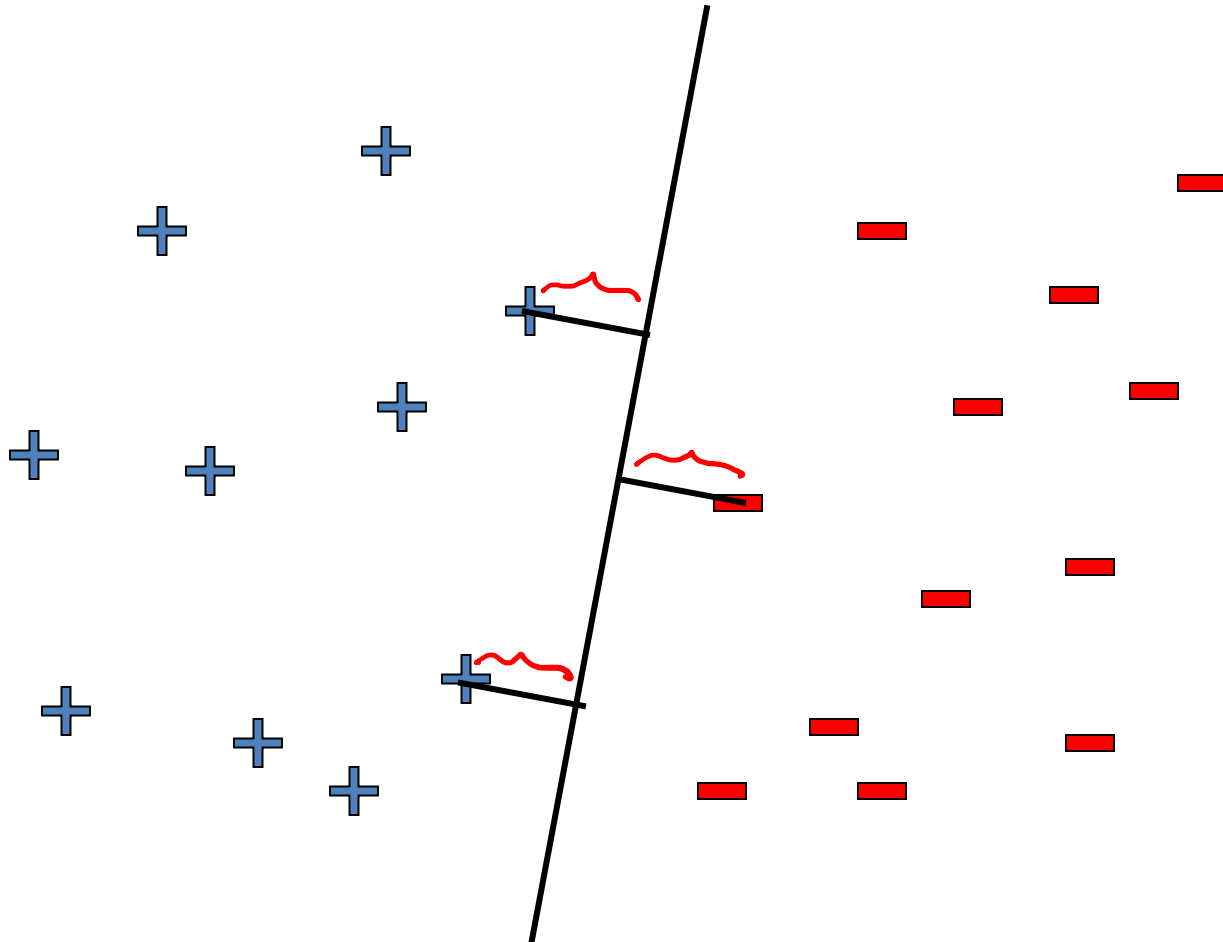
$X = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(d)} \end{bmatrix}$

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

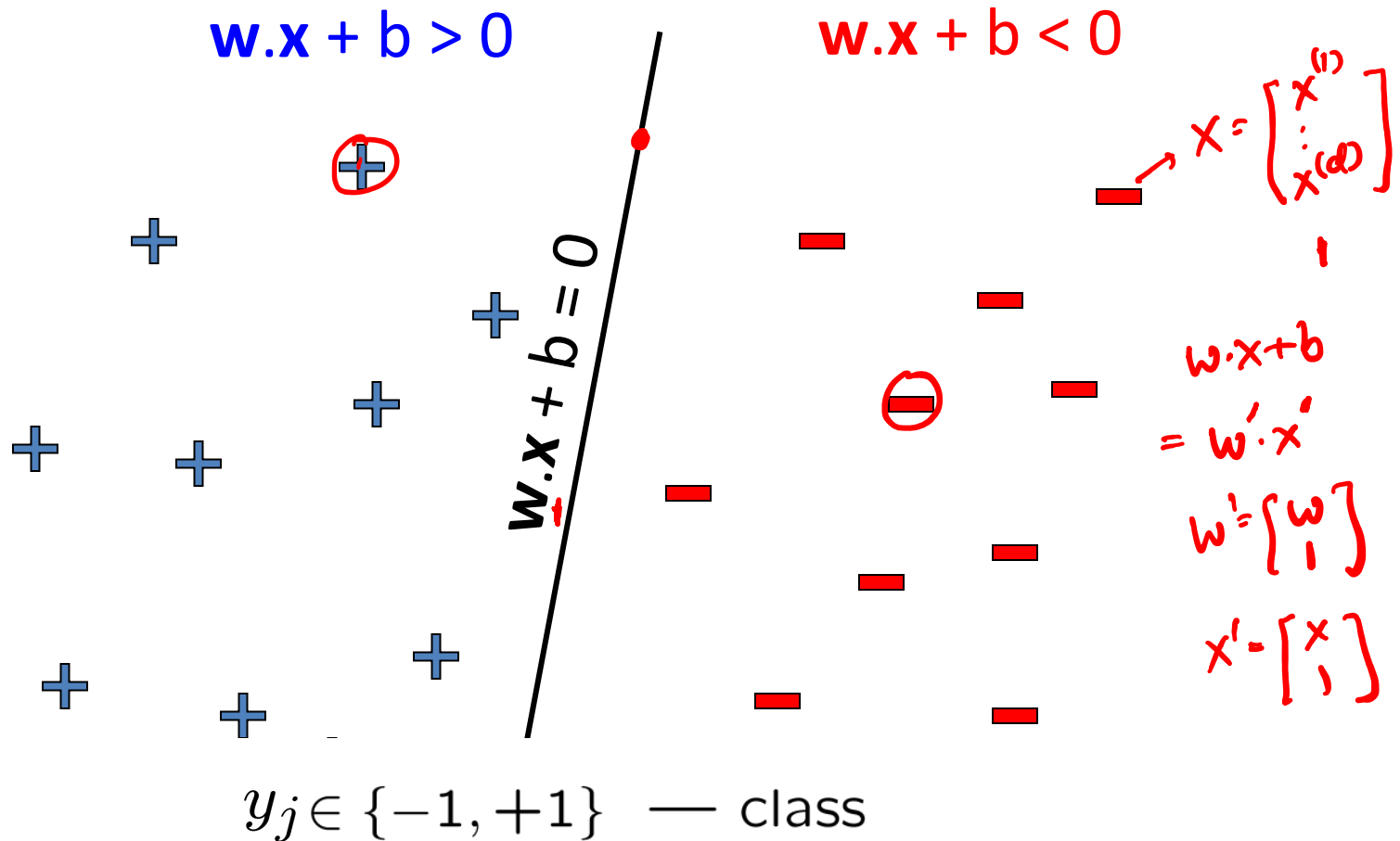
Linear classifiers – which line is better?



