



Support Vector Machines

Machine Learning – 10701/15781

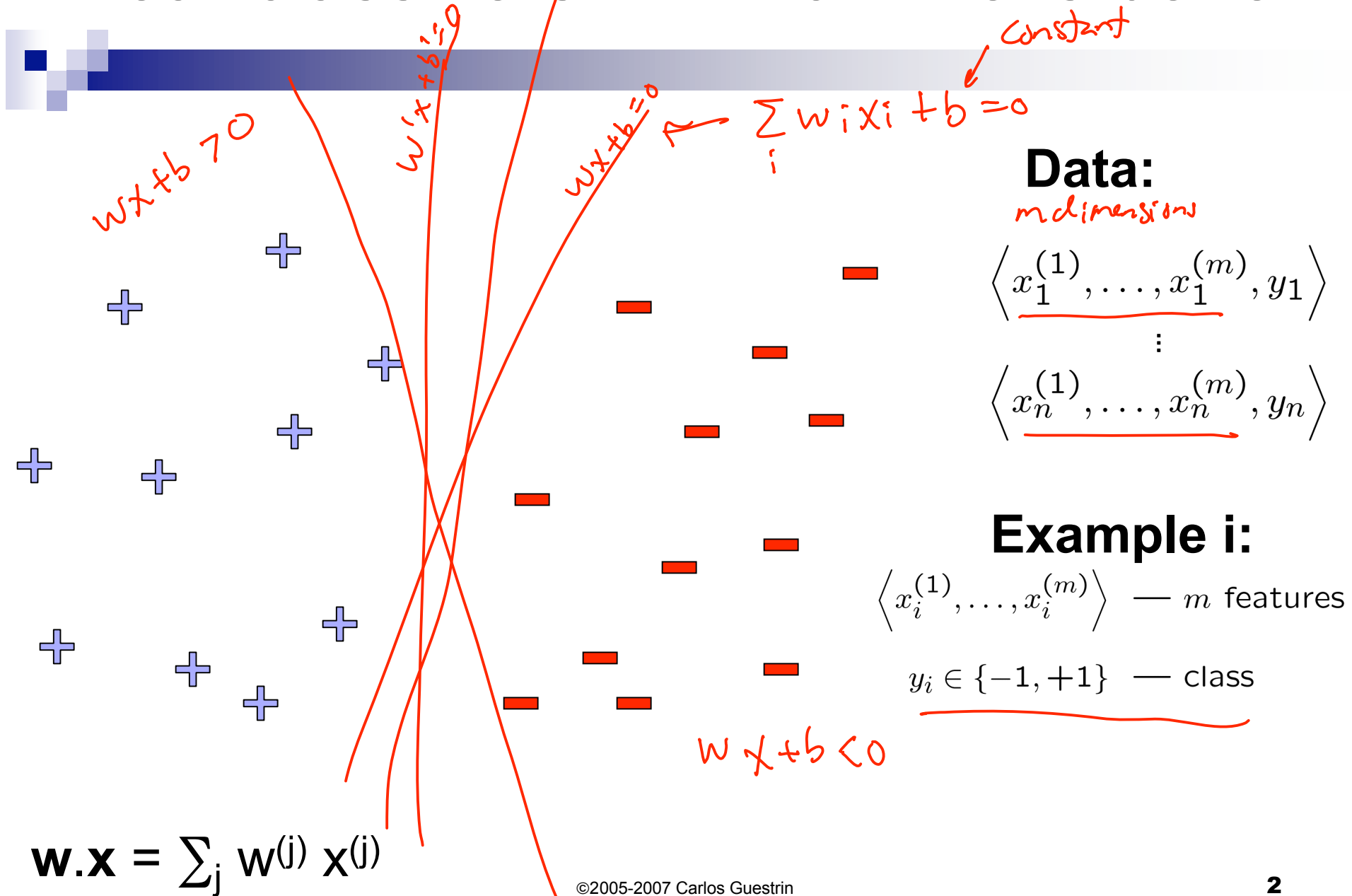
Carlos Guestrin

Carnegie Mellon University

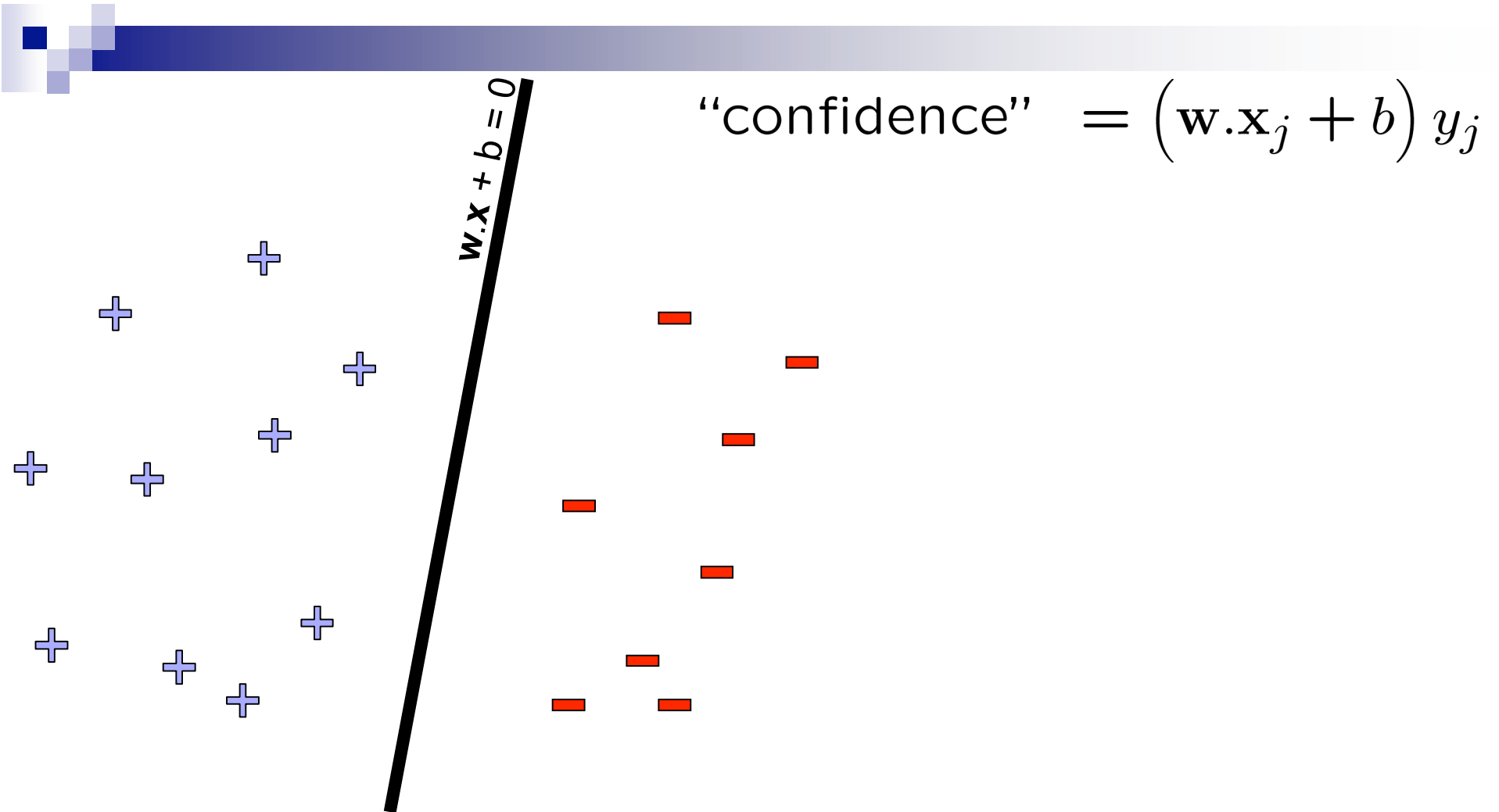
