Support Vector Machines (SVMs)

Aarti Singh

Machine Learning 10-315 Oct 20, 2021

Discriminative Classifiers

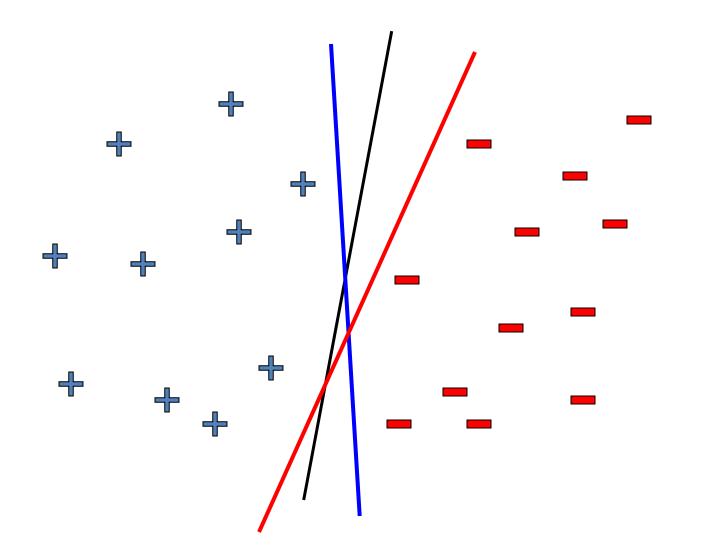
Optimal Classifier:

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

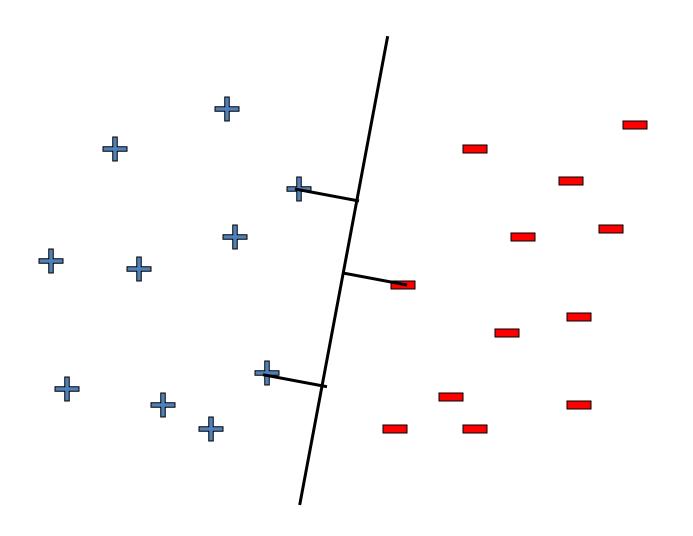
Why not learn P(Y|X) directly? Or better yet, why not learn the decision boundary directly?

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

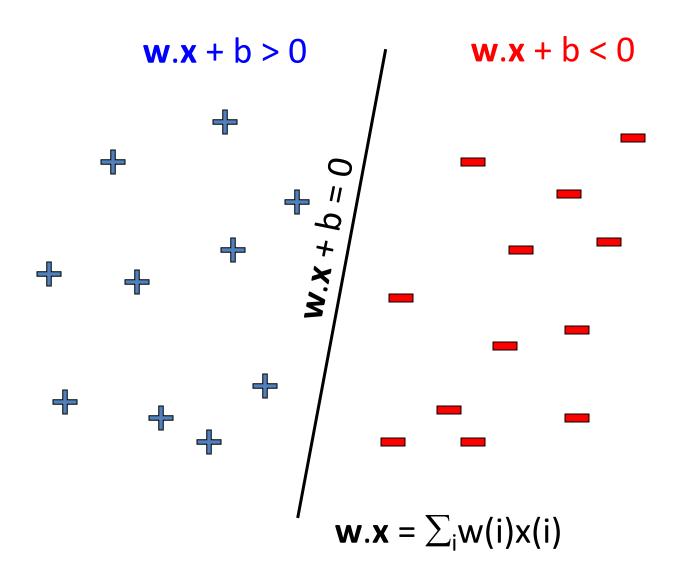
Linear classifiers – which line is better?



