



# Support Vector Machines , *SVMs*

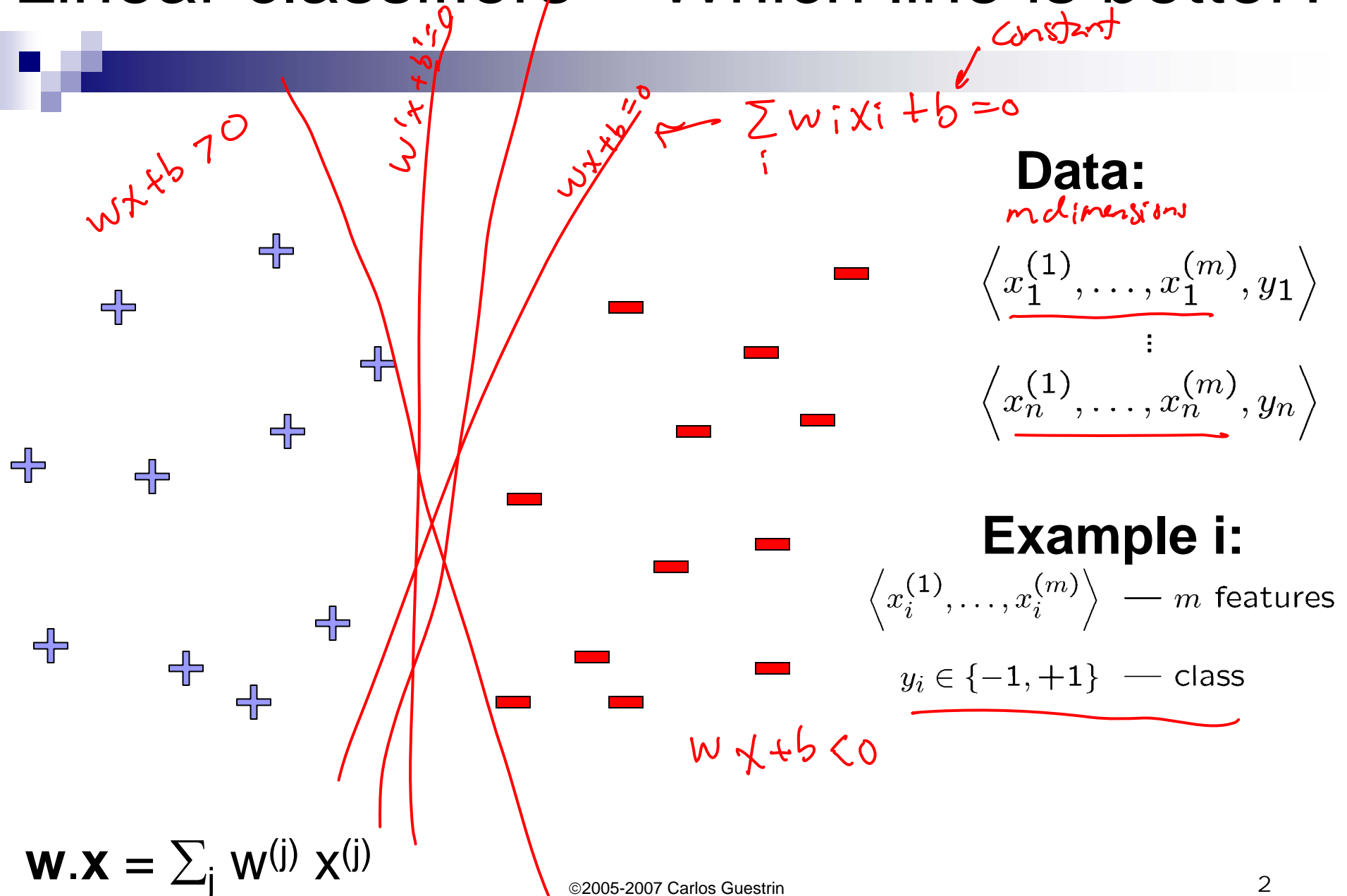
Machine Learning – 10701/15781

Carlos Guestrin

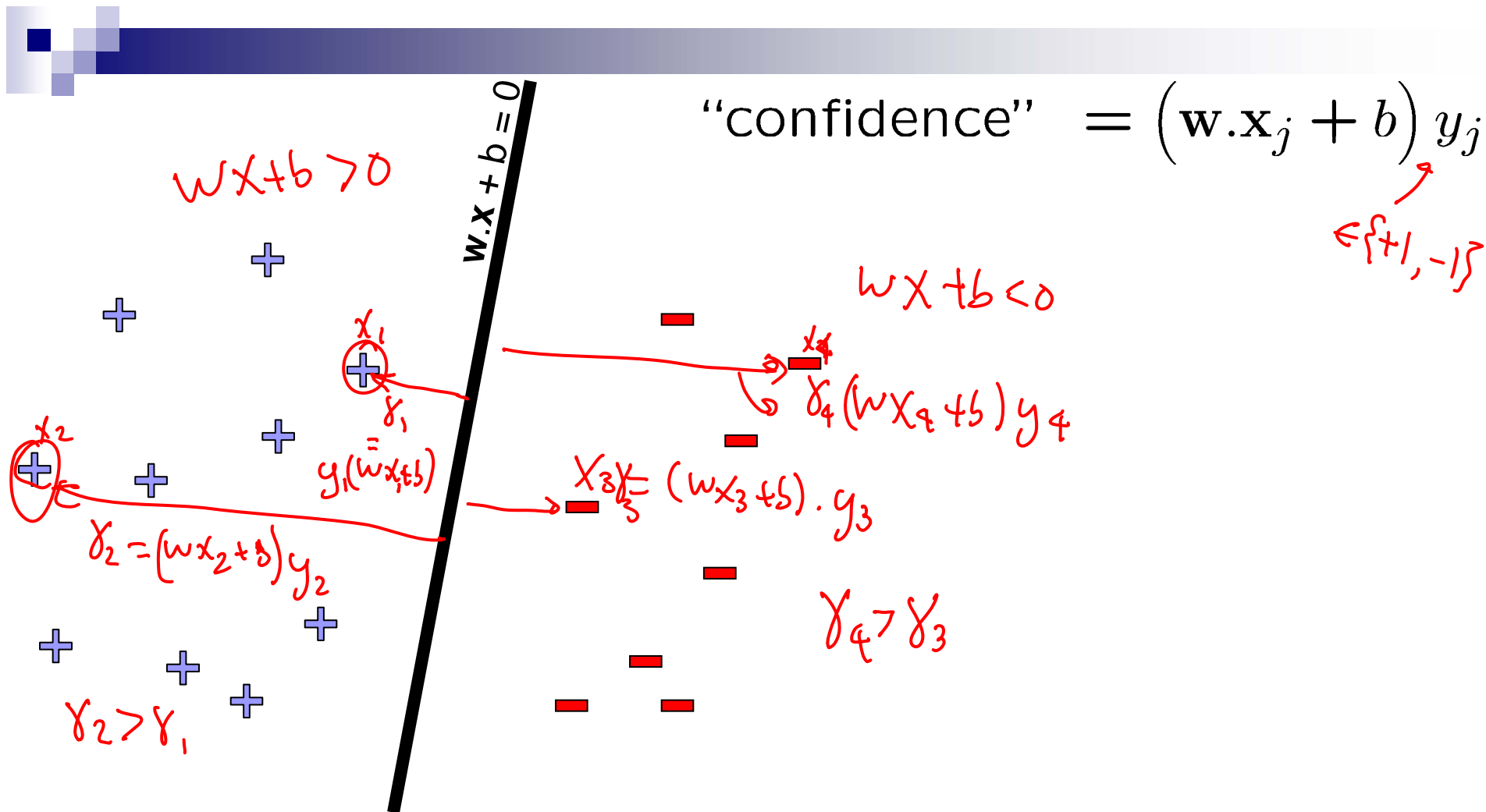
Carnegie Mellon University

