**Two SVM tutorials linked in class website (please, read both):**
- High-level presentation with applications (Hearst 1998)
- Detailed tutorial (Burges 1998)
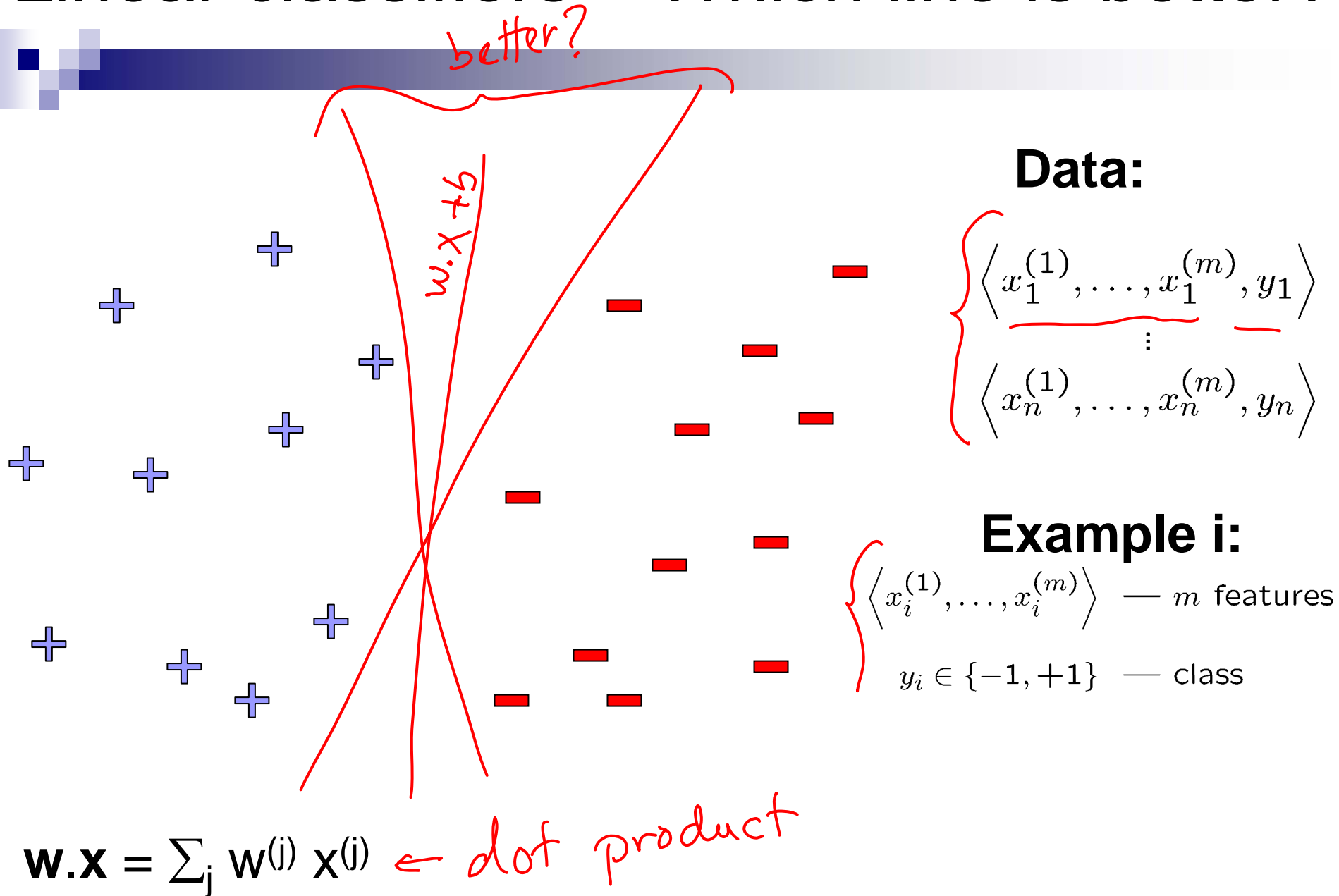
# Support Vector Machines
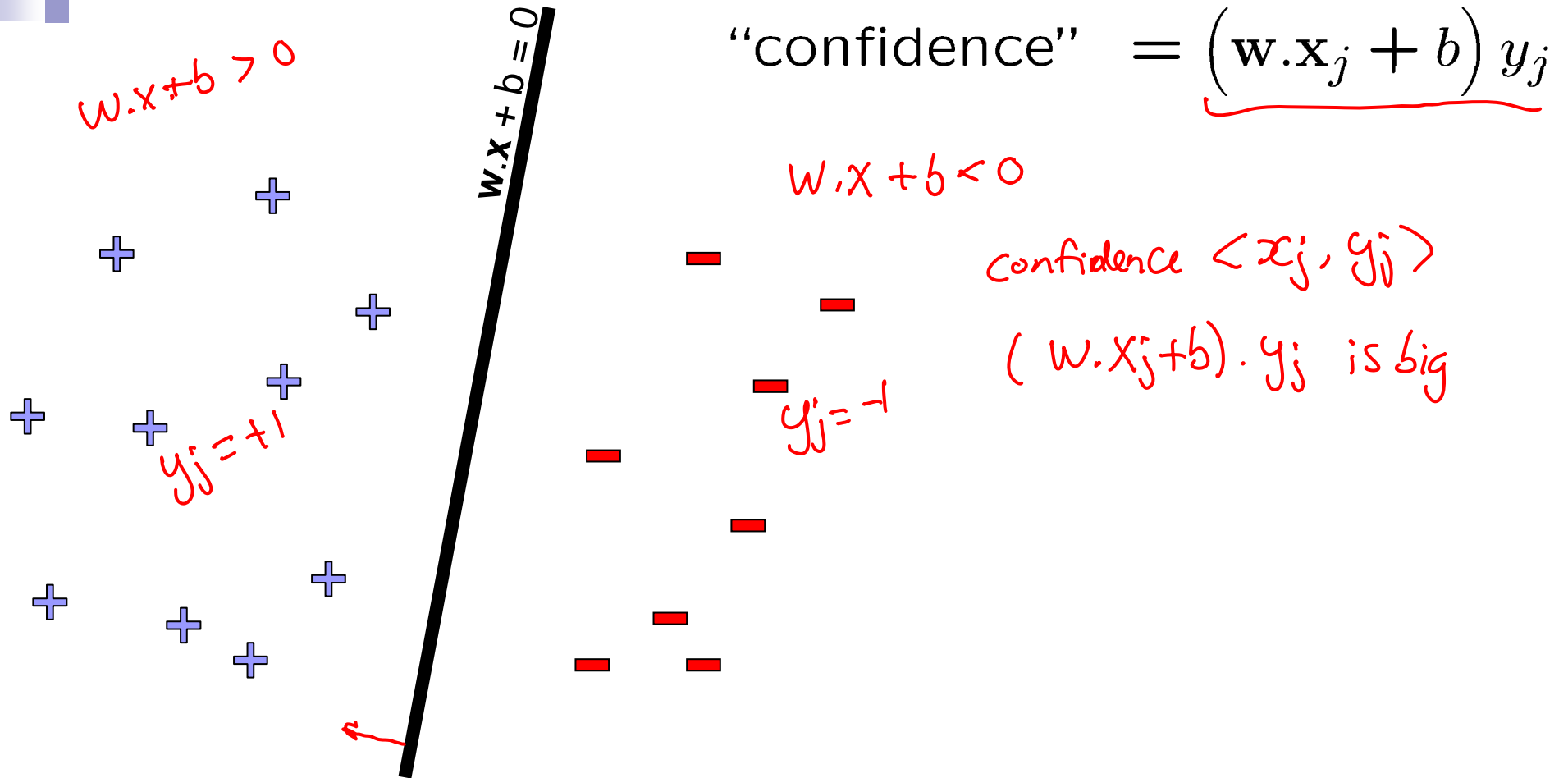
Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