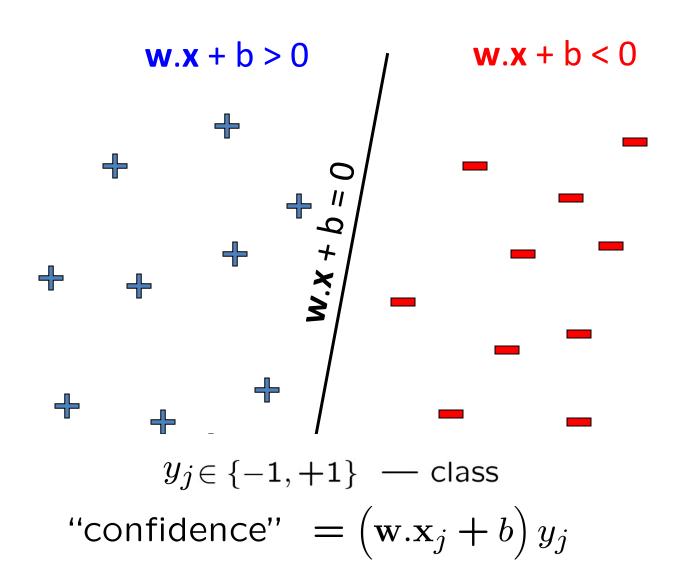
Pick the one with the largest margin!

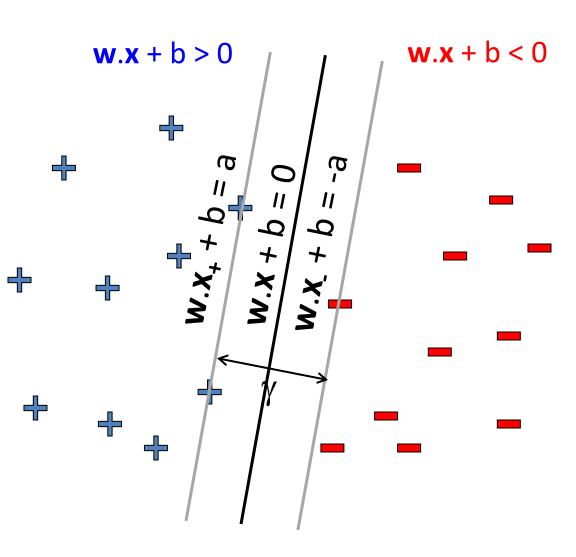


Parameterizing the decision boundary



Parameterizing the decision boundary



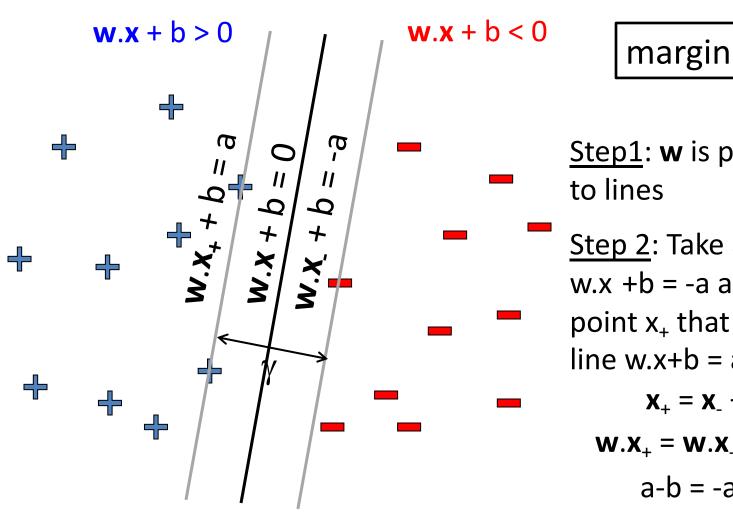


Distance of closest examples from the line/hyperplane

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

Step 1: **w** is perpendicular to lines since for any x_1 , x_2 on line **w**.($\mathbf{x}_1 - \mathbf{x}_2$) = 0