Pick the one with the largest margin!

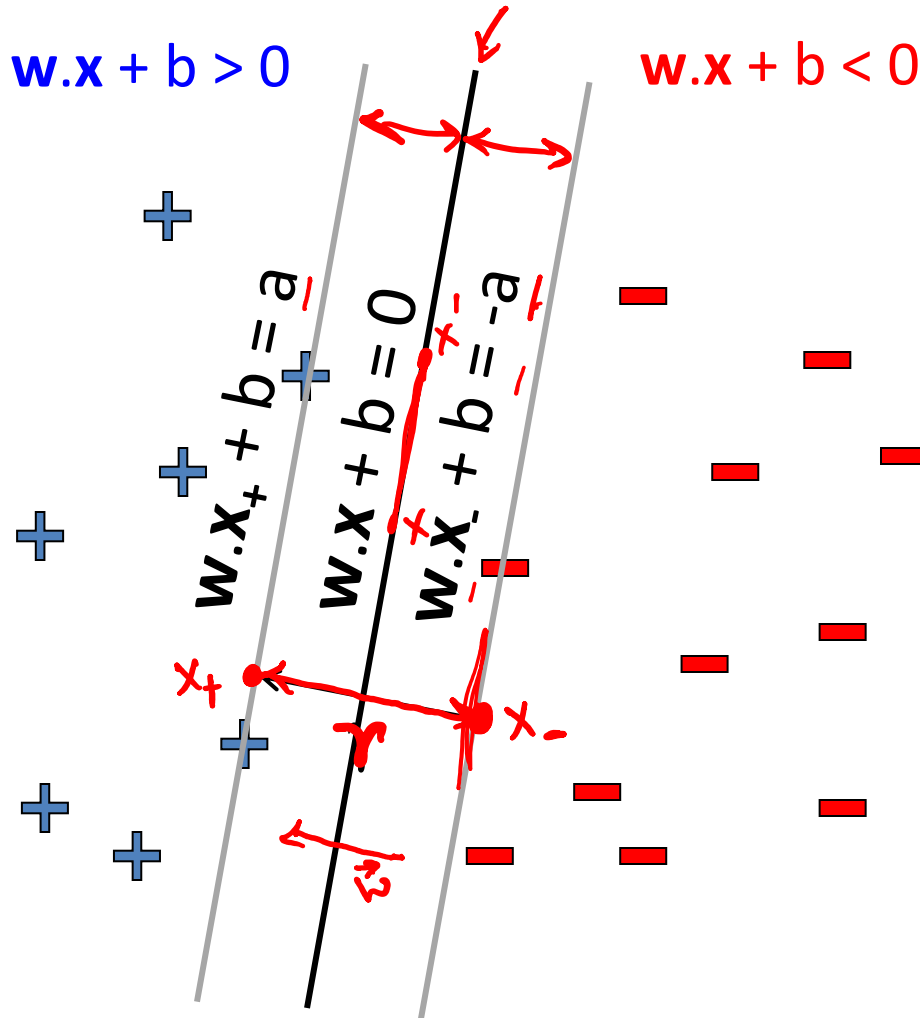


Parameterizing the decision boundary



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

Maximizing the margin



Distance of closest examples from the line/hyperplane

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

1. $w \cdot x + b = 0 \Rightarrow w \cdot (x - x') = 0$
 $w \cdot x' + b = 0$

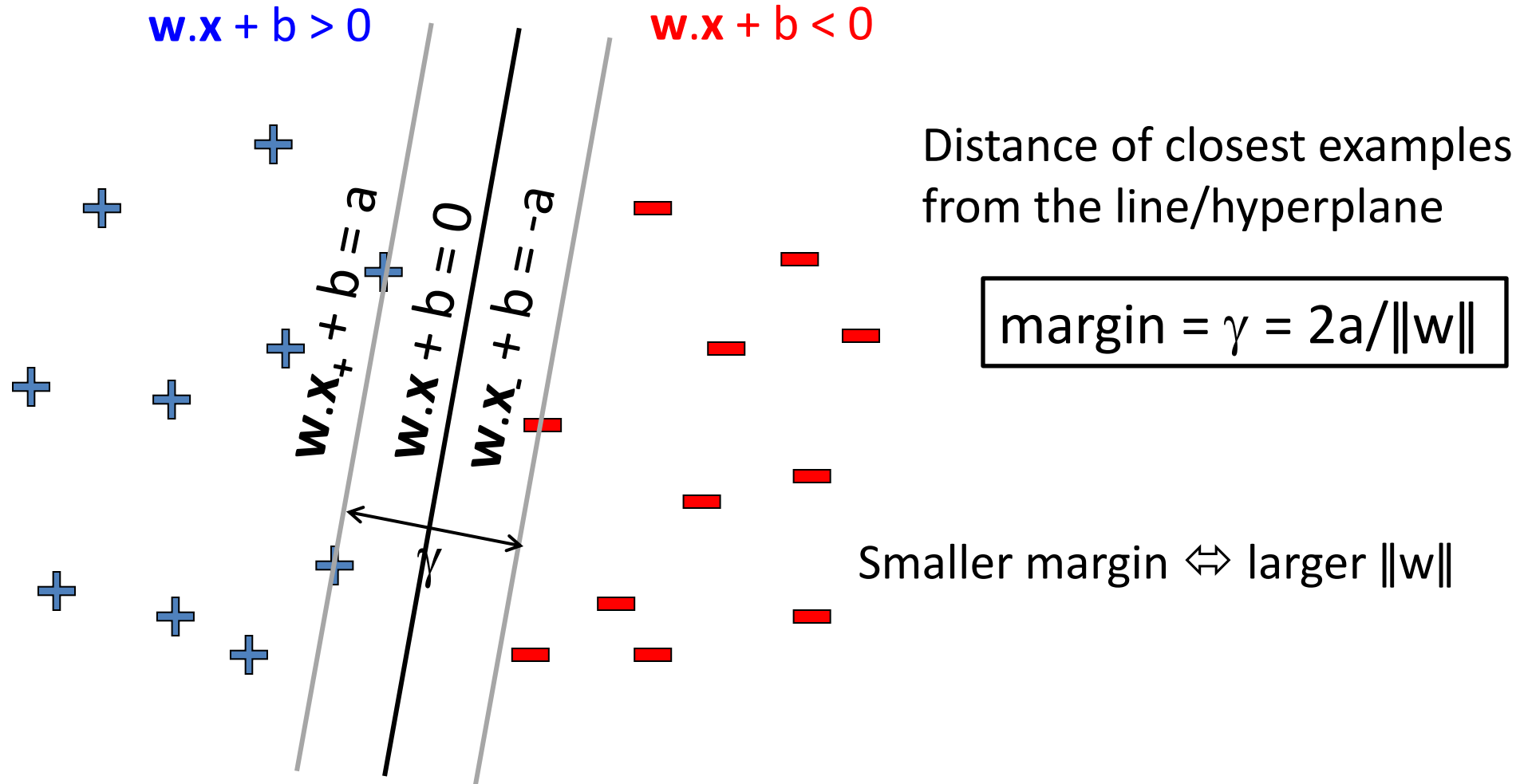
2. $x_- + \gamma \cdot \frac{w}{\|w\|} = x_+$