February 21st, 2007

©2005-2007 Carlos Guestrin

Linear classifiers – Which line is better?

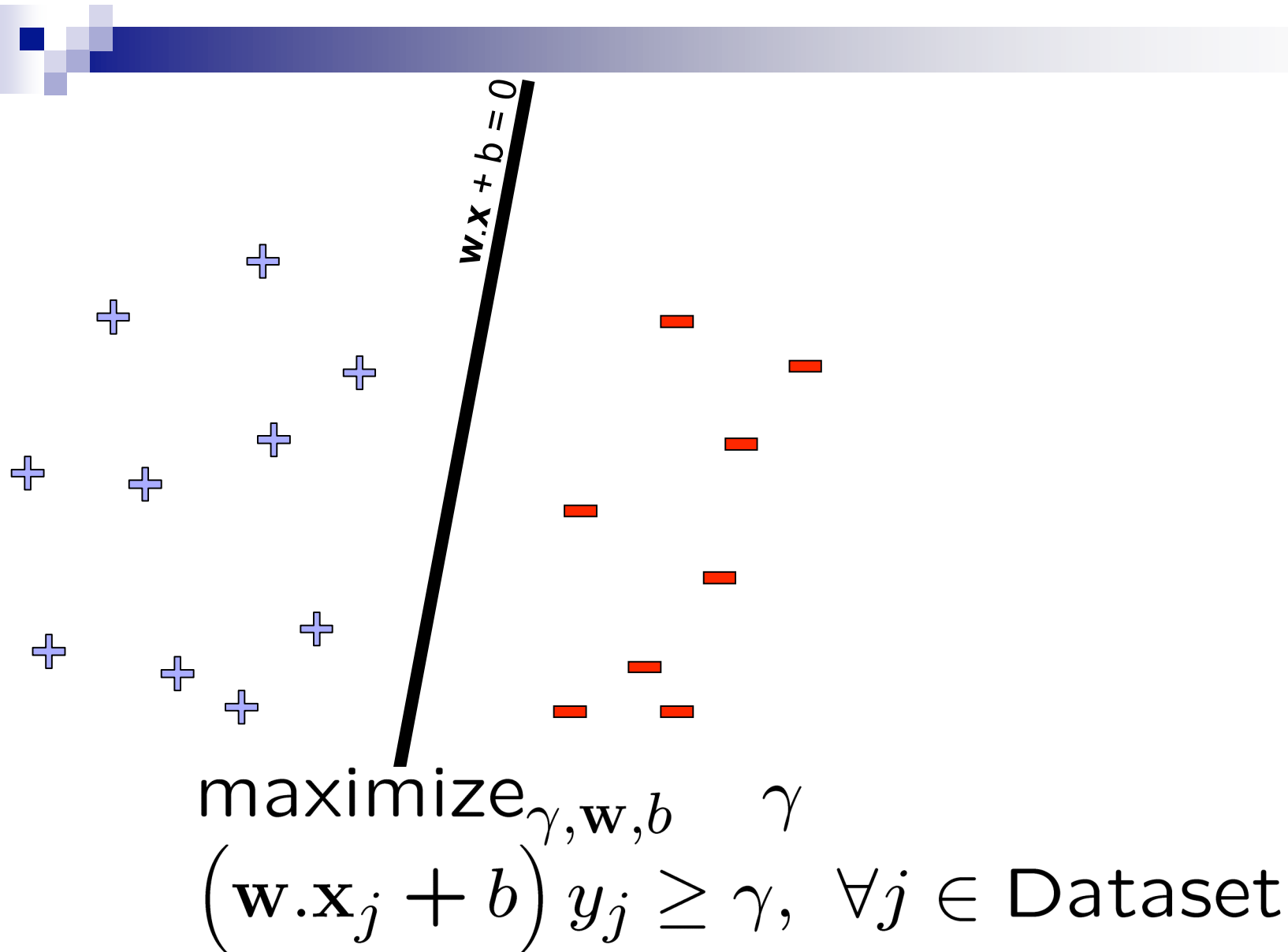


Pick the one with the largest margin!