$$0 \neq x_1$$



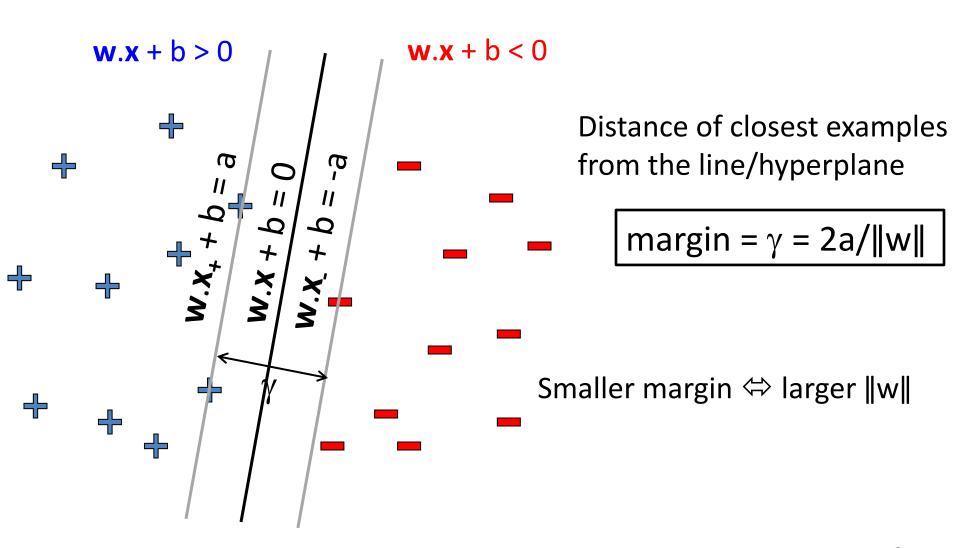
margin = γ = 2a/ $\|\mathbf{w}\|$

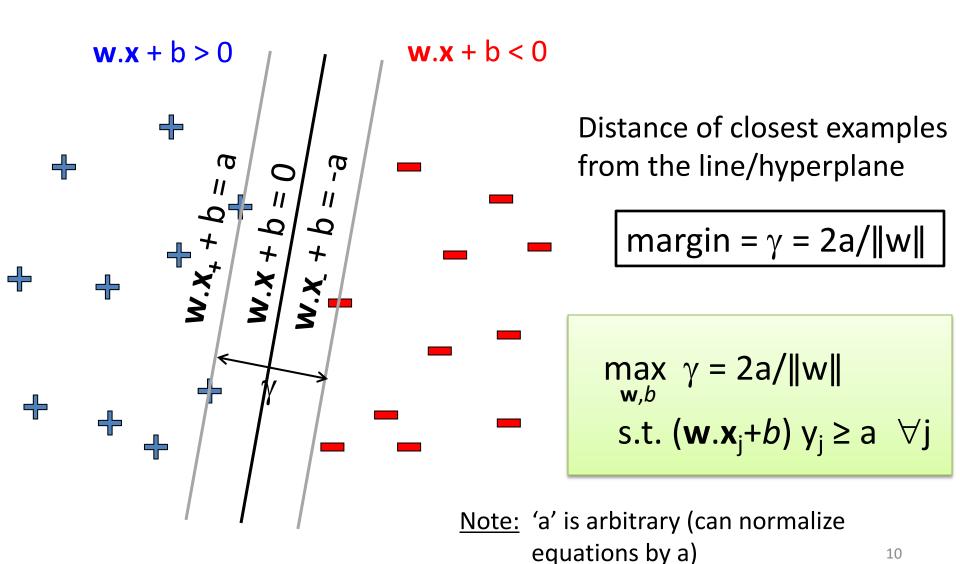
Step1: w is perpendicular

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

$$\mathbf{x}_{+} = \mathbf{x}_{-} + \gamma \mathbf{w} / \| \mathbf{w} \|$$

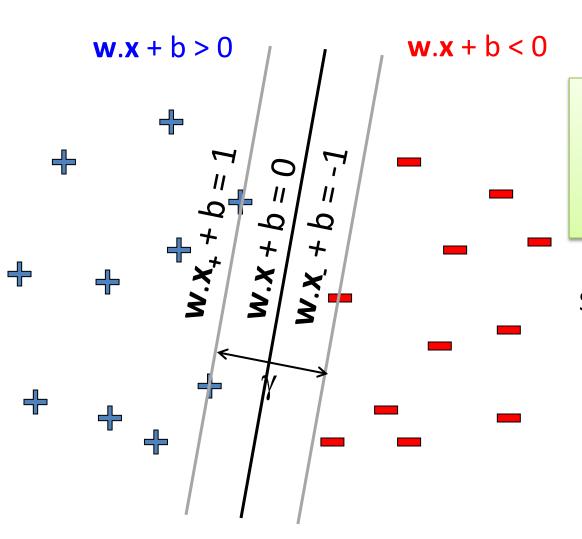
 $\mathbf{w}.\mathbf{x}_{+} = \mathbf{w}.\mathbf{x}_{-} + \gamma \mathbf{w}. \mathbf{w} / \| \mathbf{w} \|$
 $\mathbf{a}-\mathbf{b} = -\mathbf{a}-\mathbf{b} + \gamma \| \mathbf{w} \|$
 $2\mathbf{a} = \gamma \| \mathbf{w} \|$





10

Support Vector Machines

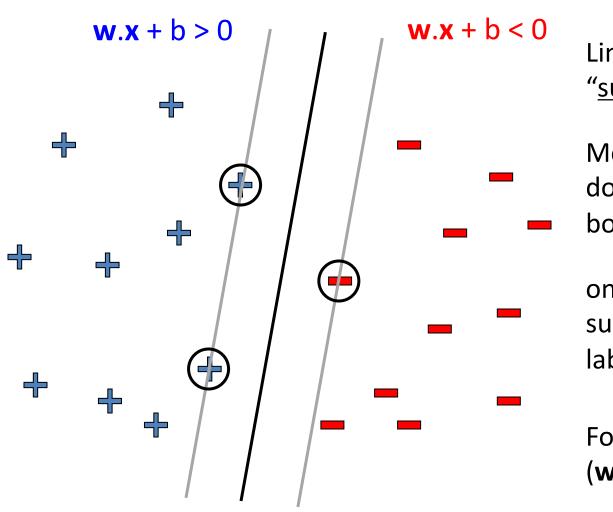