February 21<sup>st</sup>, 2007

# Linear classifiers – Which line is better?

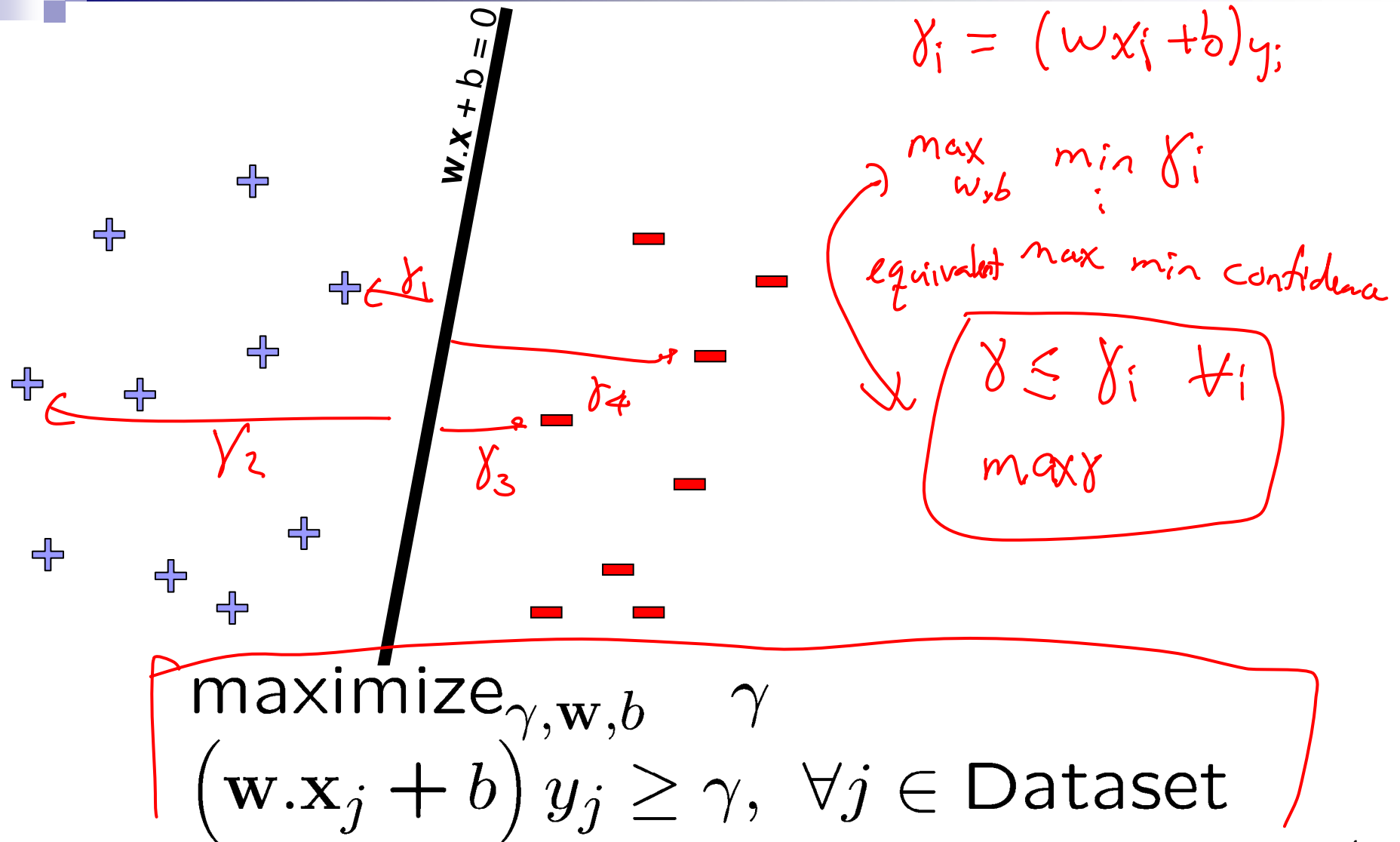


# Pick the one with the largest margin!

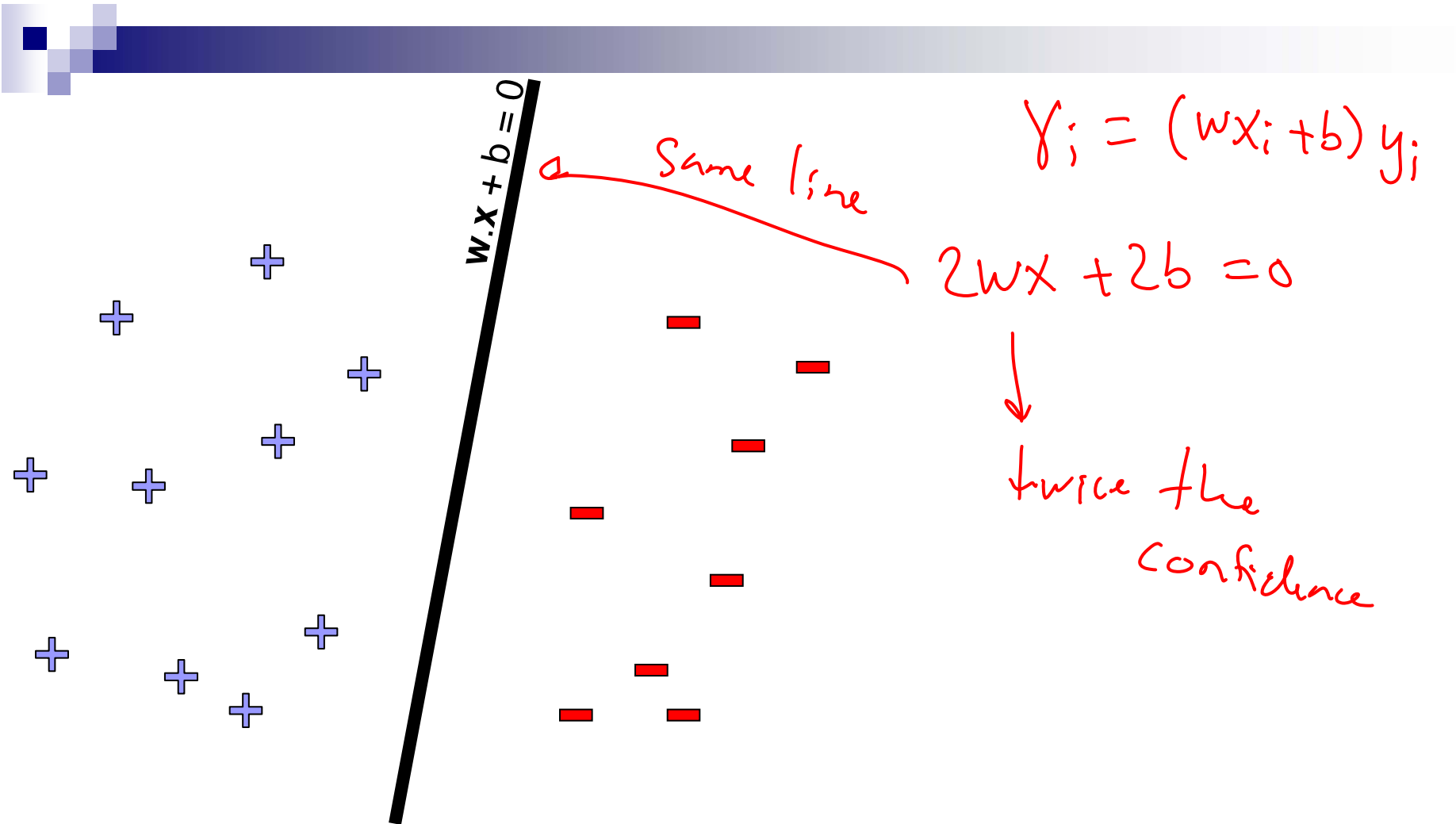