$w \cdot x_- + b + \gamma \cdot \frac{w \cdot w}{\|w\|} = w \cdot x_+ + b$

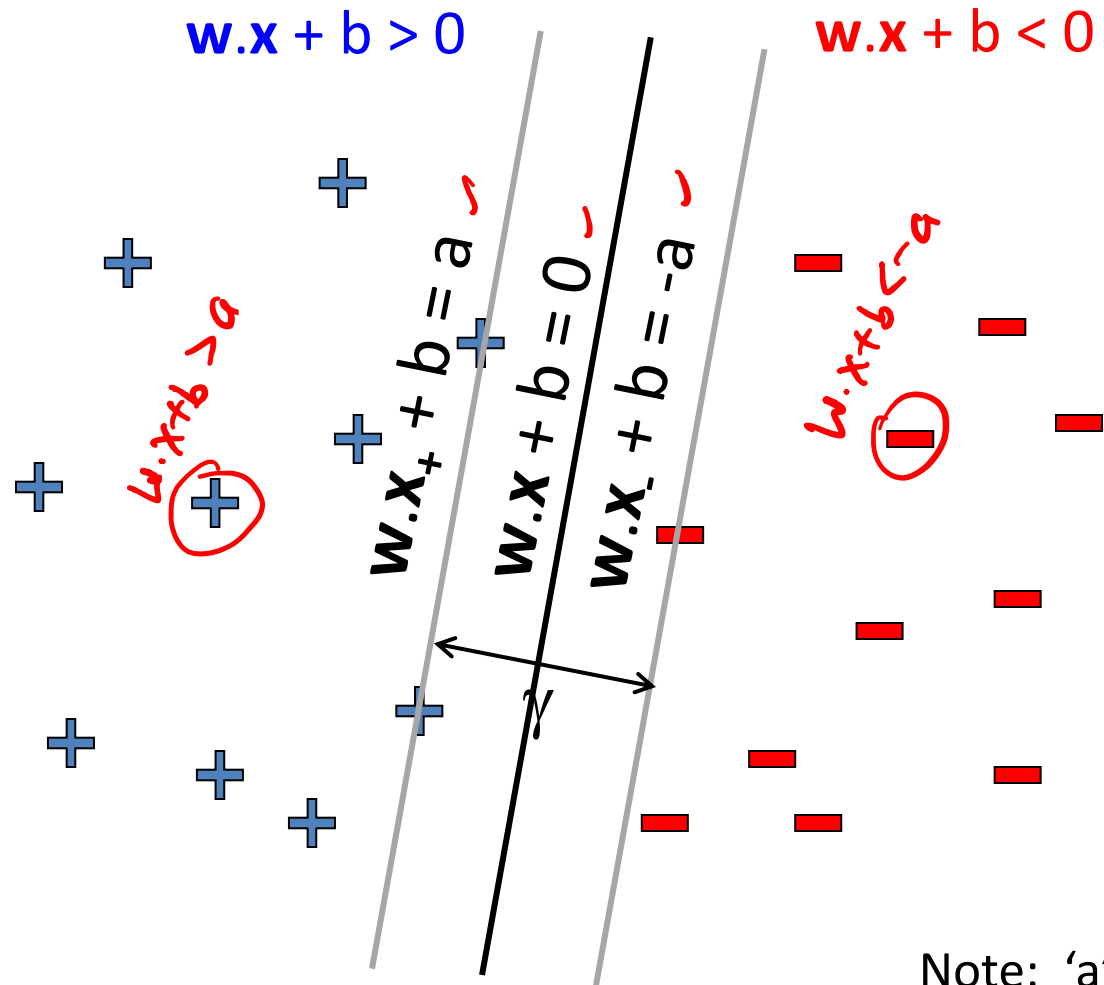
$-a + \gamma \cdot \|w\| = a$

$\gamma = \frac{2a}{\|w\|}$

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

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

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

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

γ

$$\min w \cdot w$$

$$w, b$$

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

Solve efficiently by quadratic programming (QP)