min w.w s.t. $(\mathbf{w}.\mathbf{x}_j+b)$ $y_j \ge 1 \quad \forall j$

Solve efficiently by quadratic programming (QP)

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

Support Vectors



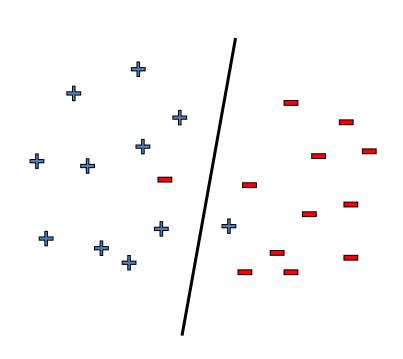
Linear hyperplane defined by "support vectors"

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 $(\mathbf{w}.\mathbf{x}_j+b)$ $\mathbf{y}_j=1$

What if data is not linearly separable?



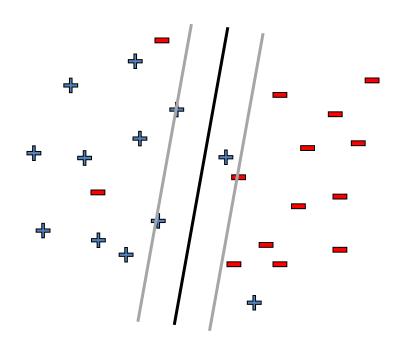
Use features of features of features of features....

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

But run risk of overfitting!

What if data is still not linearly separable?