$$w \cdot x = \sum_j w^{(j)} x^{(j)}$$

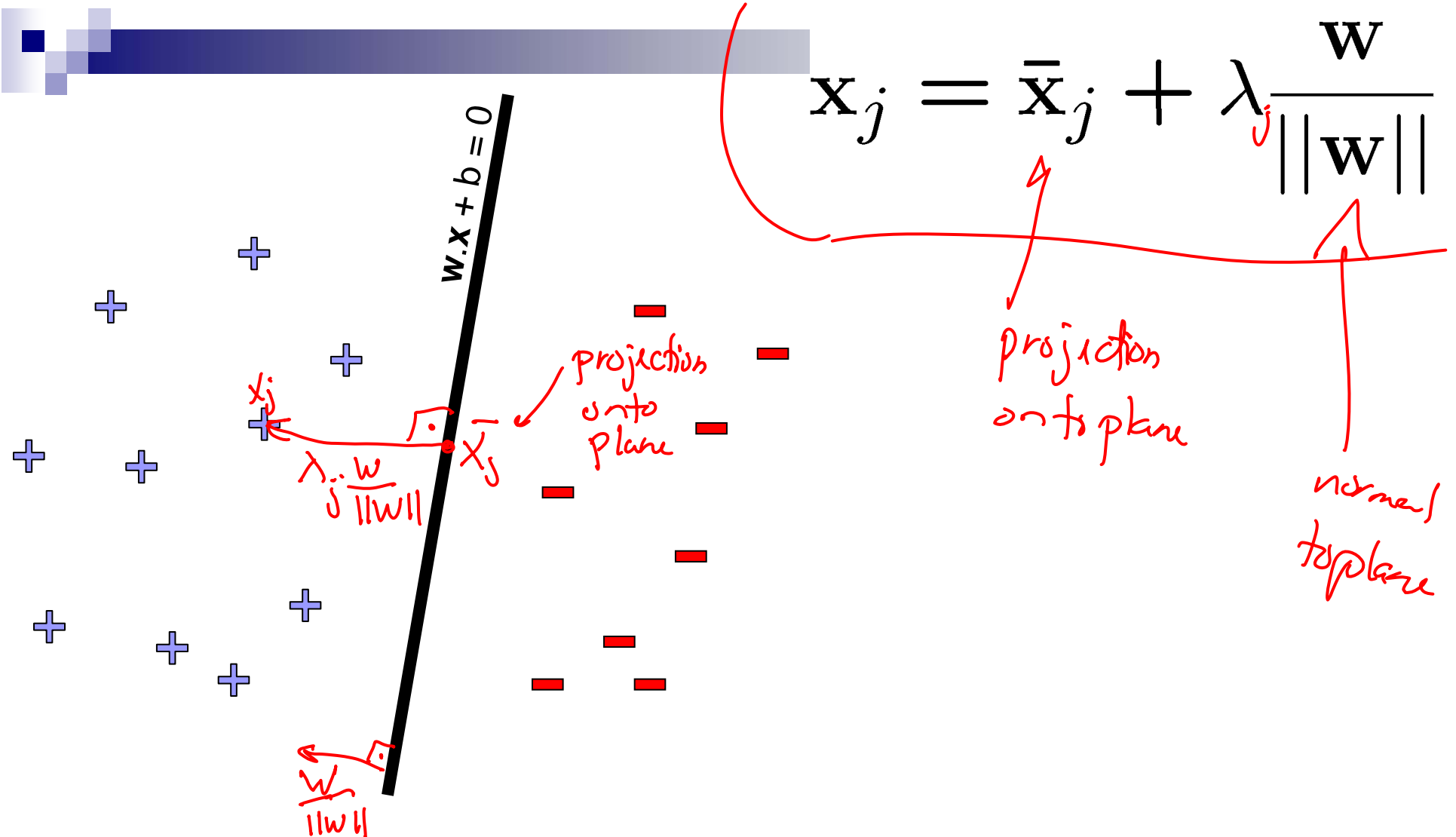
# Maximize the margin



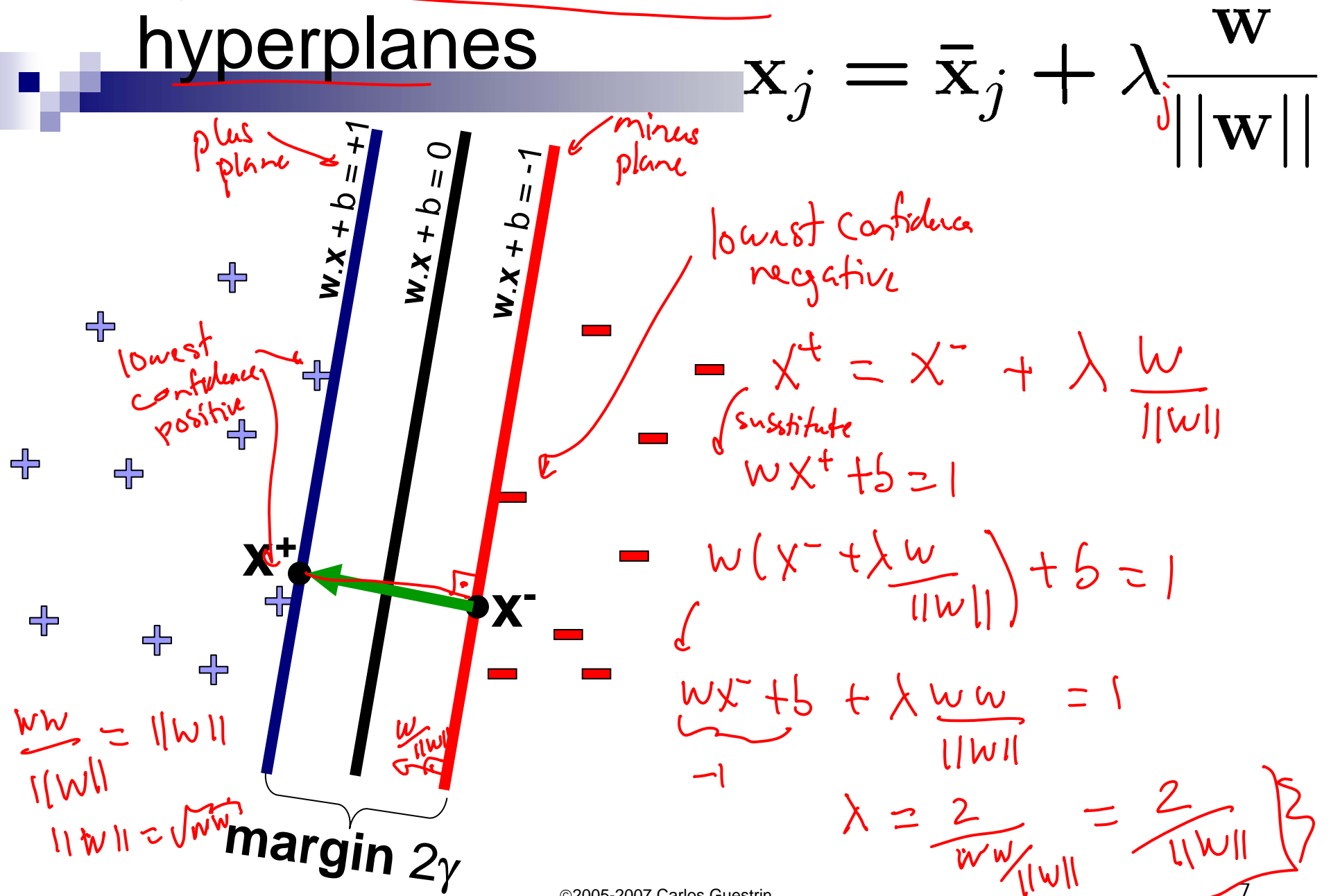
# But there are a many planes...