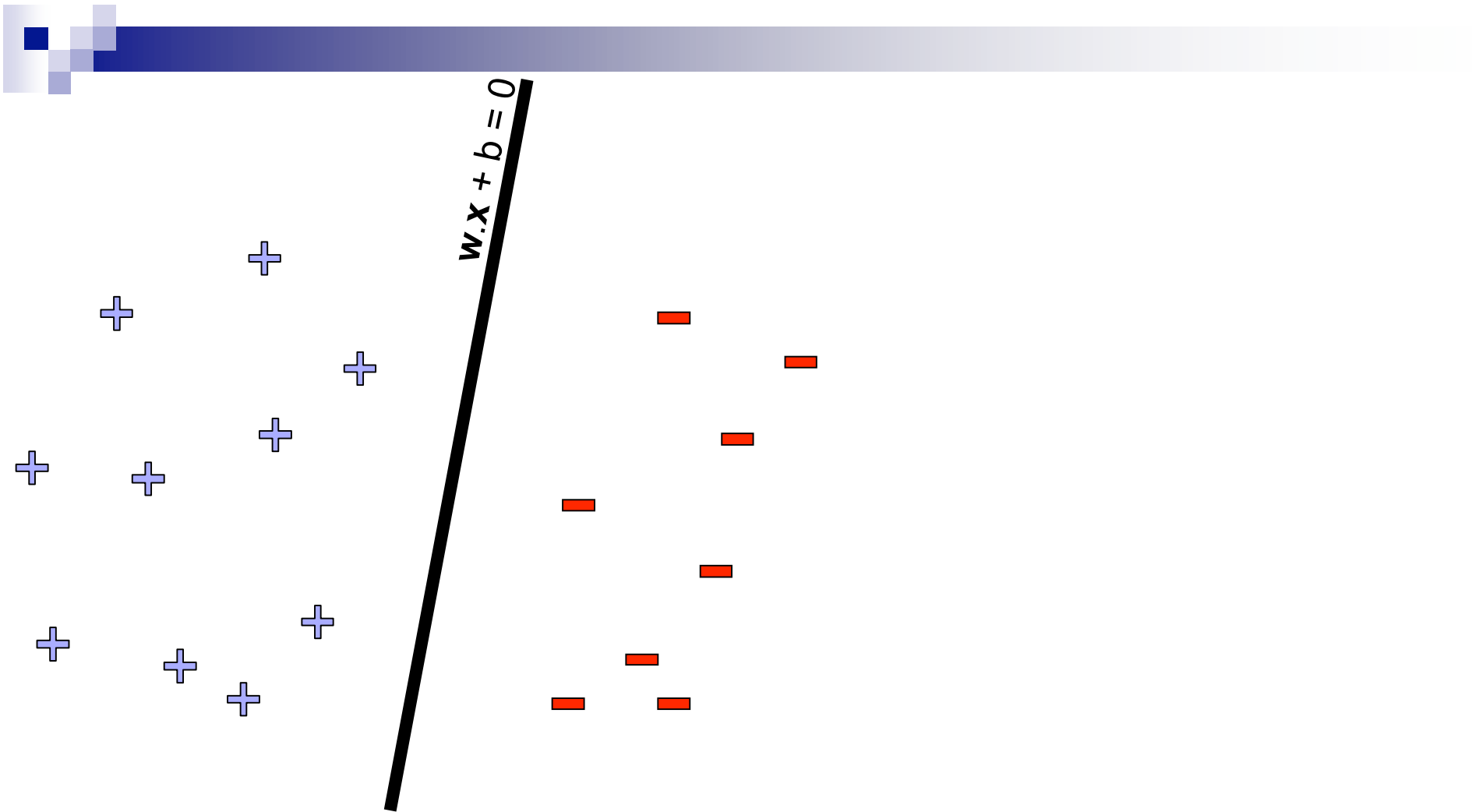


$$\mathbf{w} \cdot \mathbf{x} = \sum_j w^{(j)} x^{(j)}$$

Maximize the margin



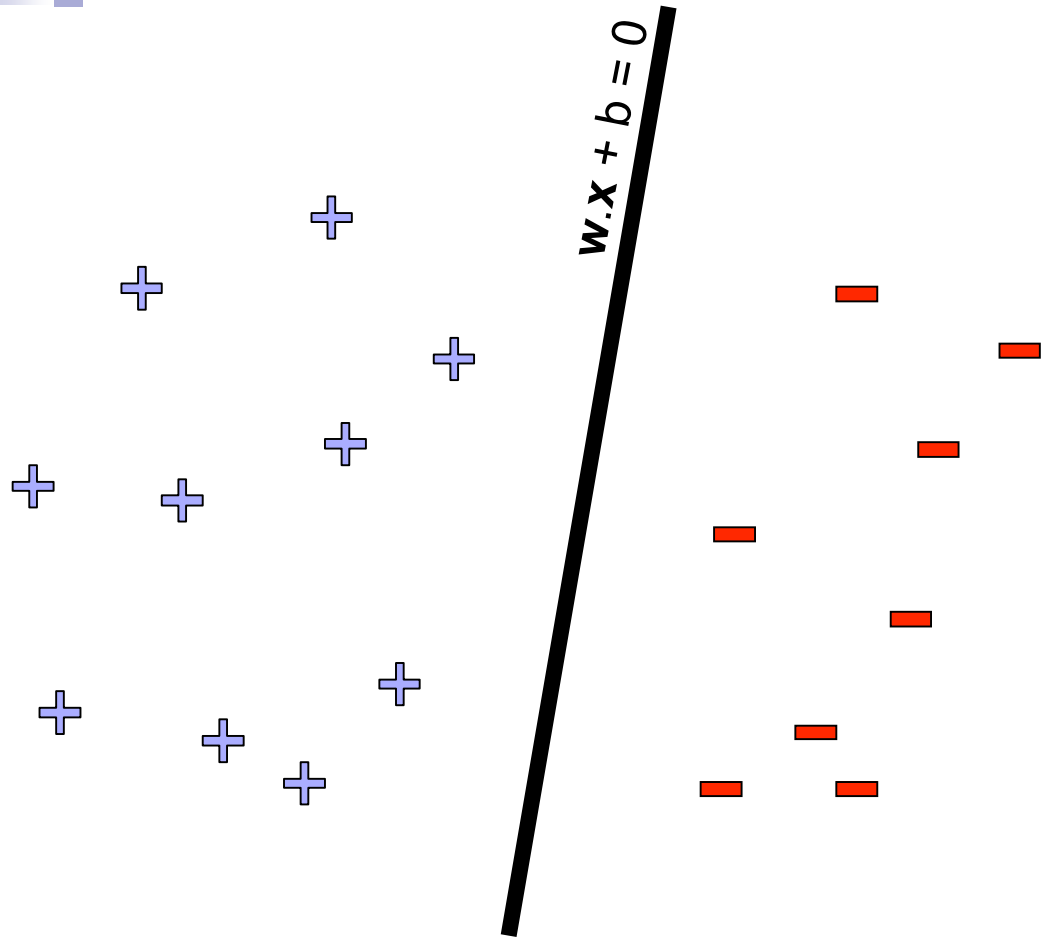
But there are a many planes...