Allow "error" in classification



Smaller margin ⇔ larger ||w||

min
$$\mathbf{w}.\mathbf{w} + C$$
 #mistakes s.t. $(\mathbf{w}.\mathbf{x}_j+b)$ $y_j \ge 1 \quad \forall j$

Maximize margin and minimize # mistakes on training data

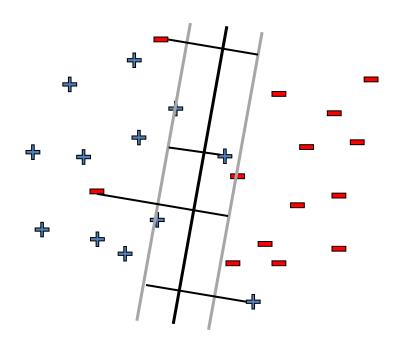
C - tradeoff parameter

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

$$\min_{\mathbf{w},b,\{\xi_{j}\}} \mathbf{w}.\mathbf{w} + C \sum_{j} \xi_{j}$$

$$s.t. (\mathbf{w}.\mathbf{x}_{j}+b) y_{j} \ge 1-\xi_{j} \quad \forall j$$

$$\xi_{j} \ge 0 \quad \forall j$$

$$\xi_j$$
 - "slack" variables
= (>1 if x_i misclassifed)

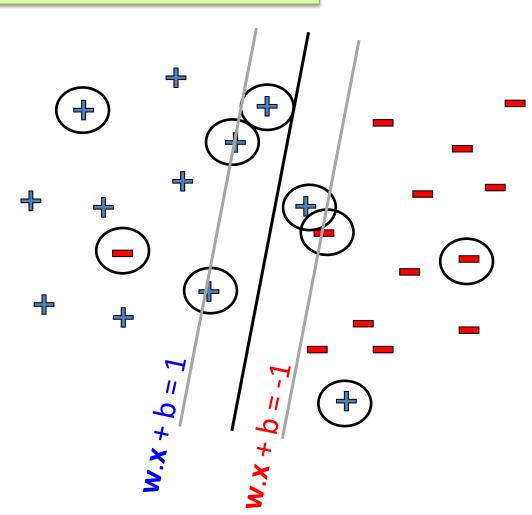
pay linear penalty if mistake

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

min
$$\mathbf{w}.\mathbf{w} + C \Sigma \xi_j$$

 $\mathbf{w},b,\{\xi_j\}$
s.t. $(\mathbf{w}.\mathbf{x}_j+b) y_j \ge 1-\xi_j \quad \forall j$
 $\xi_j \ge 0 \quad \forall j$

s.t. $(\mathbf{w}.\mathbf{x}_j+b)$ $\mathbf{y}_j \geq 1-\xi_j \quad \forall j$ Slack variables

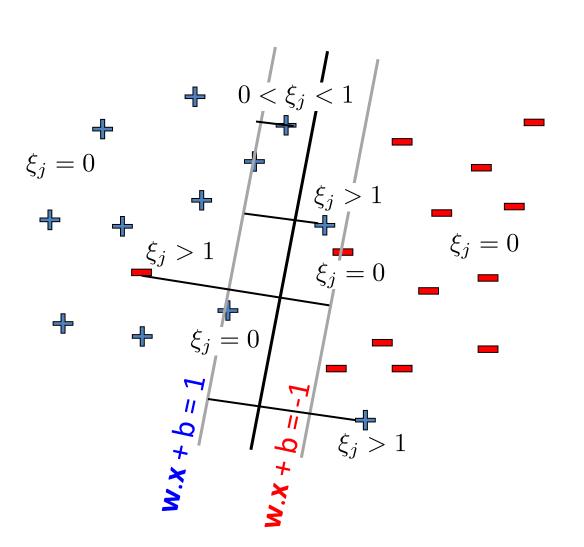


$$(\mathbf{w}.\mathbf{x}_i+b) \mathbf{y}_i \geq 1-\xi_i \quad \forall \mathbf{j}$$

What is the slack ξ_j for the following points?