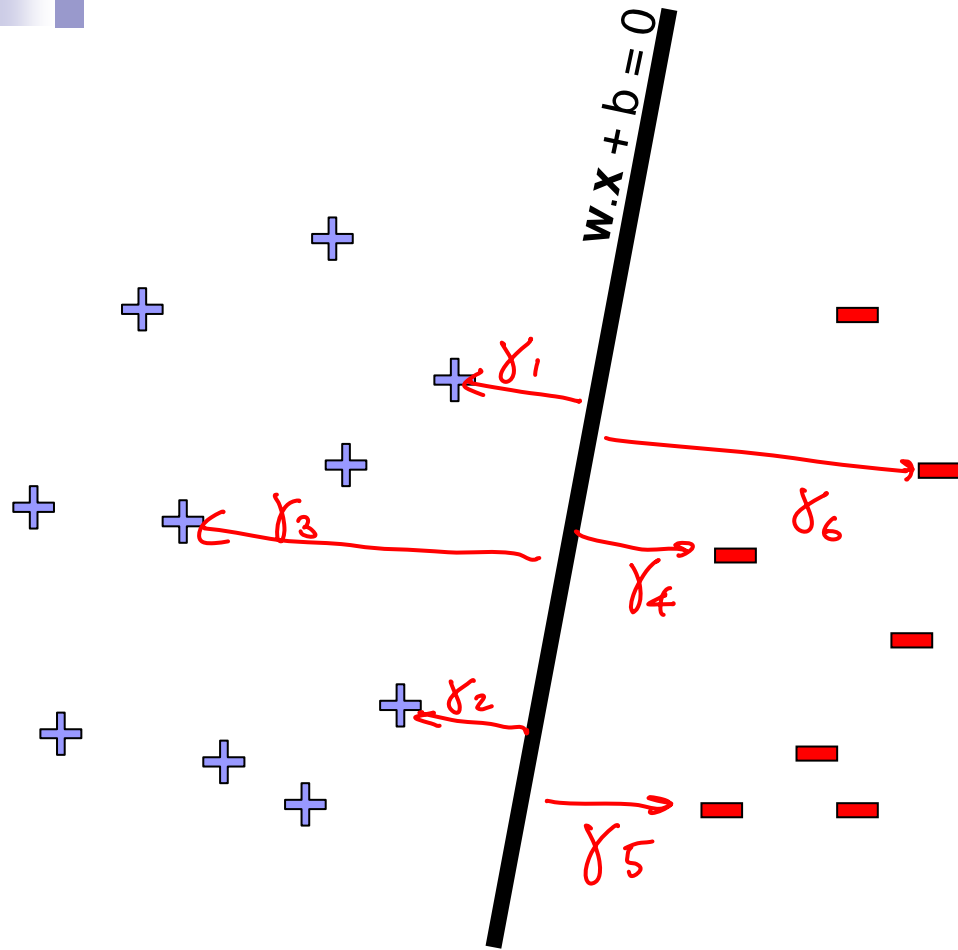February 16th, 2005

# Linear classifiers – Which line is better?

better?

$w.x + b$

**Data:**

$$\left\langle x_1^{(1)}, \ldots, x_1^{(m)}, y_1 \right\rangle$$
$$\vdots$$
$$\left\langle x_n^{(1)}, \ldots, x_n^{(m)}, y_n \right\rangle$$

**Example i:**

$$\left\langle x_i^{(1)}, \ldots, x_i^{(m)} \right\rangle \; — \; m \text{ features}$$

$$y_i \in \{-1, +1\} \; — \; \text{class}$$

$$\mathbf{w}.\mathbf{x} = \sum_j w^{(j)} x^{(j)} \quad \leftarrow \text{ dot product}$$