Review: Normal to a plane

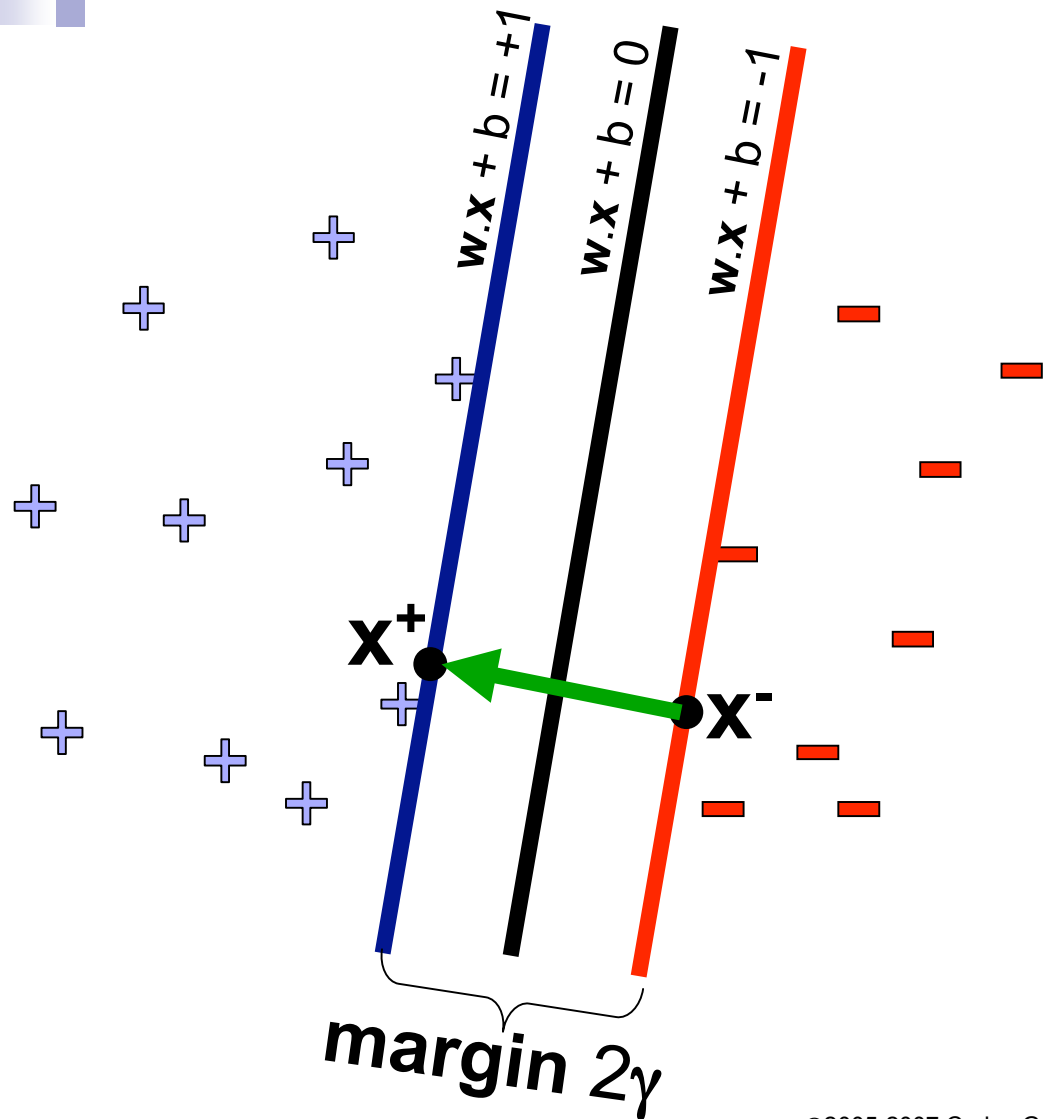


$$\mathbf{x}_j = \bar{\mathbf{x}}_j + \lambda \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