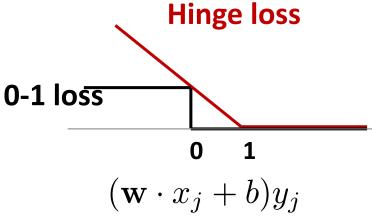
Confidence | Slack

Slack variables – Hinge loss



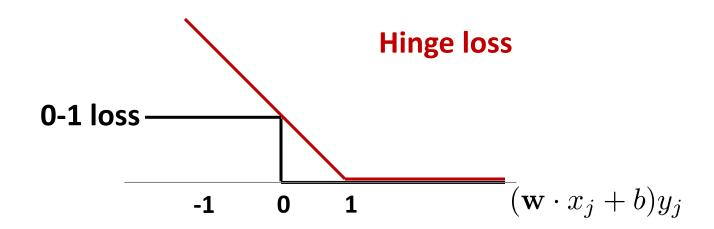
Notice that

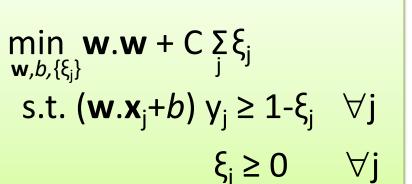
$$\xi_j = (1 - (\mathbf{w} \cdot x_j + b)y_j))_+$$



Slack variables – Hinge loss

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







Regularized hinge loss

$$\min_{\mathbf{w},b} \mathbf{w}.\mathbf{w} + C \sum_{j} (1-(\mathbf{w}.x_j+b)y_j)_+$$

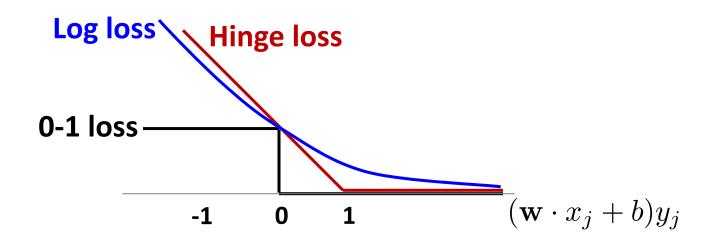
SVM vs. Logistic Regression

SVM: **Hinge loss**

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

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

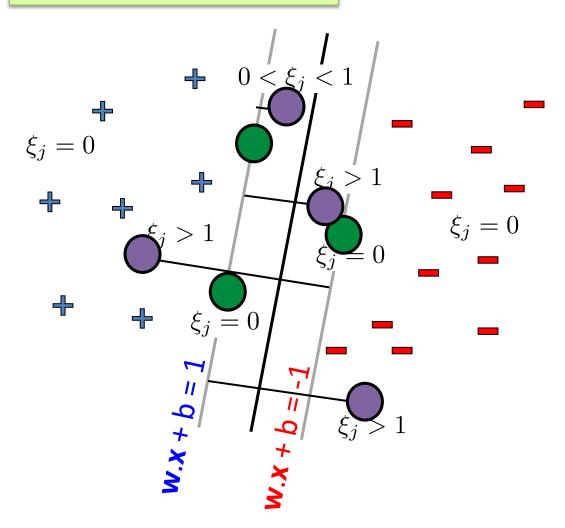
$$loss(f(x_j), y_j) = -\log P(y_j \mid x_j, \mathbf{w}, b) = \log(1 + e^{-(\mathbf{w} \cdot x_j + b)y_j})$$



min
$$\mathbf{w}.\mathbf{w} + C \Sigma \xi_j$$

 $\mathbf{w},b,\{\xi_j\}$
s.t. $(\mathbf{w}.\mathbf{x}_j+b) y_j \ge 1-\xi_j \quad \forall j$
 $\xi_j \ge 0 \quad \forall j$