# Pick the one with the largest margin!

$w.x + b > 0$

**w.x + b = 0**

"confidence" $= \left(\mathbf{w}.\mathbf{x}_j + b\right) y_j$

$w.x + b < 0$

confidence $\langle x_j, y_j \rangle$

$(w.x_j + b) \cdot y_j$ is big

$y_j = -1$

$y_j = +1$

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

# Maximize the margin

$$w.x + b = 0$$

$$(x_1 . w + b) y_1 = \gamma_1$$

$$(x_j . w + b) . y_j = \gamma_j$$

$$(x_n . w + b) y_n = \gamma_n$$

$$\gamma = \min \gamma_j$$

$$\underset{w,b}{\text{maximize }} \gamma$$

$$(x_j w + b) y_j \geq \gamma, \forall j$$

$\gamma_1$

$\gamma_3$

$\gamma_4$

$\gamma_6$

$\gamma_2$

$\gamma_5$

# But there are a many planes…

w.x + b = 0

maximize $\gamma$

$2wx + 2b = 0$

$\frac{1}{2}w.x + \frac{1}{2}b = 0$

Redundancy

$\max_{w,b} \gamma$

$(x_j w + b) y_j \geq \gamma \; \forall j$

# *Review*: Normal to a plane

$$x_j = \bar{x_j} + \lambda \cdot w$$

any point in the space

$$\text{w.x} + \text{b} = 0$$

$\bar{x_j}$

$x_j$

$\dfrac{w}{\|w\|}$

normal

# Normalized margin – Canonical hyperplanes

Plus plane

Minus plane

w.x + b = +1

w.x + b = 0

w.x + b = -1

$X^+$

$X^-$

$\frac{w}{||w||}$

$\lambda w$

**margin** $\gamma$

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

↳ Plug into Plus plane:

$w X^+ + b = +1$

$w(X^- + \lambda w) + b = +1$

$\underbrace{w X^- + b} + \lambda w.w = +1$

$X^-$ in Minus plane

$-1$

$-1 + \lambda w.w = +1$

$\lambda = \dfrac{2}{w.w}$

$\gamma = \dfrac{2||w||}{w.w}$

$||w|| = \sqrt{w.w}$

Margin of conf.

Max: $\gamma = \dfrac{2}{\sqrt{\mathbf{w.w}}}$

# Margin maximization using canonical hyperplanes

$$\frac{2}{\sqrt{w.w}} \leq 1$$

$$\text{w.x} + b = +1$$
$$\text{w.x} + b = 0$$
$$\text{w.x} + b = -1$$

**margin** $\gamma$

$$\text{Max } \frac{2}{\sqrt{w.w}}$$

$$\equiv$$

$$\text{Min } \frac{\sqrt{w.w}}{2}$$

$$\equiv$$

$$\text{Min } w.w$$

Because
$w.w > 0$

$\sqrt{\cdot}$ is
monotonic

$$\text{Maximize}_{w,b} \quad \gamma$$

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

subject $\gamma \leq 1$

Might as well $\gamma = 1$

$$\text{Maximize}_{w,b} \quad \frac{2}{\sqrt{w.w}}$$

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

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

# Support vector machines (SVMs)



Support vector

$\mathbf{w}.\mathbf{x} + b = +1$

$\mathbf{w}.\mathbf{x} + b = 0$

$\mathbf{w}.\mathbf{x} + b = -1$

**margin** $\gamma$

minimal set of points

$$\text{minimize}_{\mathbf{w}} \quad \mathbf{w}.\mathbf{w}$$

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

- Solve efficiently by quadratic programming (QP)
  - Well-studied solution algorithms

- Hyperplane defined by support vectors

# What if the data is not linearly separable?

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

$2d$ $\langle x^{(1)}, x^{(2)}, y \rangle$

$\downarrow$ Feed SVM:

$\langle x^{(1)}, x^{(1)}.x^{(1)}, x^{(2)}, x^{(2)}.x^{(2)}, y \rangle$

polynomial features

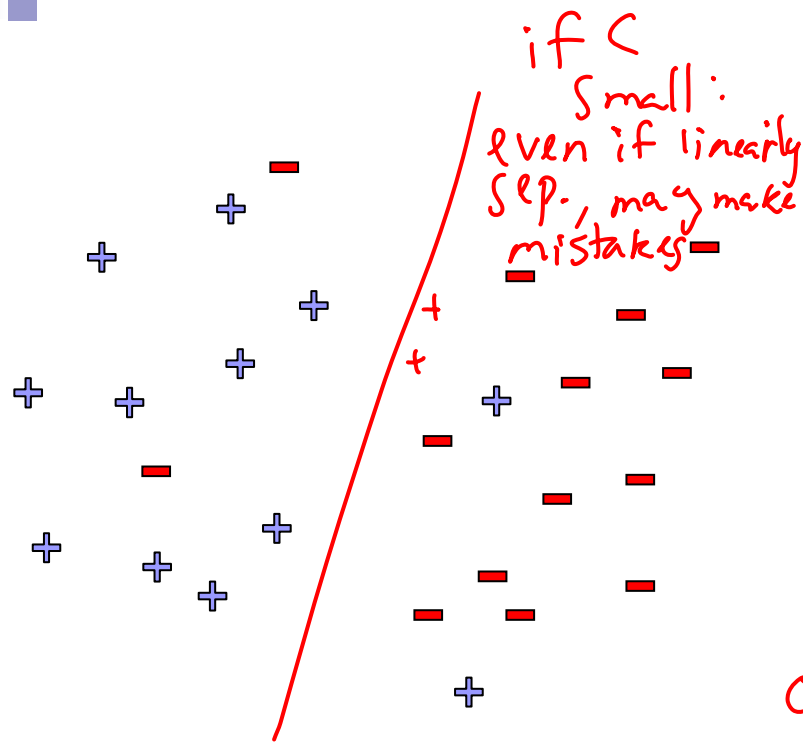not linearly seperable
$\Rightarrow \nexists$ hyperplane
$\gamma > 0$