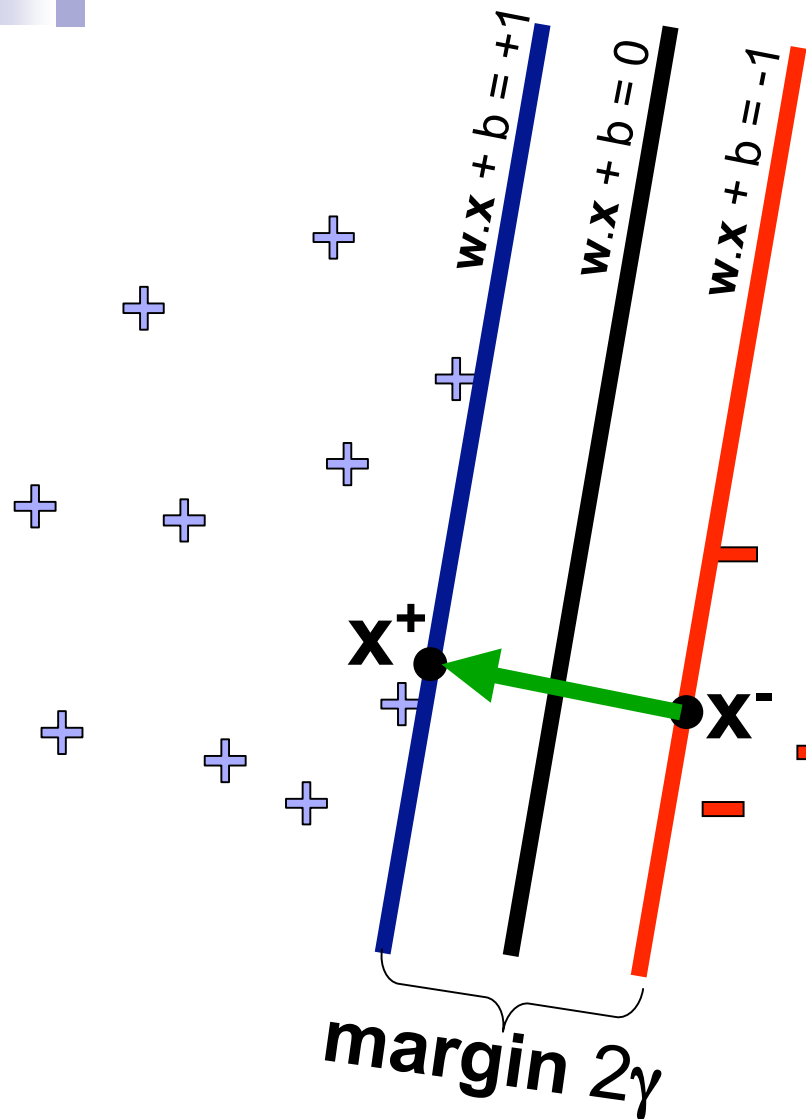


Normalized margin – Canonical hyperplanes

$$\mathbf{x}_j = \bar{\mathbf{x}}_j + \lambda \frac{\mathbf{w}}{\|\mathbf{w}\|}$$



Normalized margin – Canonical hyperplanes



$$\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$$

$$\mathbf{w} \cdot \mathbf{x}^+ + b = 1$$

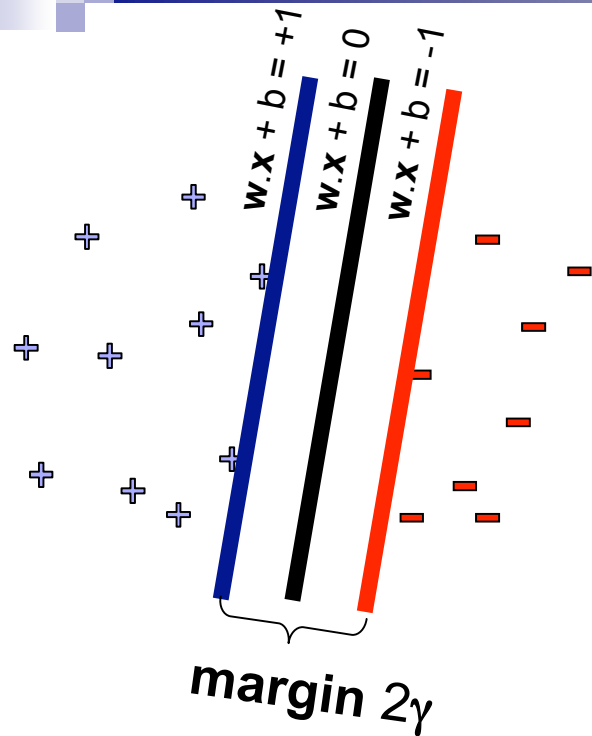
$$\mathbf{w} \cdot \left(\mathbf{x}^- + \lambda \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + b = 1$$

$$\lambda = \frac{2}{\|\mathbf{w}\|}$$

$$\gamma = \frac{1}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$

Margin maximization using canonical hyperplanes

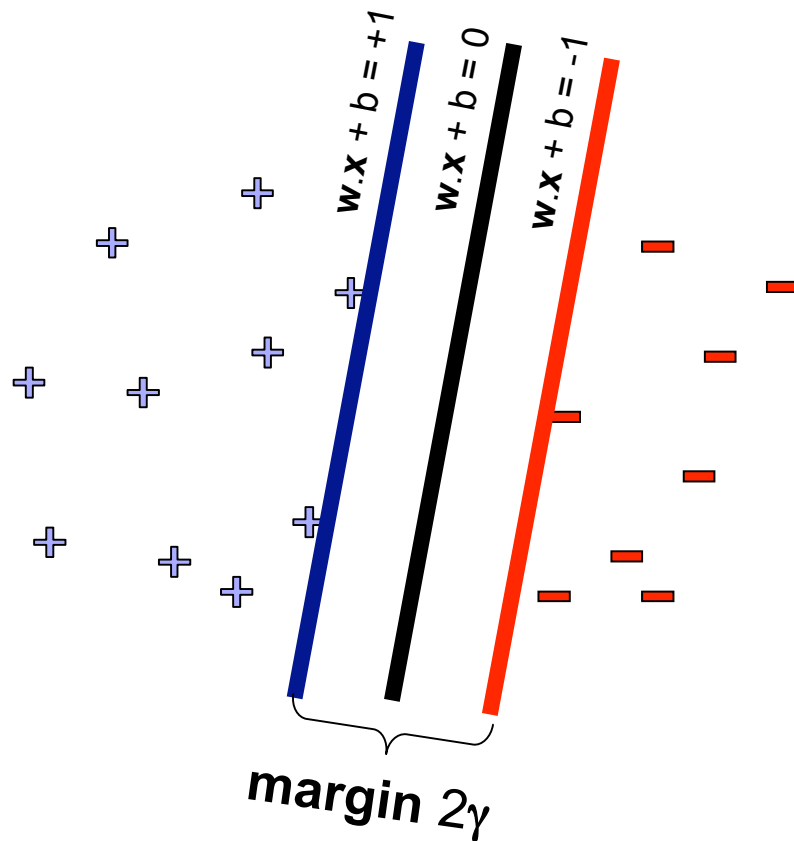
$$\gamma = \frac{1}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$