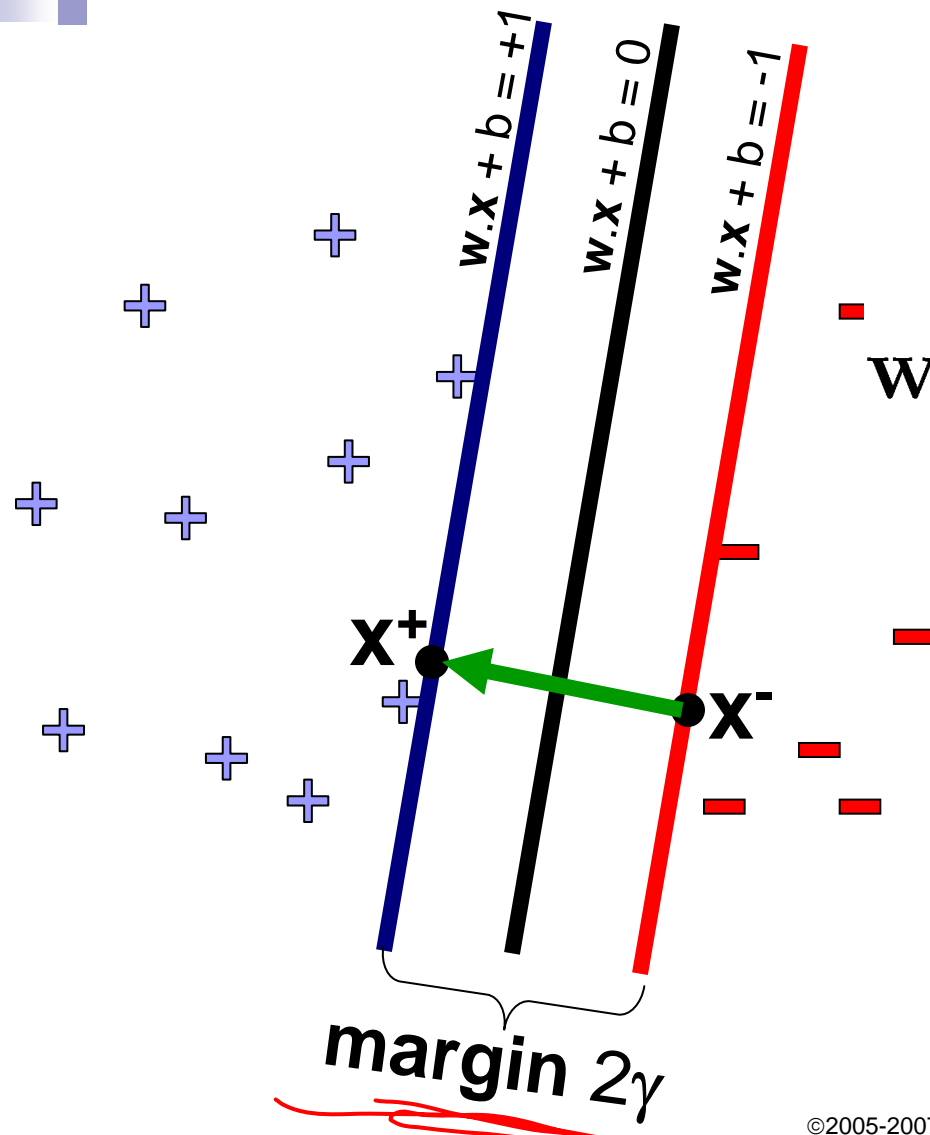
# Review: Normal to a plane



# Normalized margin – Canonical hyperplanes



# Normalized margin – Canonical hyperplanes



$$x^+ = x^- + \lambda w$$

*$x^+$  belongs to + plane*

$$w \cdot x^+ + b = 1$$

*substitute*

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

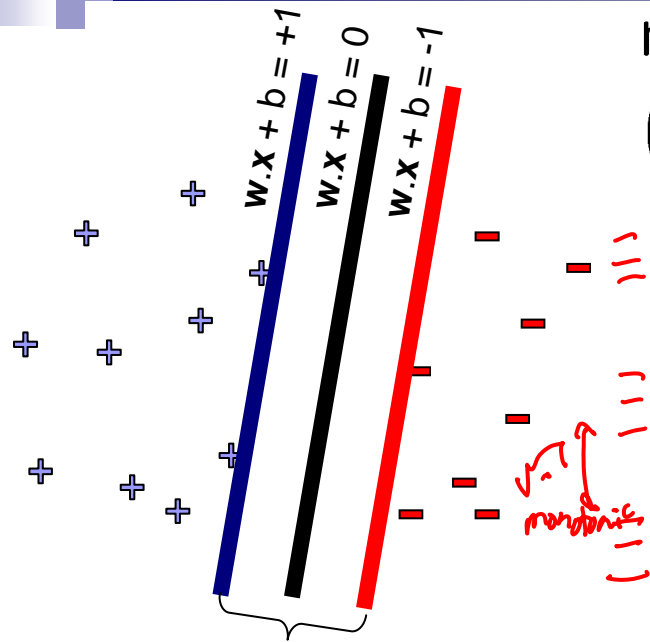
$$\lambda = \frac{2}{||w||}$$

$$\gamma = \frac{1}{\sqrt{w \cdot w}} \approx \frac{1}{||w||}$$



# Margin maximization using canonical hyperplanes

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



$$\text{maximize}_{\gamma, \mathbf{w}, b} \quad \gamma$$

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

$\gamma \leq 1$

$$\max \frac{1}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$

$$\min \sqrt{\mathbf{w} \cdot \mathbf{w}}$$

$$\min \mathbf{w} \cdot \mathbf{w}$$

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

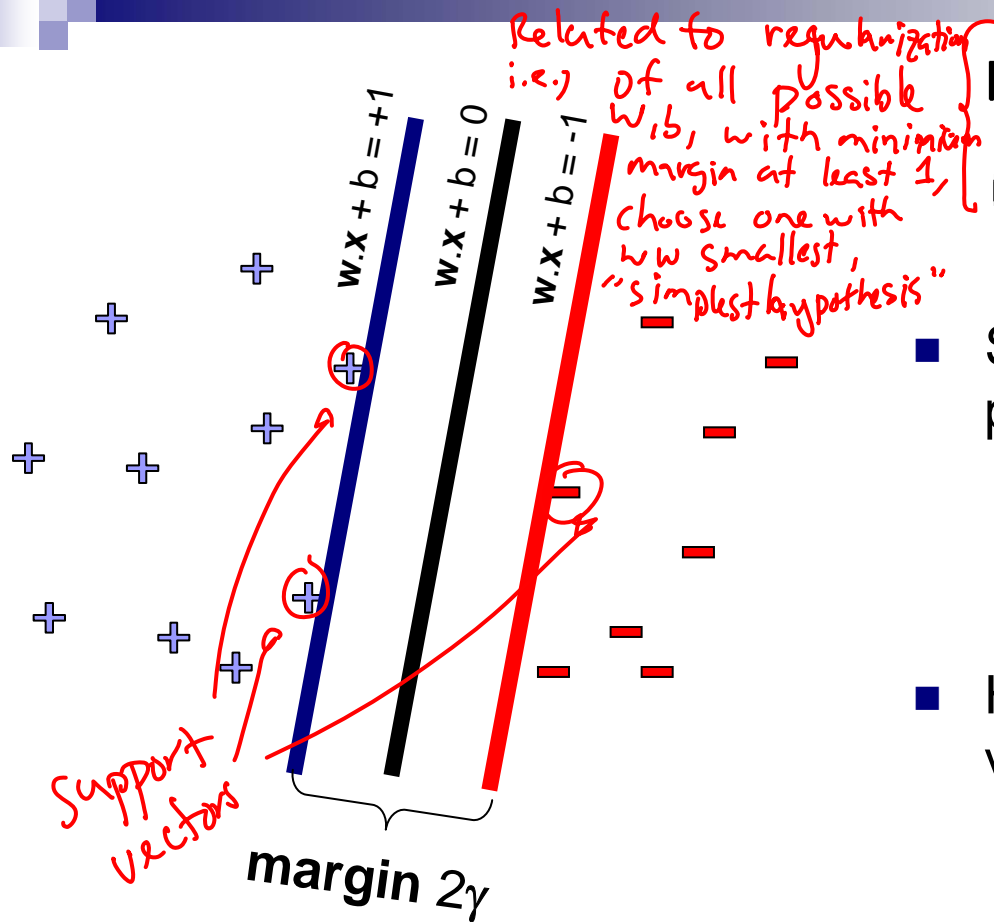