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

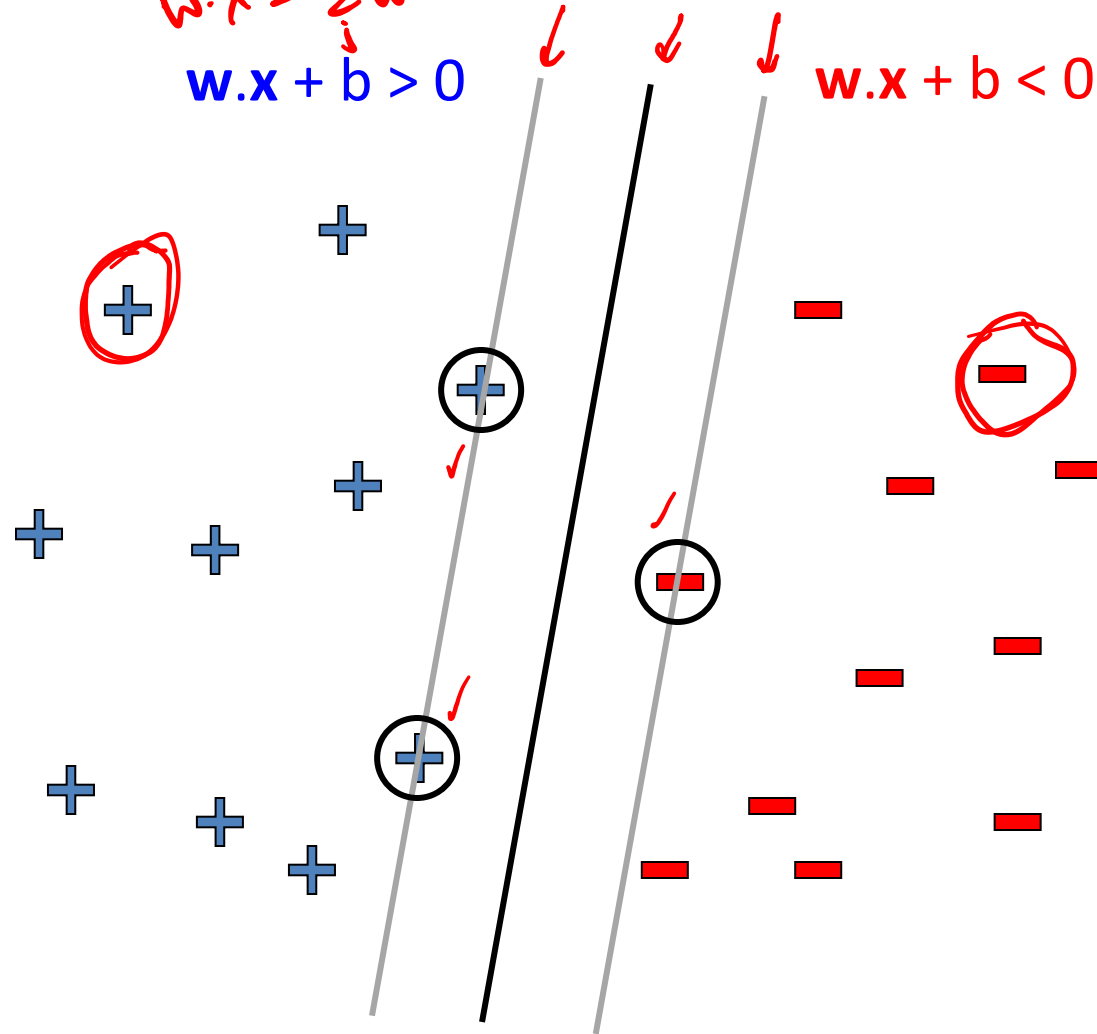
Support Vectors

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

$$w \cdot x = \sum_j w^{(j)} x^{(j)}$$

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

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



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

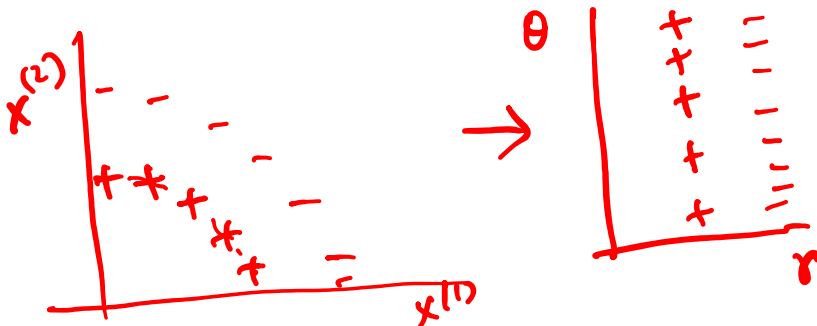
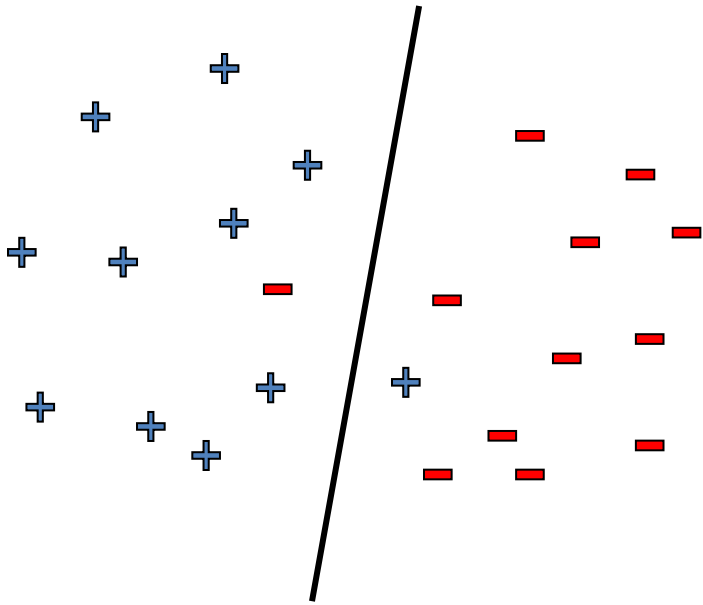
For support vectors $(w \cdot x_j + b) y_j = 1$

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!

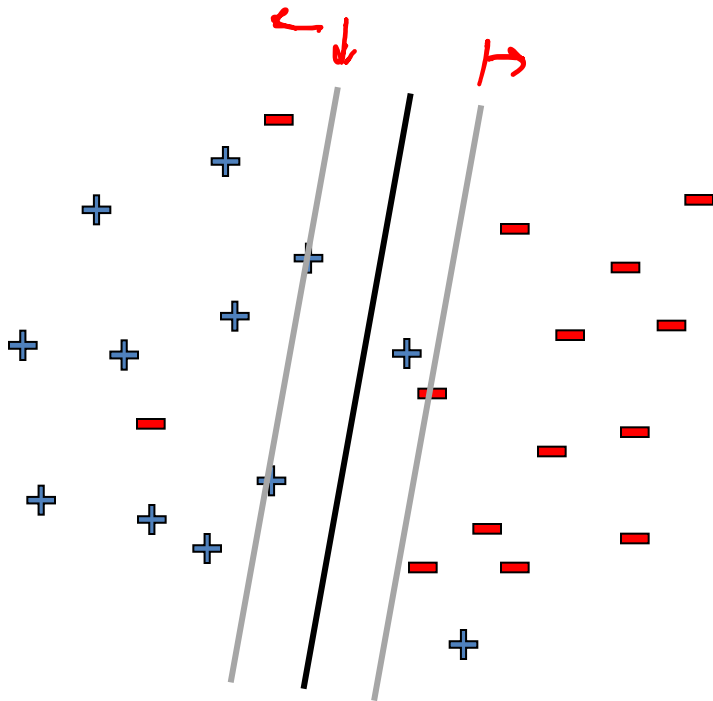


What if data is still not linearly separable?

$U \cdot w \rightarrow C \sum_j ()^{-1}$

$\sum_{i=1}^n 1_{w \cdot x_i + b \neq y_i}$

Allow "error" in classification



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

$$\min_{w,b} w \cdot w + C \# \text{mistakes}$$

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

Maximize margin and minimize # mistakes on training data

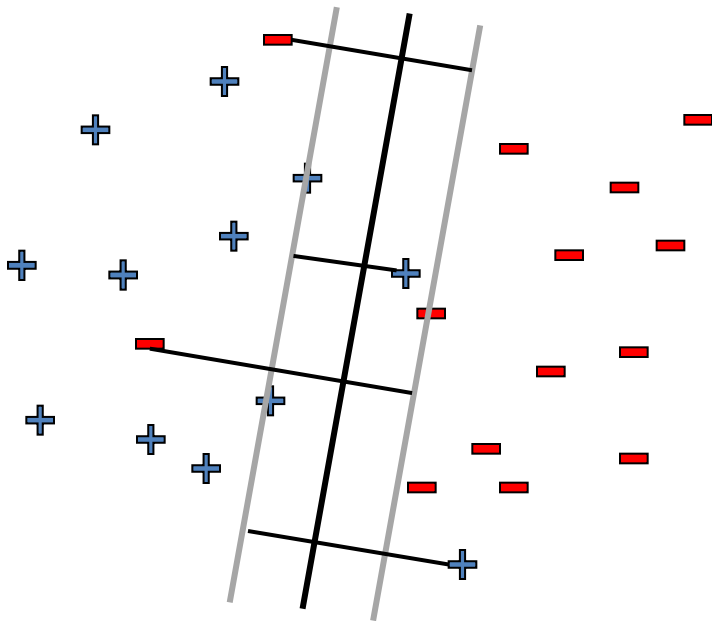
C - tradeoff parameter

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

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

ξ_j - “slack” variables
= (>1 if x_j misclassified)

pay linear penalty if mistake

C - tradeoff parameter (C = ∞
recovers hard margin SVM)