$$\text{maximize}_{\gamma, \mathbf{w}, b} \quad \gamma$$
$$\left(\mathbf{w} \cdot \mathbf{x}_j + b \right) y_j \geq \gamma, \quad \forall j \in \text{Dataset}$$

$$\text{minimize}_{\mathbf{w}, b} \quad \mathbf{w} \cdot \mathbf{w}$$
$$\left(\mathbf{w} \cdot \mathbf{x}_j + b \right) y_j \geq 1, \quad \forall j \in \text{Dataset}$$

Support vector machines (SVMs)



minimize _{w, b} $w \cdot w$

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

- Solve efficiently by quadratic programming (QP)
 - Well-studied solution algorithms
- Hyperplane defined by support vectors

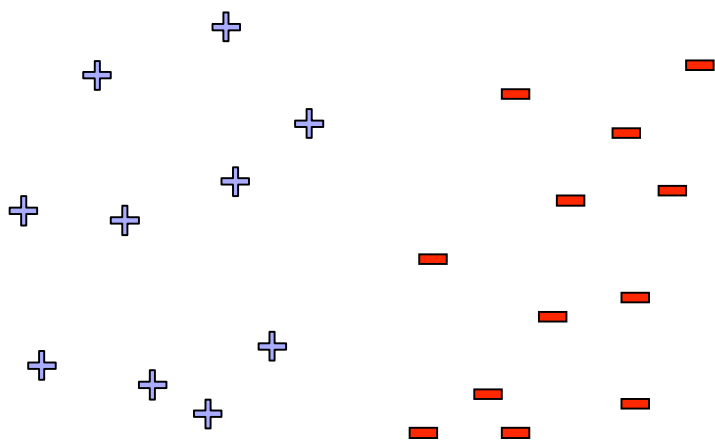
Announcements



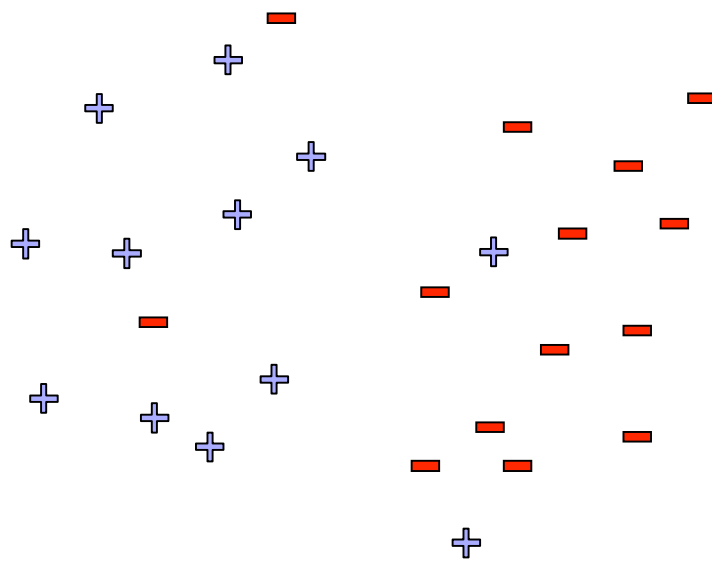
- Third homework out later today
- This one is shorter!!!! :)
- Due on Monday March 5th
- No late days allowed
 - so we can give solutions before midterm

What if the data is not linearly separable?

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



What if the data is still not linearly separable?

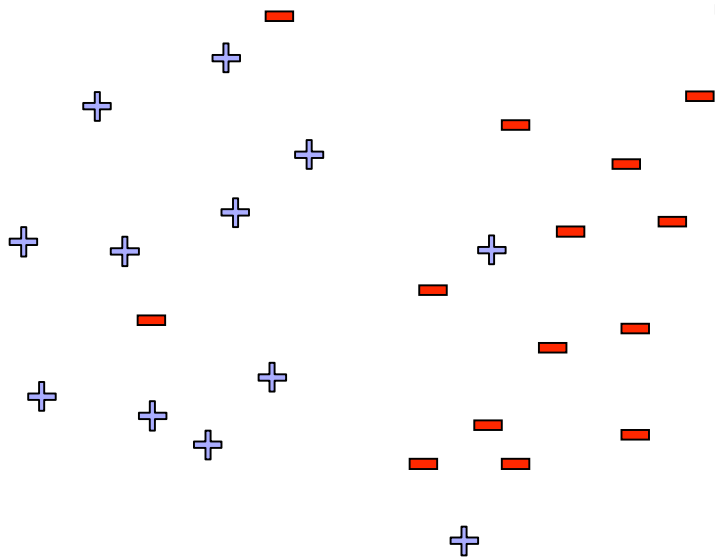


$$\text{minimize}_{\mathbf{w}, b} \quad \mathbf{w} \cdot \mathbf{w}$$
$$\left(\mathbf{w} \cdot \mathbf{x}_j + b \right) y_j \geq 1 \quad , \forall j$$

- Minimize $\mathbf{w} \cdot \mathbf{w}$ and number of training mistakes
 - Tradeoff two criteria?
- Tradeoff #(mistakes) and $\mathbf{w} \cdot \mathbf{w}$
 - 0/1 loss
 - Slack penalty C
 - Not QP anymore
 - Also doesn't distinguish near misses and really bad mistakes

Slack variables – Hinge loss

$$\text{minimize}_{w,b} \quad w \cdot w$$
$$\left(w \cdot x_j + b \right) y_j \geq 1 \quad , \forall j$$



- If margin ≥ 1 , don't care
- If margin < 1 , pay linear penalty

Side note: What's the difference between SVMs and logistic regression?