$\gamma \leq 1$

will maximize  $\gamma$  until hit  $\gamma=1$

$$\text{minimize}_{\mathbf{w}, b} \quad \mathbf{w} \cdot \mathbf{w}$$

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

# Support vector machines (SVMs)



$$\text{minimize}_{w,b} \quad w.w$$

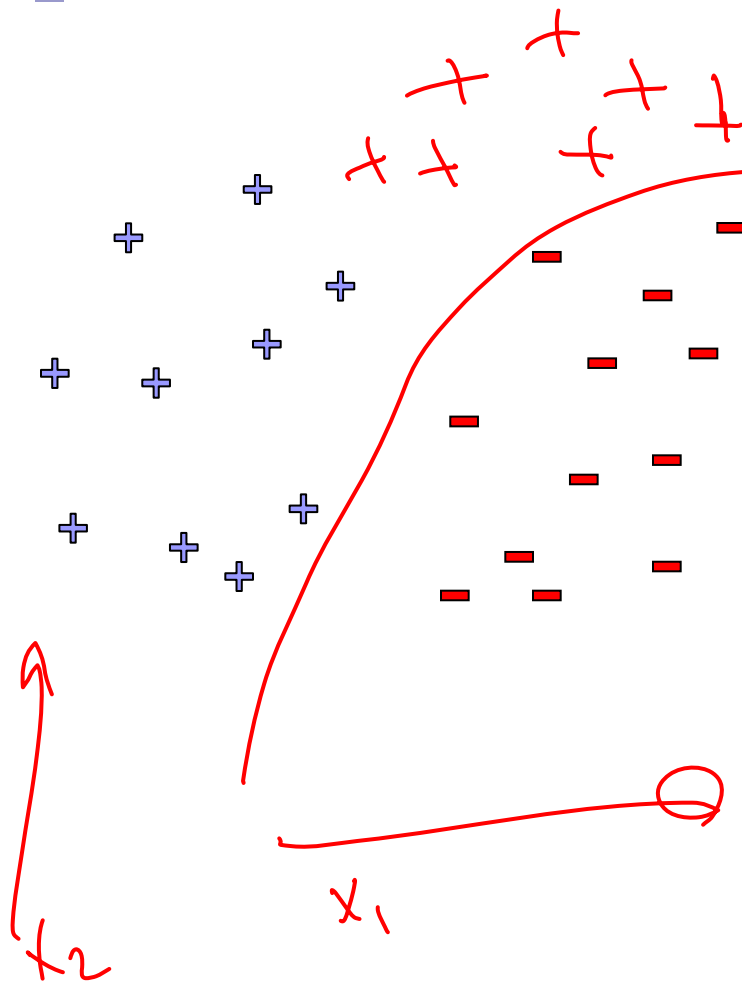
$$(w.x_j + b) y_j \geq 1, \quad \forall j$$

- Solve efficiently by quadratic programming (QP) → quadratic objective  
→ linear constraints
  - 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....**

features:

$$\{x_1, x_2, x_1^2, x_2^2, x_1x_2, \dots\}$$

e.g., polynomial

# What if the data is still not linearly separable?

$$\text{minimize}_{\mathbf{w}, b} \quad \mathbf{w} \cdot \mathbf{w} + C \cdot \#(\text{misclassified points})$$

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

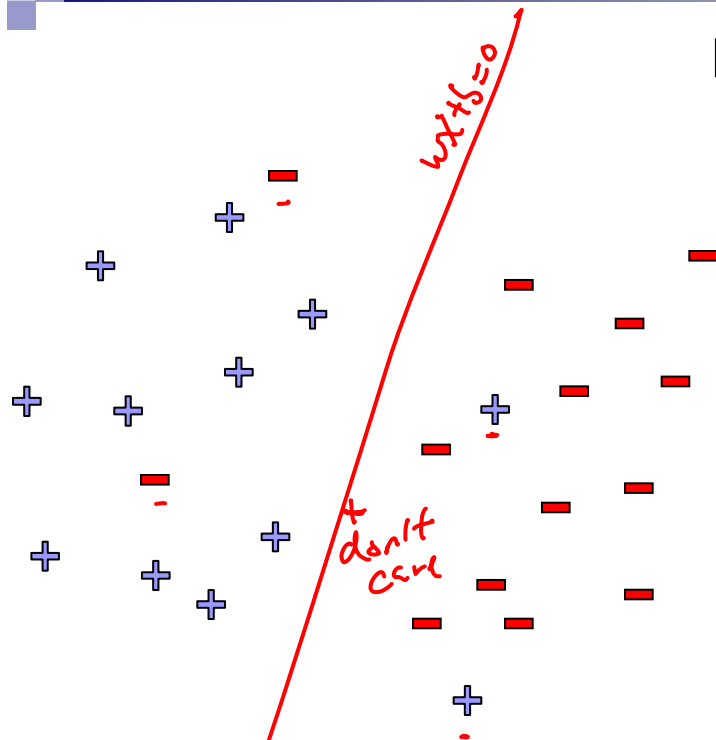
- Minimize  $\mathbf{w} \cdot \mathbf{w}$  and number of training mistakes

- Tradeoff two criteria?