Still QP 😊

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

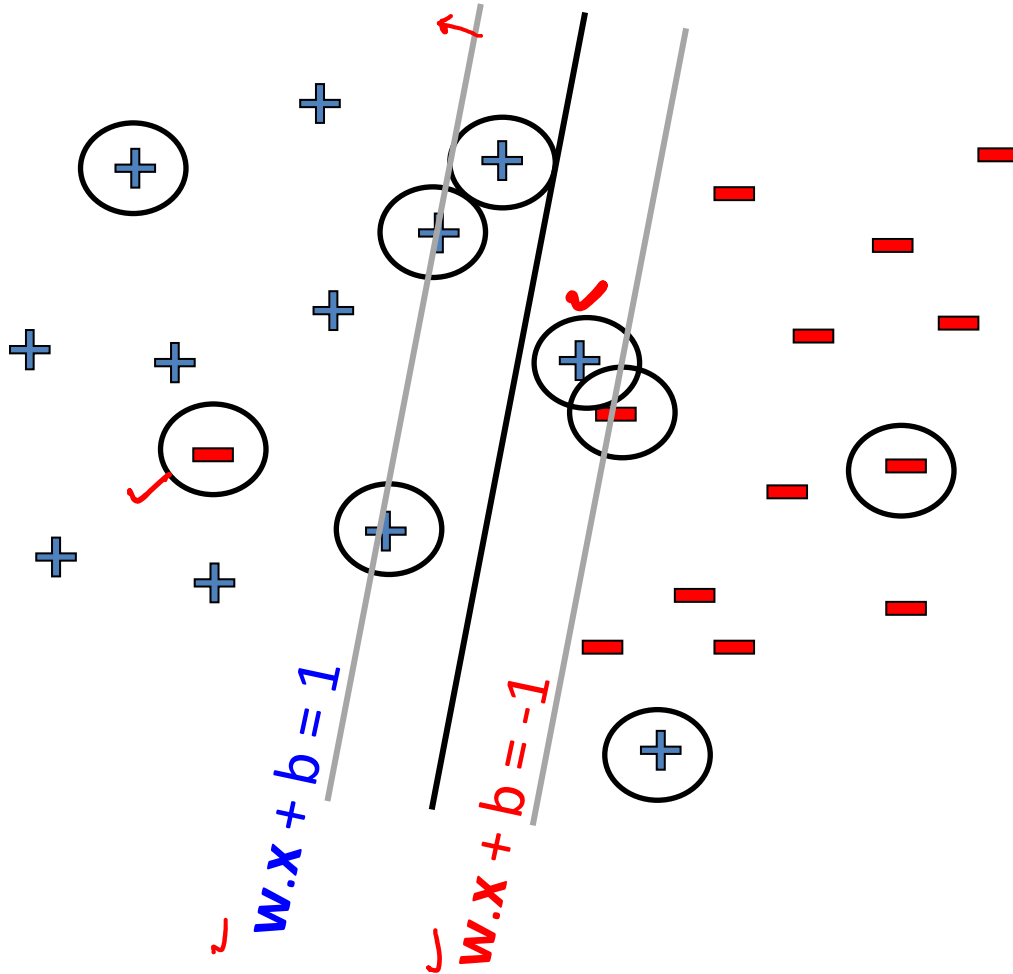
$$\text{s.t. } (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j$$

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

Variables – Hinge loss

$$1 \geq 1 - \xi_j$$

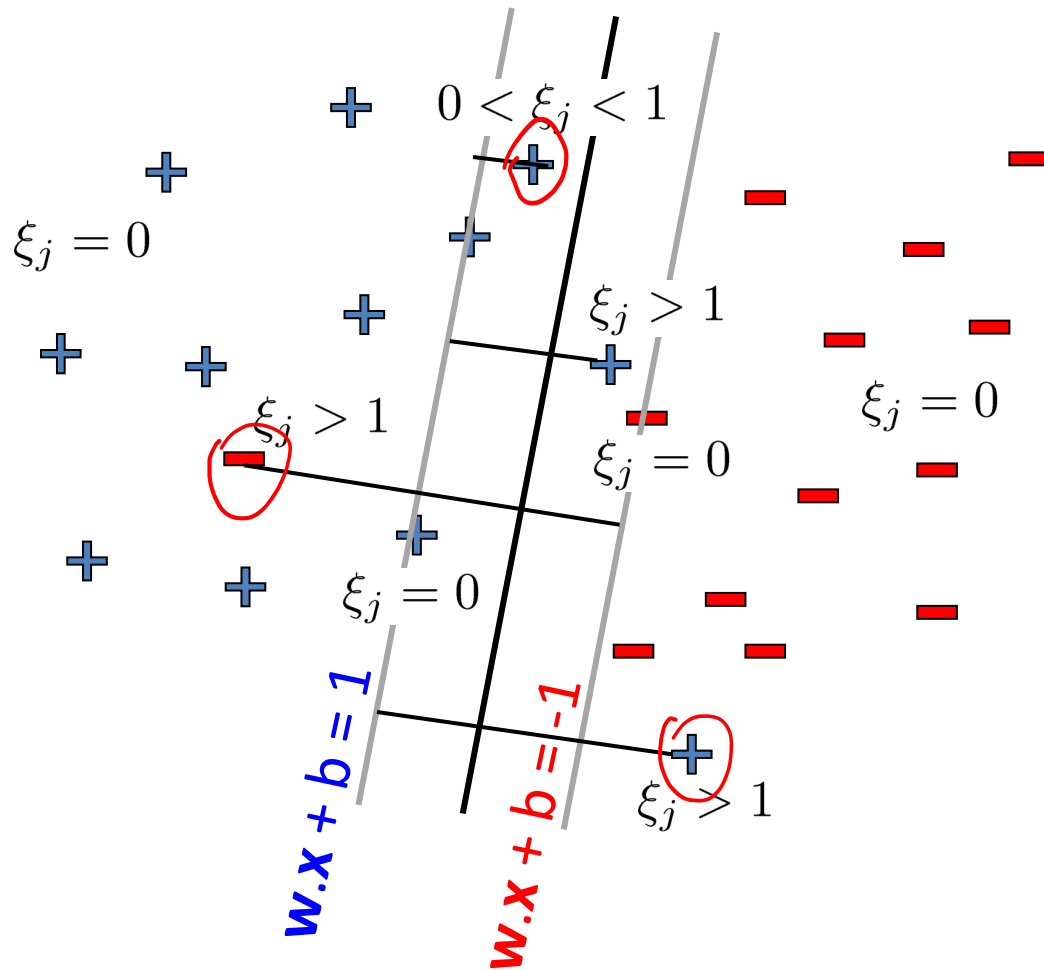
$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j$$



What is the slack ξ_j for the following points?