SVM:

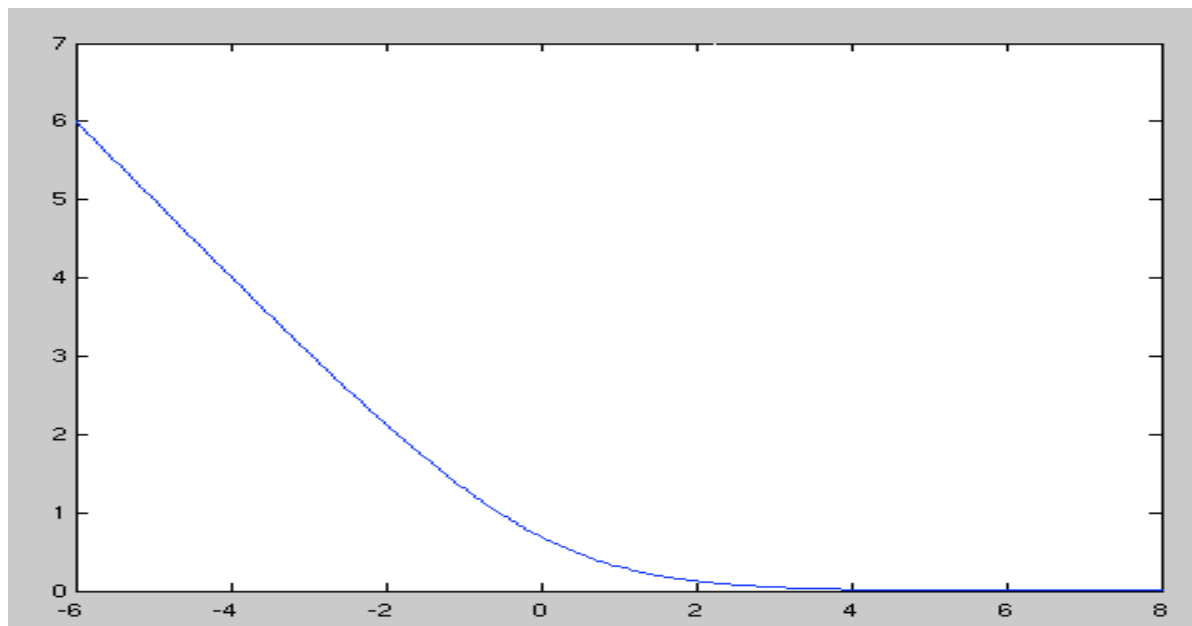
$$\begin{aligned} \text{minimize}_{\mathbf{w}, b} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j \\ & \xi_j \geq 0, \quad \forall j \end{aligned}$$

Logistic regression:

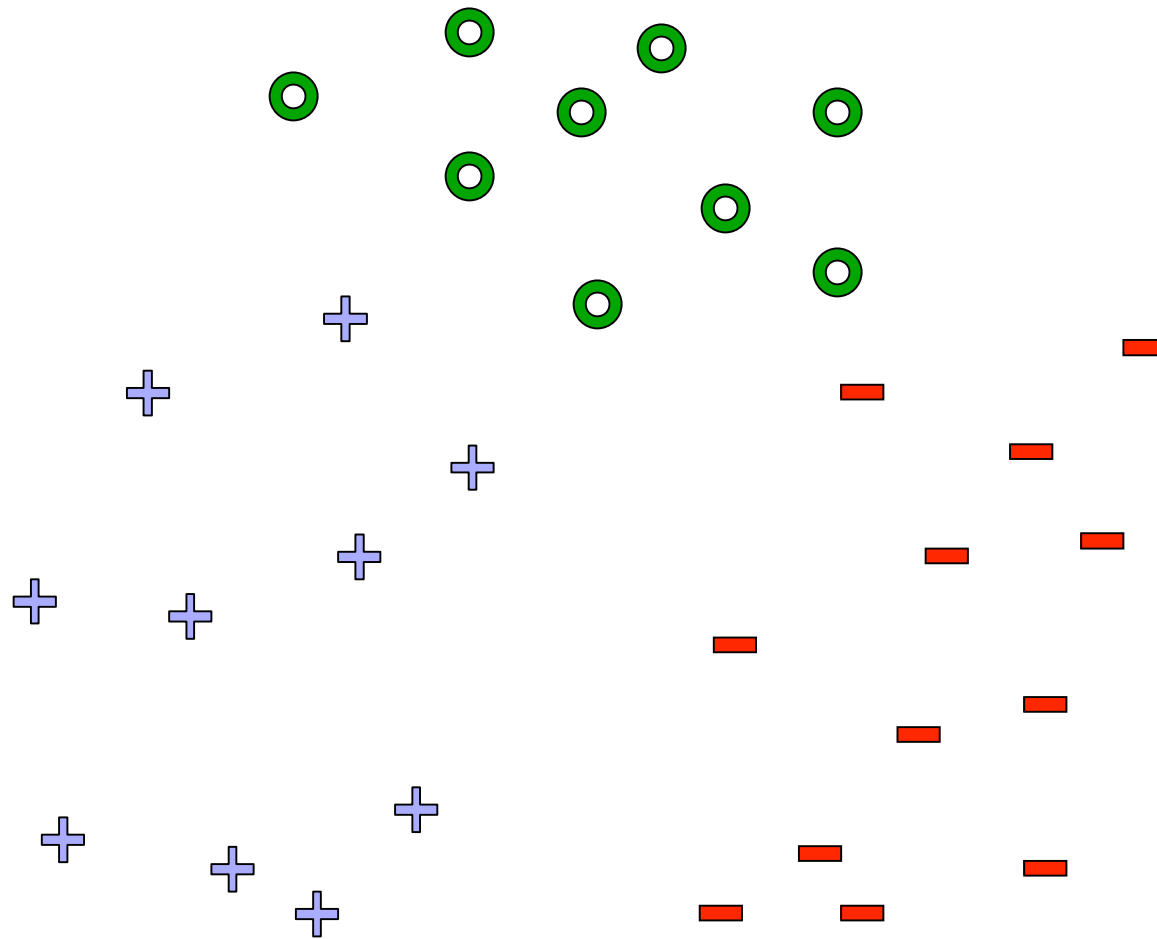
$$P(Y = 1 \mid x, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}}$$