$C$  is trade off parameter

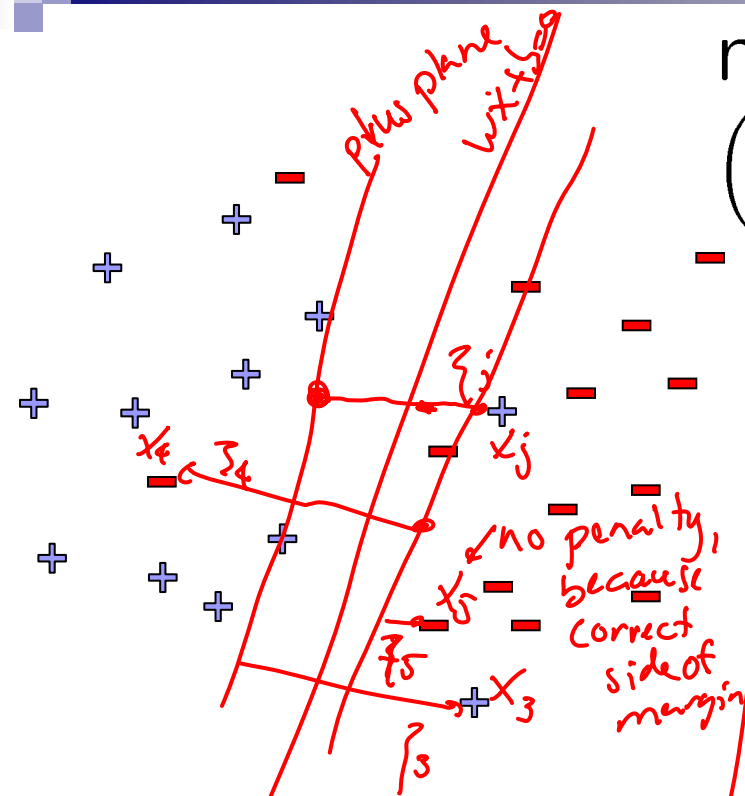
- Tradeoff  $\#(\text{mistakes})$  and  $\mathbf{w} \cdot \mathbf{w}$

- 0/1 loss
- Slack penalty  $C$
- Not QP anymore optimization hard
- Also doesn't distinguish near misses and really bad mistakes