Confidence	Slack
$\rightarrow 1$	$\xi_j = 0$
$\rightarrow > 1$	$\xi_j = 0$
$\rightarrow 0 < < 1$	$0 < \xi_j < 1$
$-ve$	$\xi_j > 1$

Slack variables – Hinge loss



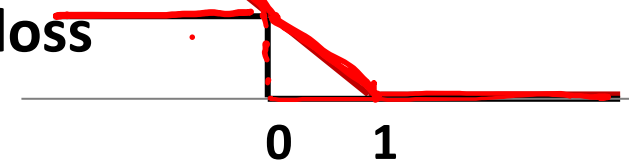
Notice that

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

$$\max(1 - (\mathbf{w} \cdot \mathbf{x}_j + b)y_j, 0)$$

Hinge loss

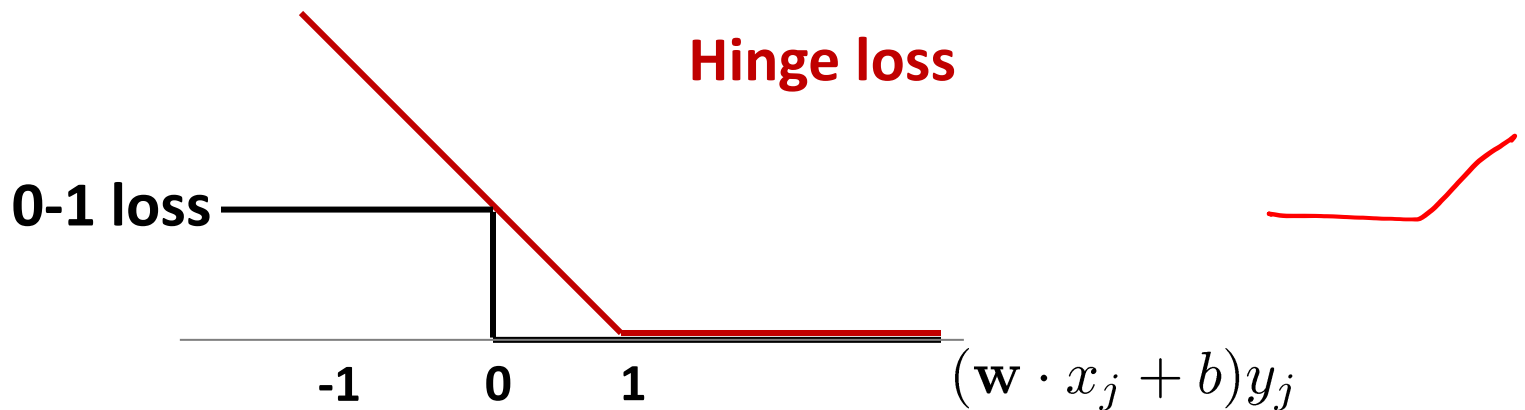
0-1 loss



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

Slack variables – Hinge loss

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



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



$$\min_{\mathbf{w}, b} \quad \underbrace{\mathbf{w} \cdot \mathbf{w}}_{\|\mathbf{w}\|^2} + C \sum_j \underbrace{(1 - (\mathbf{w} \cdot \mathbf{x}_j + b)y_j)_+}_{\text{Hinge loss}}$$

SVM vs. Logistic Regression

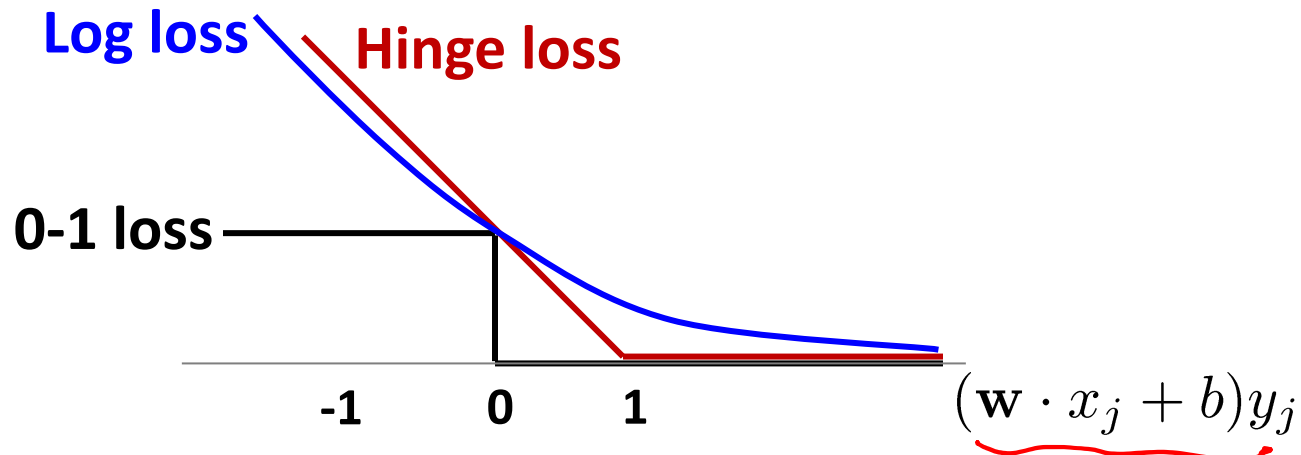
$$\max_{\mathbf{w}, b} \prod_{j=1}^n P(y_j | x_j, \mathbf{w}, b)$$

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)

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

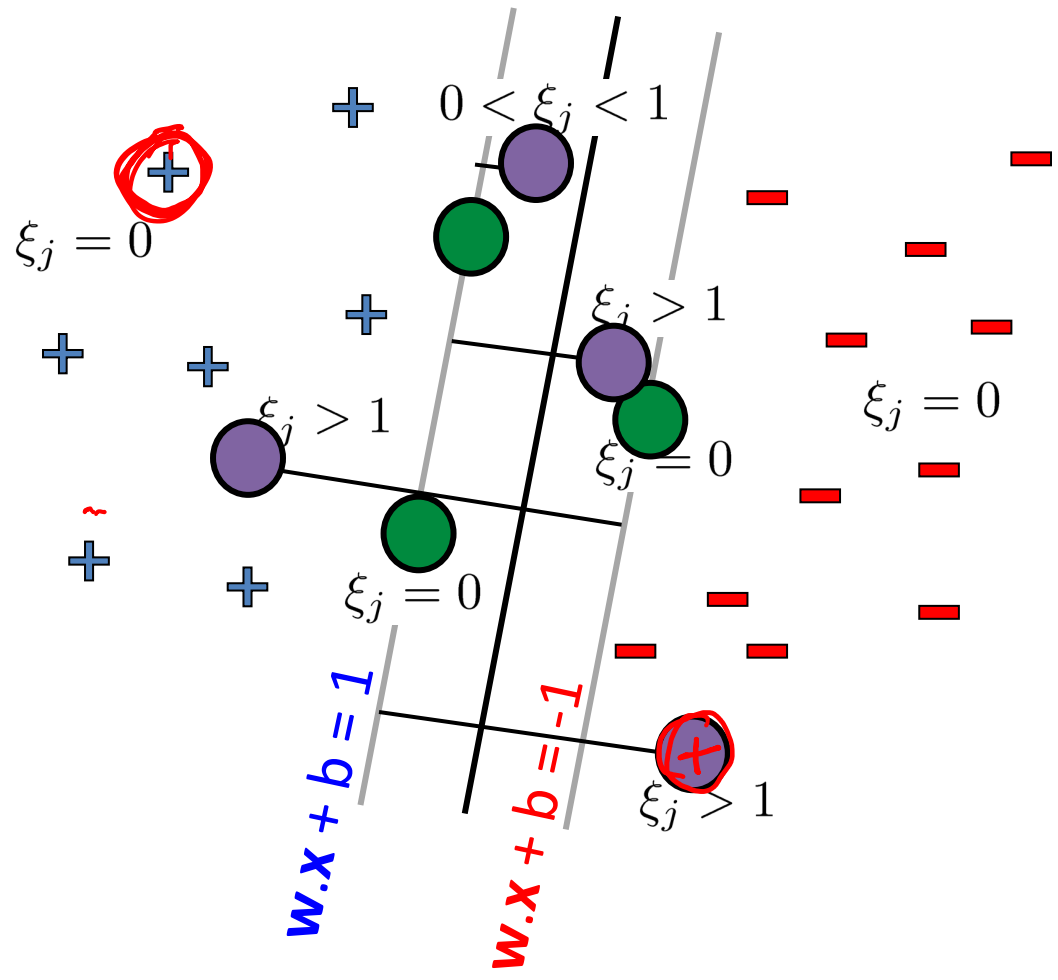


$$\min_{\mathbf{w}, b, \{\xi_j\}} \mathbf{w} \cdot \mathbf{w} + C \sum \xi_j$$

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

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

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)