Log loss:

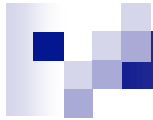
$$-\ln P(Y = 1 \mid x, \mathbf{w}) = \ln(1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)})$$



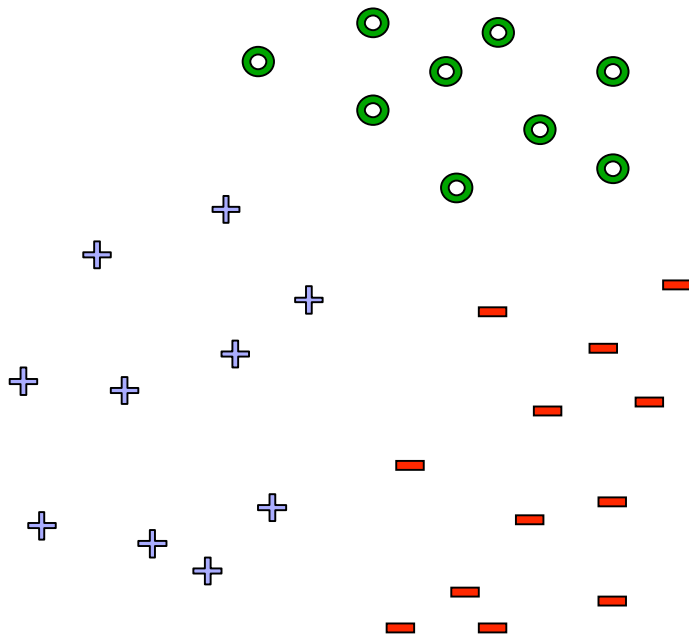
What about multiple classes?



One against All

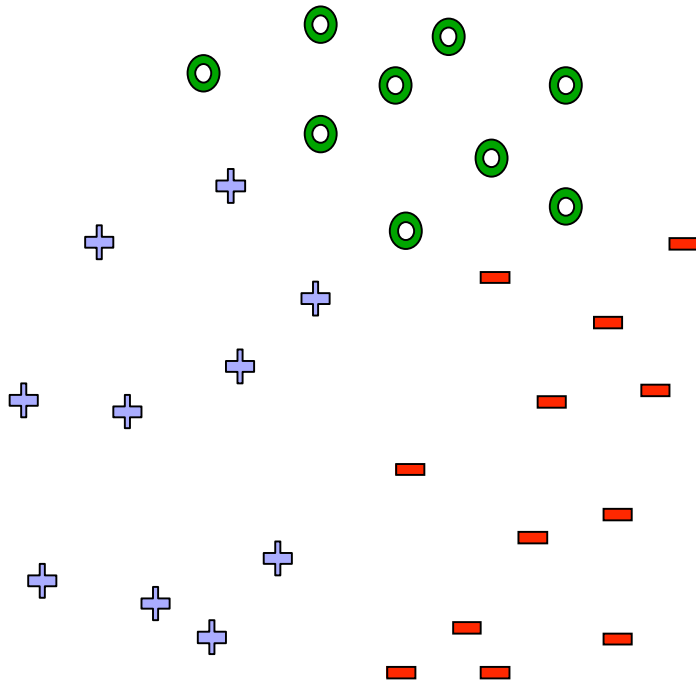


Learn 3 classifiers:



Learn 1 classifier: Multiclass SVM

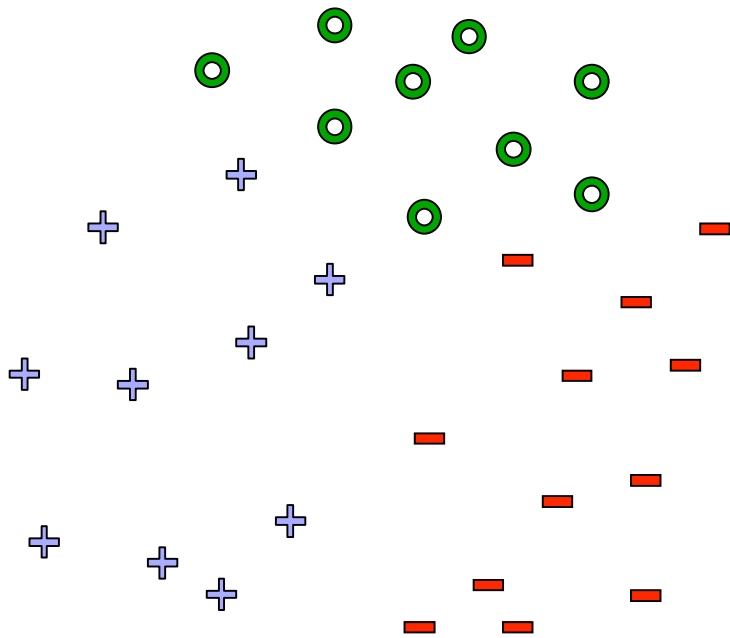
Simultaneously learn 3 sets of weights



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

Learn 1 classifier: Multiclass SVM

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



What you need to know



- Maximizing margin
- Derivation of SVM formulation
- Slack variables and hinge loss
- Relationship between SVMs and logistic regression
 - 0/1 loss
 - Hinge loss
 - Log loss
- Tackling multiple class
 - One against All
 - Multiclass SVMs