penalize for misclassified points

# Slack variables – Hinge loss

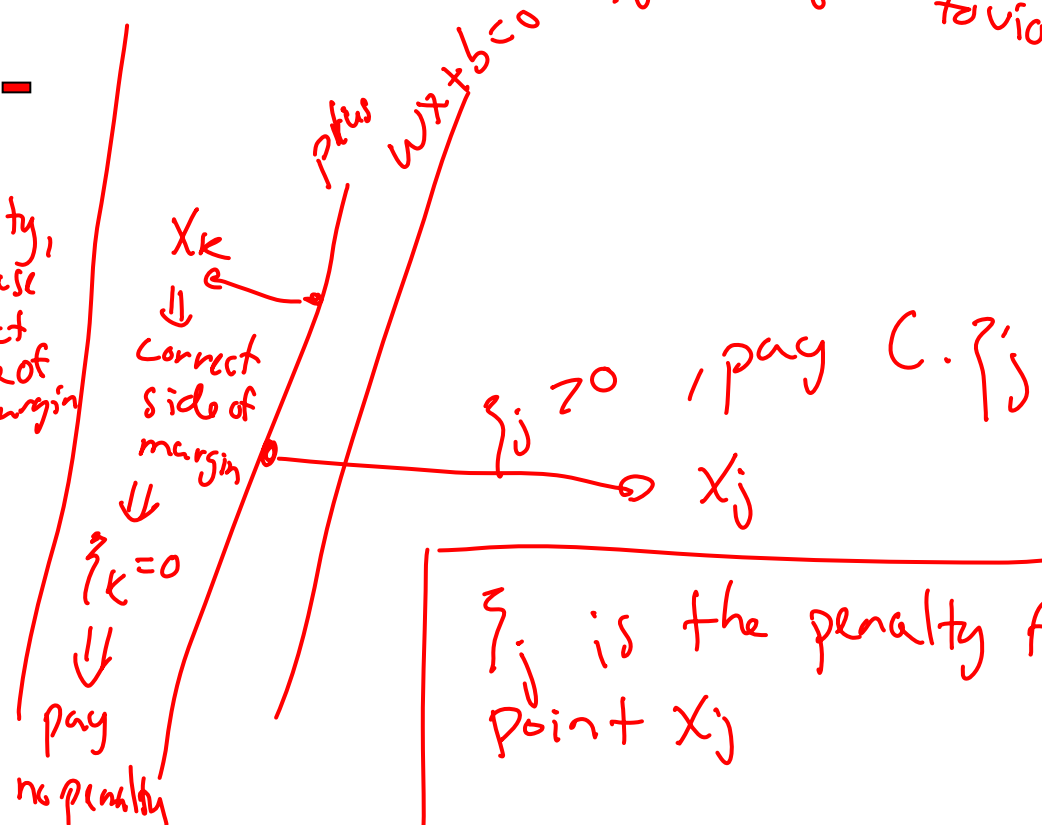


- If margin  $\geq 1$ , don't care
- If margin  $< 1$ , pay linear penalty

$$\text{minimize}_{w,b,\xi} \quad w \cdot w + C \sum_j \xi_j$$

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

$\xi_j \geq 0 \forall j$  ← allow you to violate margin



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

**SVM:**

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

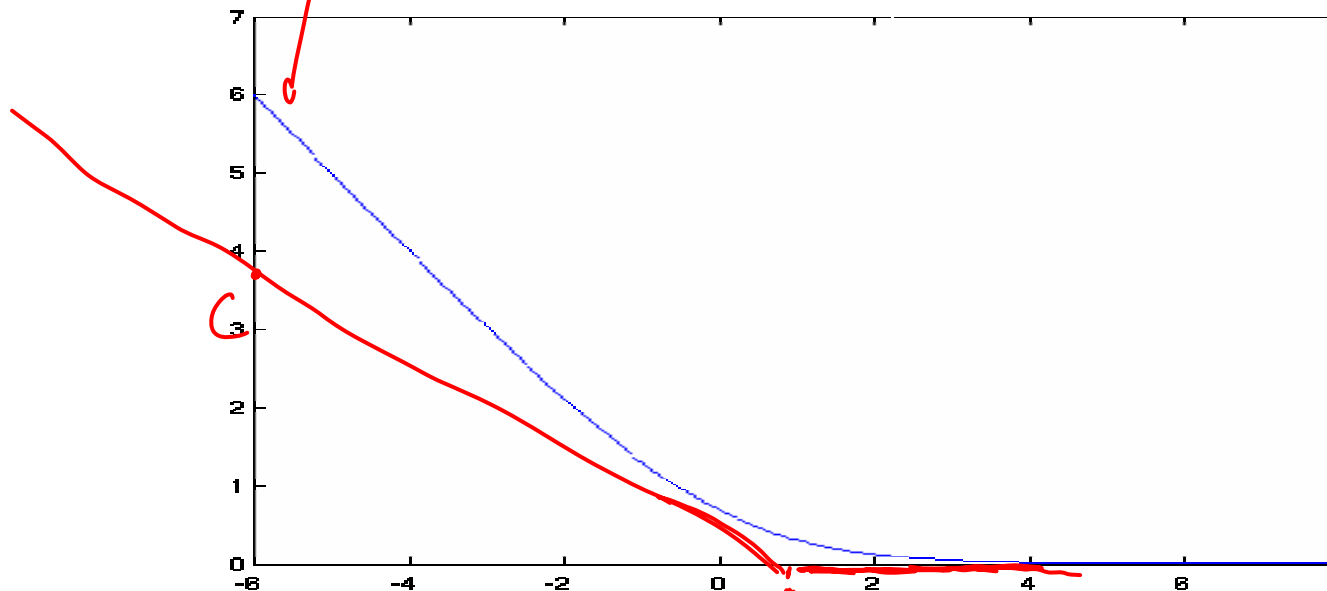
**Logistic regression:**

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

**Log loss:**

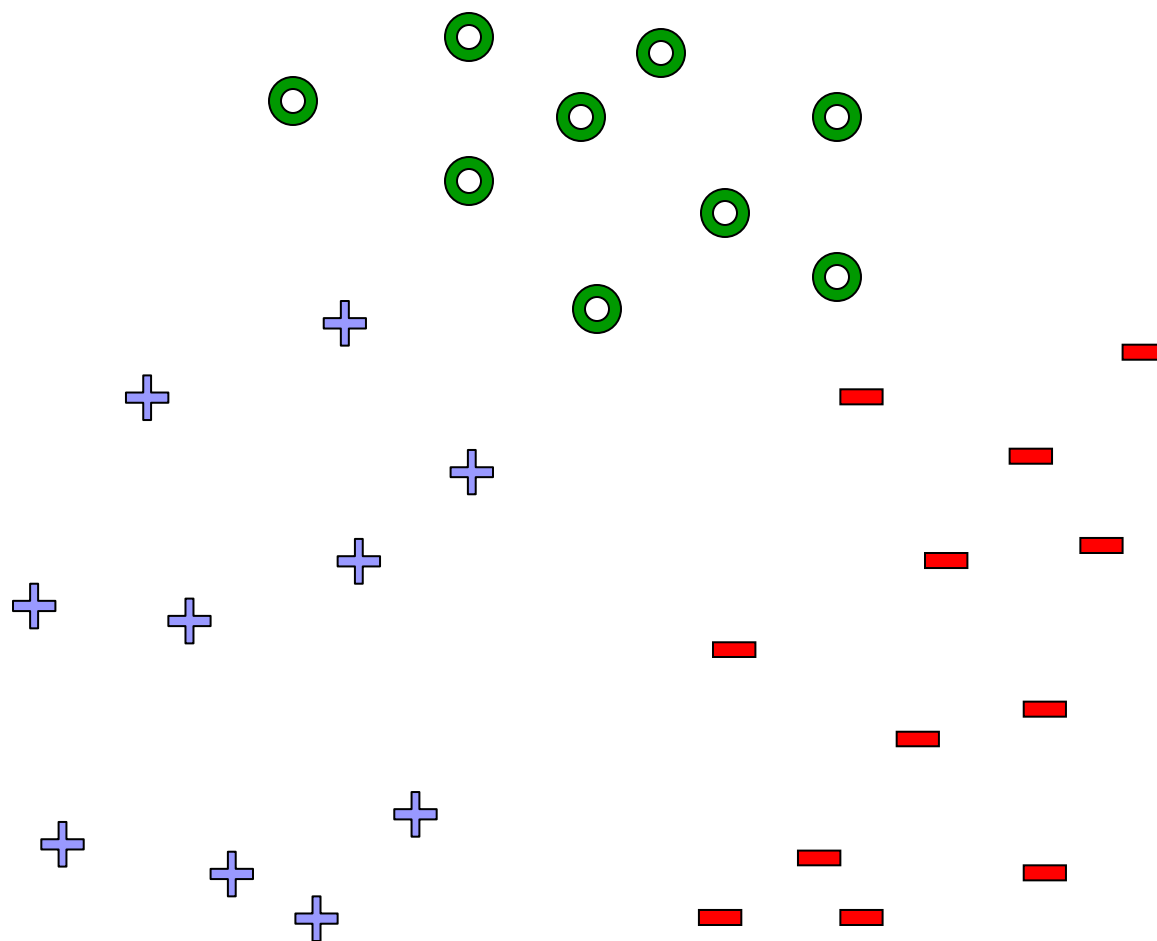
$$\min_{\mathbf{w}, b} -\ln P(Y = 1 \mid x, \mathbf{w}, b) = \ln(1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)})$$

$\ln(1 + e^{-z}) \leftarrow$  penalty function for LR



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

# What about multiple classes?





# One against All

Learn 3 classifiers:

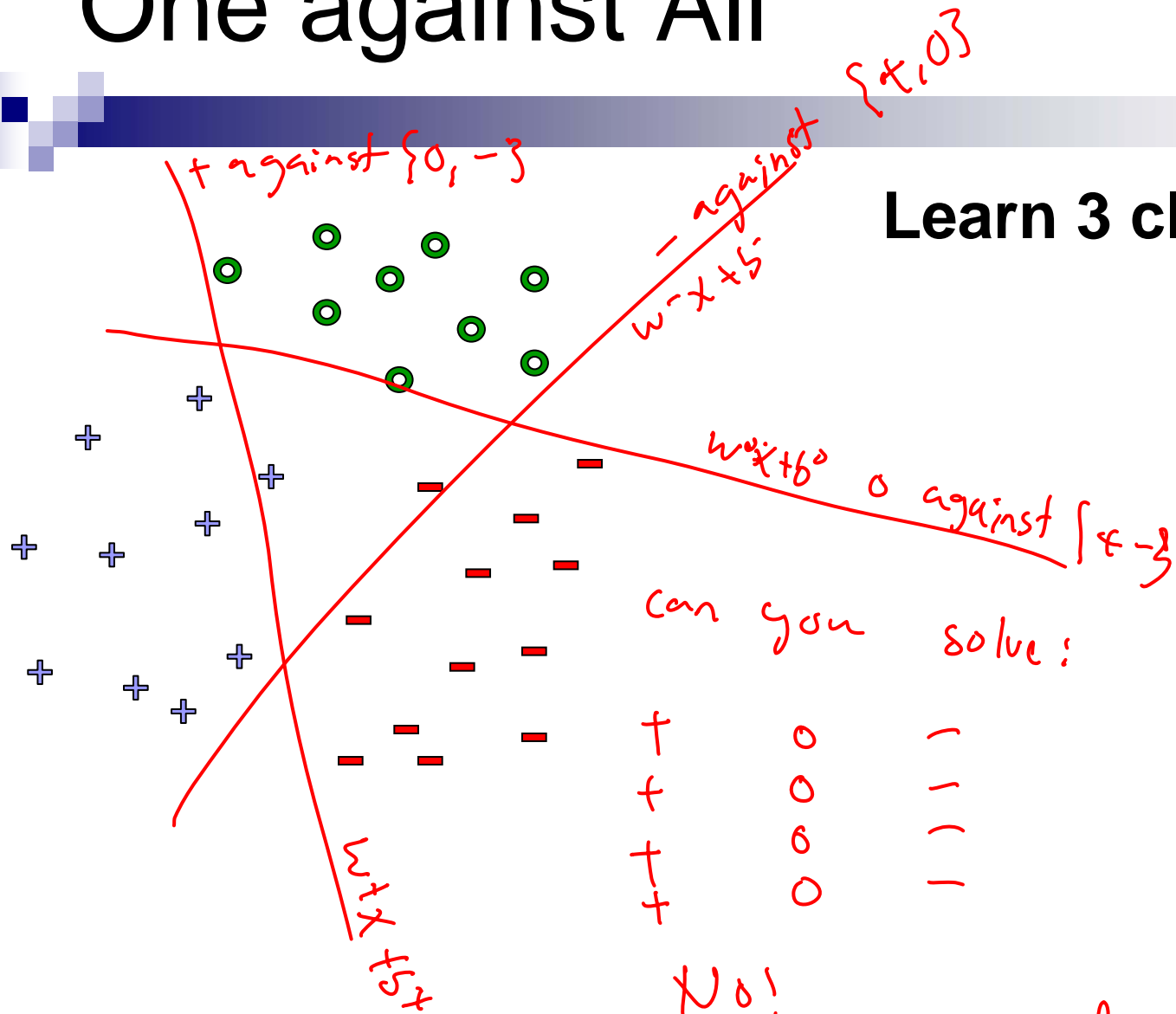
very easy to implement

new point  $x$   
e.g.

$x$  is  $t$  if

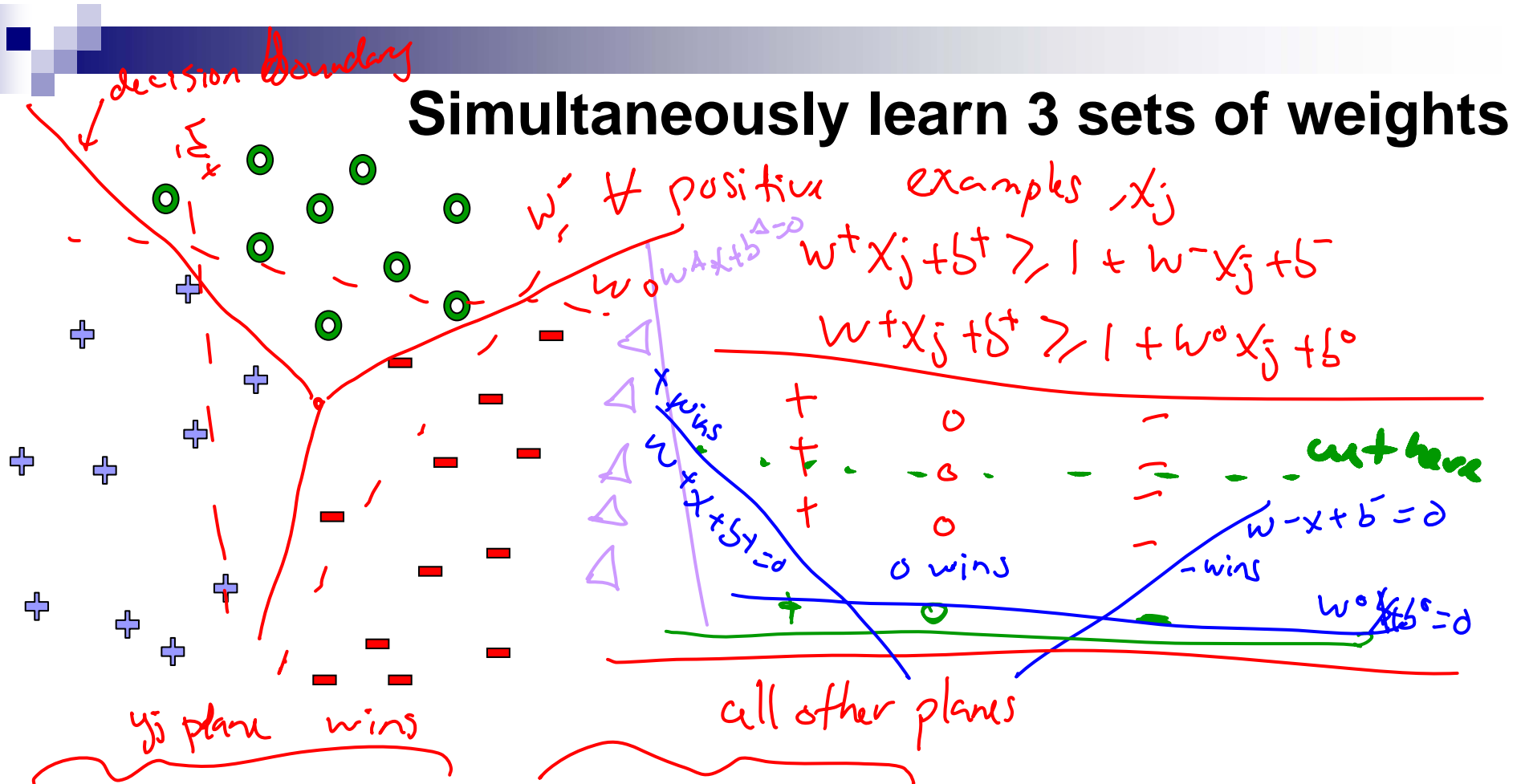
$$w^t x + b^t \geq w^+ x + b^+ \\ \geq w^- x + b^-$$

and so on



No!  
SVM  $w^+ x + b^+$  very confused!  
for

# Learn 1 classifier: Multiclass SVM

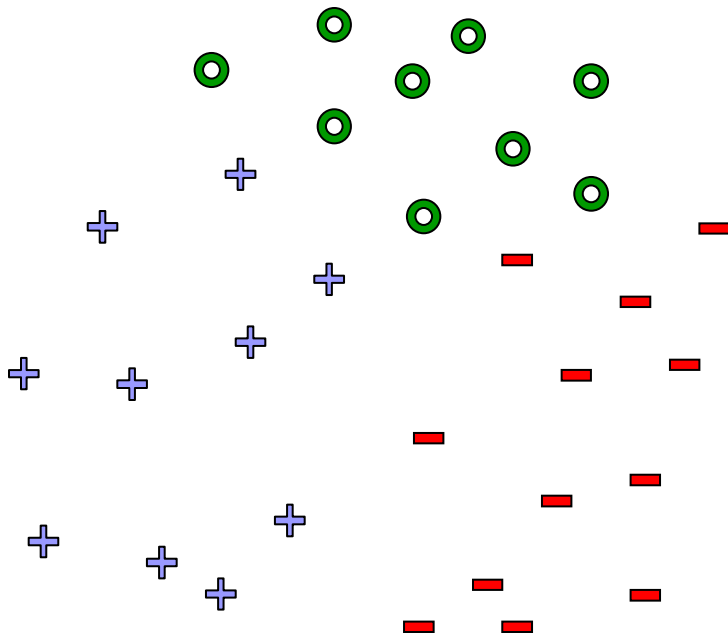


$$w^{(y_j)} \cdot x_j + b^{(y_j)} \geq w^{(y')} \cdot x_j + b^{(y')} + 1, \quad \forall y' \neq y_j, \quad \forall j$$

confidence of at least one

# 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