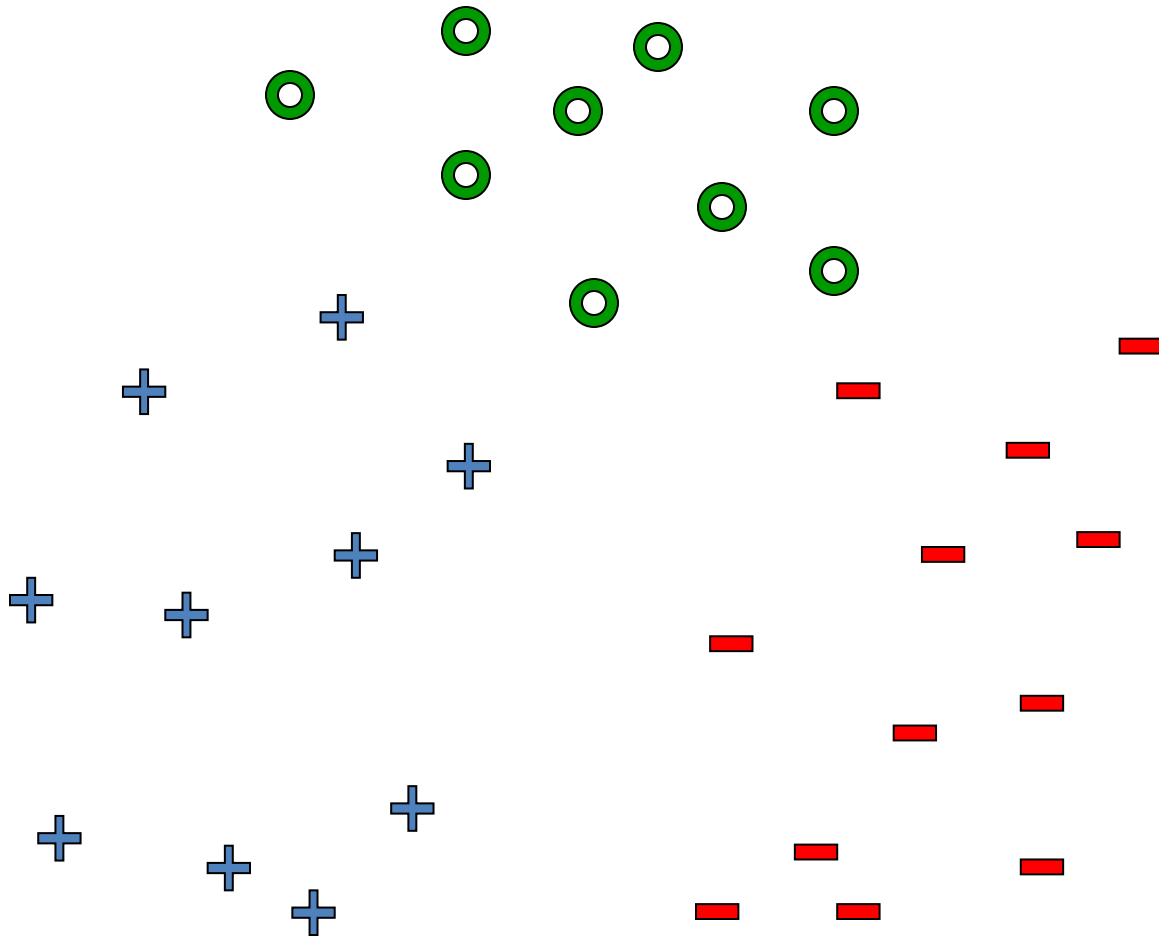
Correctly classified but on margin

Non-margin support vectors

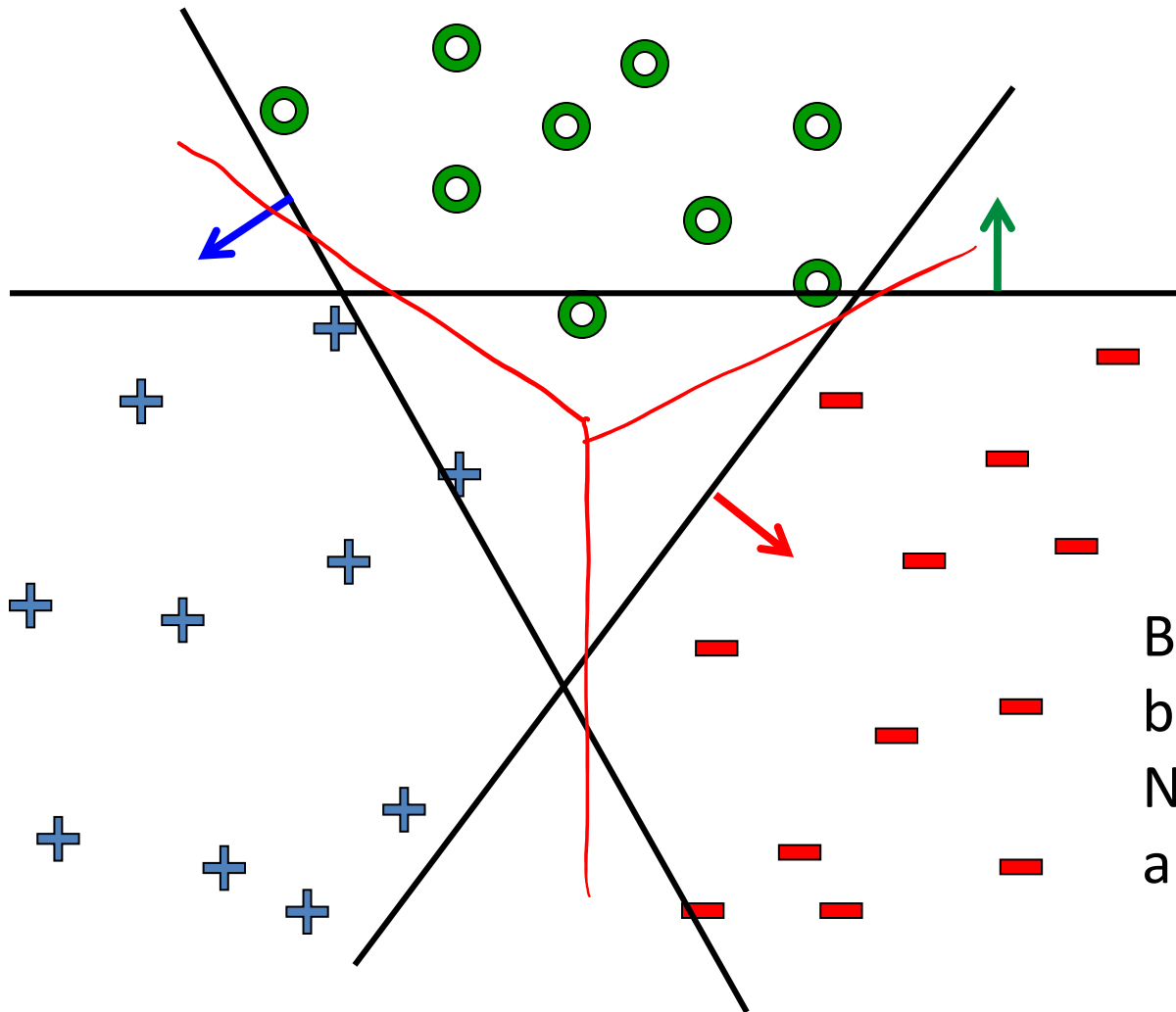
$\xi_j > 0$
 (contribute to both objective and constraints)