Support Vectors



Margin support vectors

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

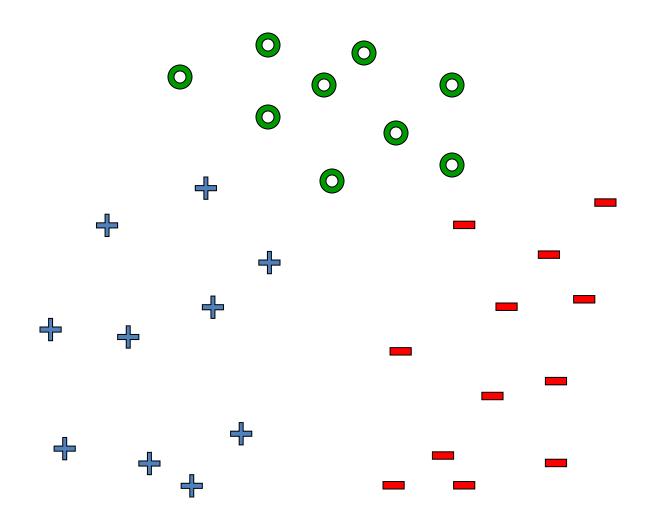
Correctly classified but on margin

Non-margin support vectors

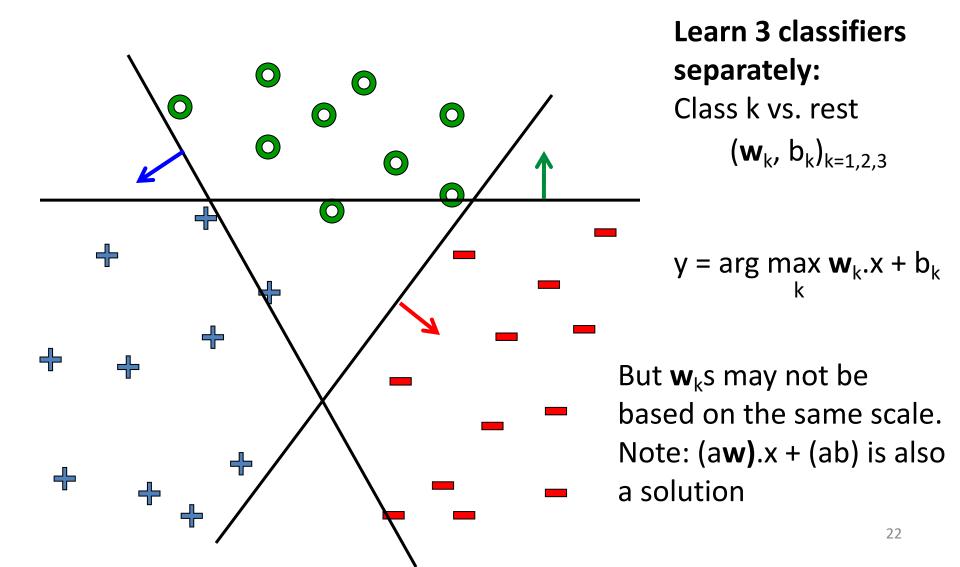
 $\xi_j > 0$ (contribute to both objective and constraints)

 $1 > \xi_j > 0$ Correctly classified but inside margin $\xi_i > 1$ Incorrectly classified 20

What about multiple classes?

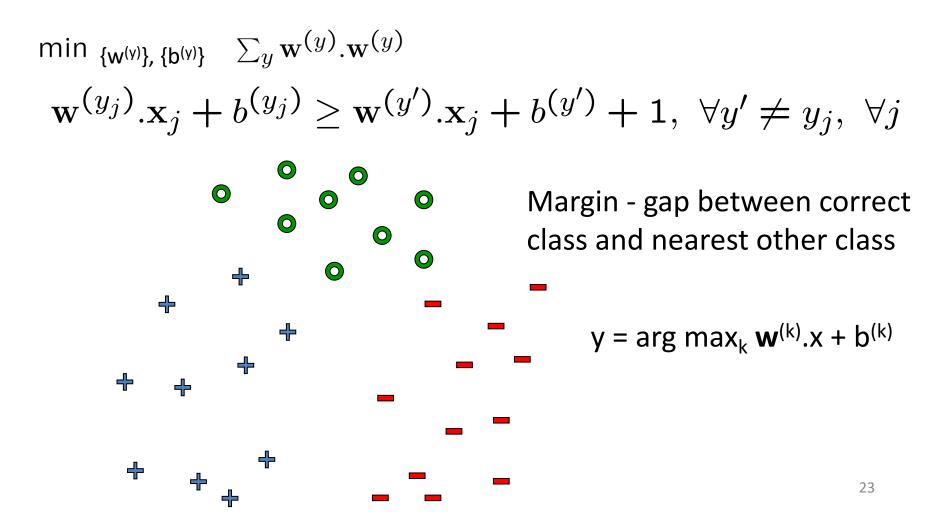


One vs. rest



Learn 1 classifier: Multi-class SVM

Simultaneously learn 3 sets of weights



Learn 1 classifier: Multi-class SVM

Simultaneously learn 3 sets of weights