- ✓ $1 > \xi_j > 0$ Correctly classified but inside margin
- ✓ $\xi_j > 1$ Incorrectly classified

What about multiple classes?



One vs. rest



Learn 3 classifiers
separately:

Class k vs. rest

$$(\mathbf{w}_k, b_k)_{k=1,2,3}$$

$$y = \arg \max_k \mathbf{w}_k \cdot x + b_k$$

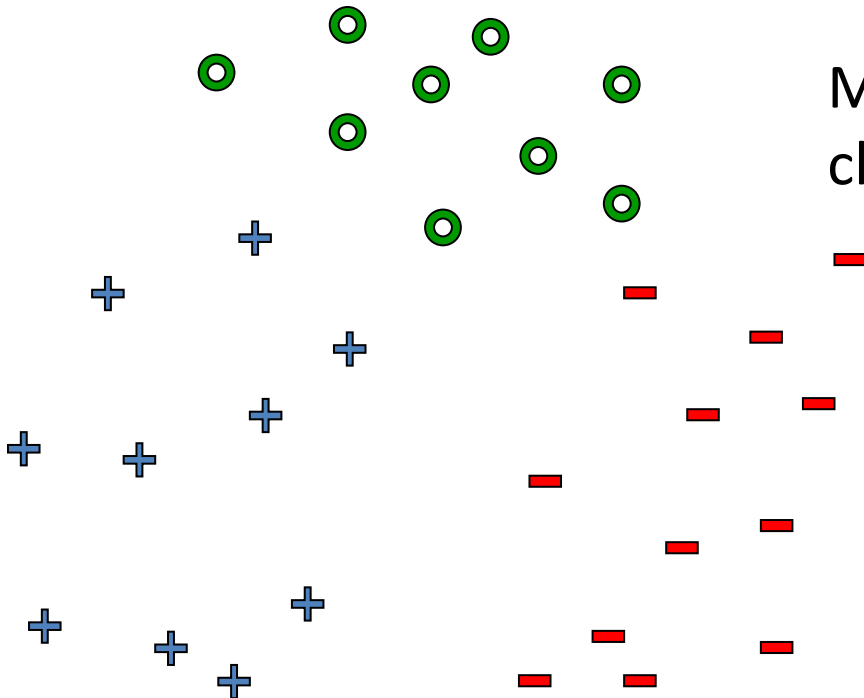
But \mathbf{w}_k s may not be
based on the same scale.
Note: $(a\mathbf{w}) \cdot x + (ab)$ is also
a solution

Learn 1 classifier: Multi-class SVM

Simultaneously learn 3 sets of weights

$$\min_{\{w^{(y)}\}, \{b^{(y)}\}} \sum_y \underline{w^{(y)} \cdot w^{(y)}}$$

$$\underline{w^{(y_j)} \cdot x_j + b^{(y_j)}} \geq \underline{w^{(y')} \cdot x_j + b^{(y')}} + \underline{1}, \quad \forall y' \neq y_j, \quad \forall j$$



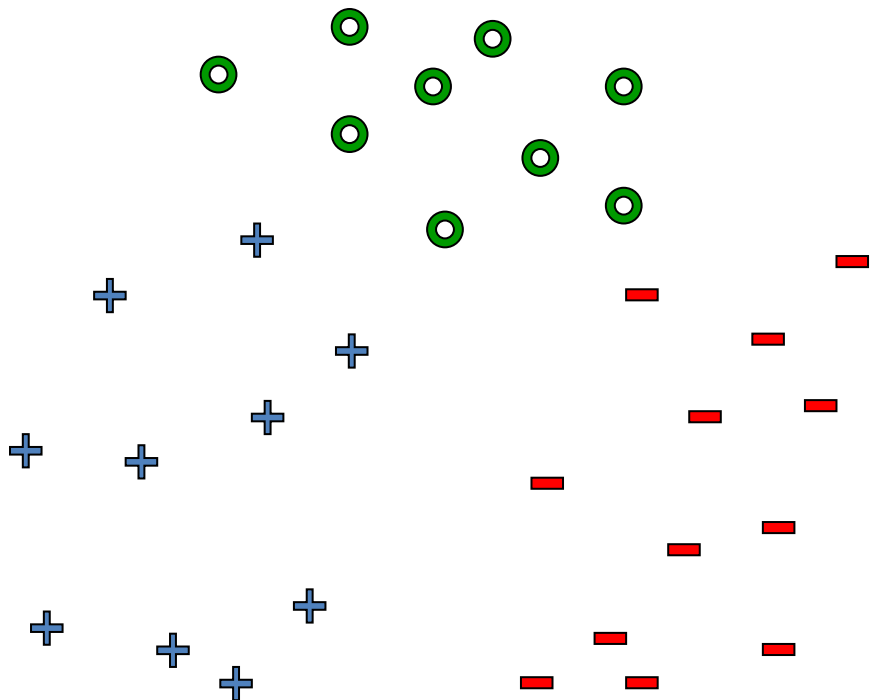
Margin - gap between correct class and nearest other class

$$y = \arg \max_k \underline{w^{(k)} \cdot x + b^{(k)}}$$

Learn 1 classifier: Multi-class SVM

Simultaneously learn 3 sets of weights

$$\begin{aligned} &\text{minimize} && \sum_y \mathbf{w}^{(y)} \cdot \mathbf{w}^{(y)} + C \sum_j \sum_{y \neq y_j} \xi_j^{(y)} && \text{over } \{\mathbf{w}^{(y)}\}, \{b^{(y)}\}, \{\xi_j^{(y)}\} \\ &\mathbf{w}^{(y_j)} \cdot \mathbf{x}_j + b^{(y_j)} \geq \mathbf{w}^{(y)} \cdot \mathbf{x}_j + b^{(y)} + \underbrace{1 - \xi_j^{(y)}}_{\text{margin}}, && \forall y \neq y_j, \forall j \\ &\xi_j^{(y)} \geq 0 && , \forall y \neq y_j, \forall j \end{aligned}$$



$$y = \arg \max \mathbf{w}^{(k)} \cdot \mathbf{x} + b^{(k)}$$

Joint optimization: \mathbf{w}_k s have the same scale.