# What if the data is still not linearly separable?

$$\text{minimize}_{\mathbf{w}} \quad \mathbf{w}.\mathbf{w} + C\,(\#\text{mistakes})$$

$$\left(\mathbf{w}.\mathbf{x}_j + b\right)y_j \geq 1 \qquad , \forall j$$

*{j : did I j wrong?}*

- Minimize **w.w** and number of training mistakes
  - Tradeoff two criteria?

- Tradeoff #(mistakes) and **w.w**
  - 0/1 loss
  - Slack penalty $C$
  - Not QP anymore
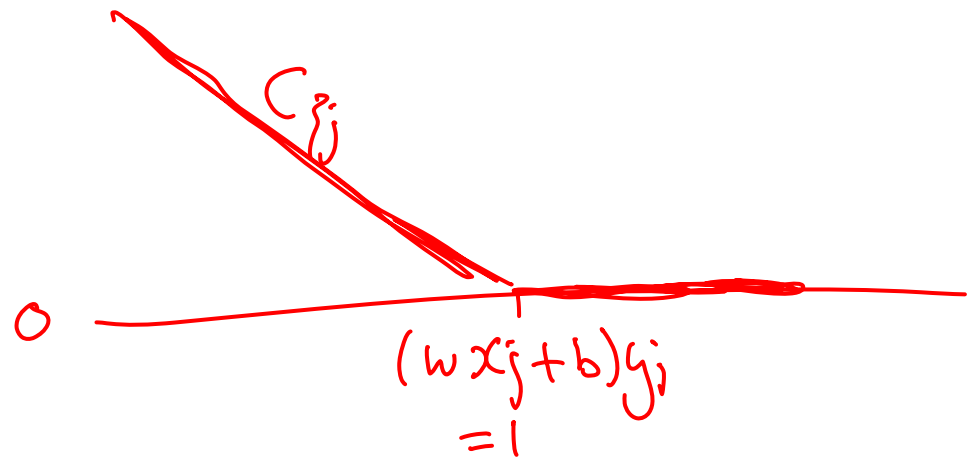  - Also doesn't distinguish near misses and really bad mistakes

$(wx_j+b)y_j$
$=1$

# Slack variables – Hinge loss



$$\text{minimize}_{\mathbf{w}} \quad \mathbf{w}.\mathbf{w} + C \sum_j \xi_j$$

$$\left( \mathbf{w}.\mathbf{x}_j + b \right) y_j \geq 1 - \xi_j \,, \forall j$$

$$\xi_j \geq 0$$

if $C$ small:
even if linearly
sep., may make
mistakes

$C \xi_j$

$0$

$(w x_j + b) y_j$
$= 1$

- If margin $\geq 1$, don't care $\Rightarrow$ $\xi_j = 0$, pay nothing

- If margin $< 1$, pay linear penalty $\Rightarrow$ $\xi_j > 0$, and pay $C \cdot \xi_j$

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

**SVM:** → *linear classifiers* → **Logistic regression:**

$$\text{minimize}_{\mathbf{w}} \quad \mathbf{w}.\mathbf{w} + C \sum_j \xi_j$$
$$\left(\mathbf{w}.\mathbf{x}_j + b\right) y_j \geq 1 - \xi_j, \ \forall j$$
$$\xi_j \geq 0, \ \forall j$$
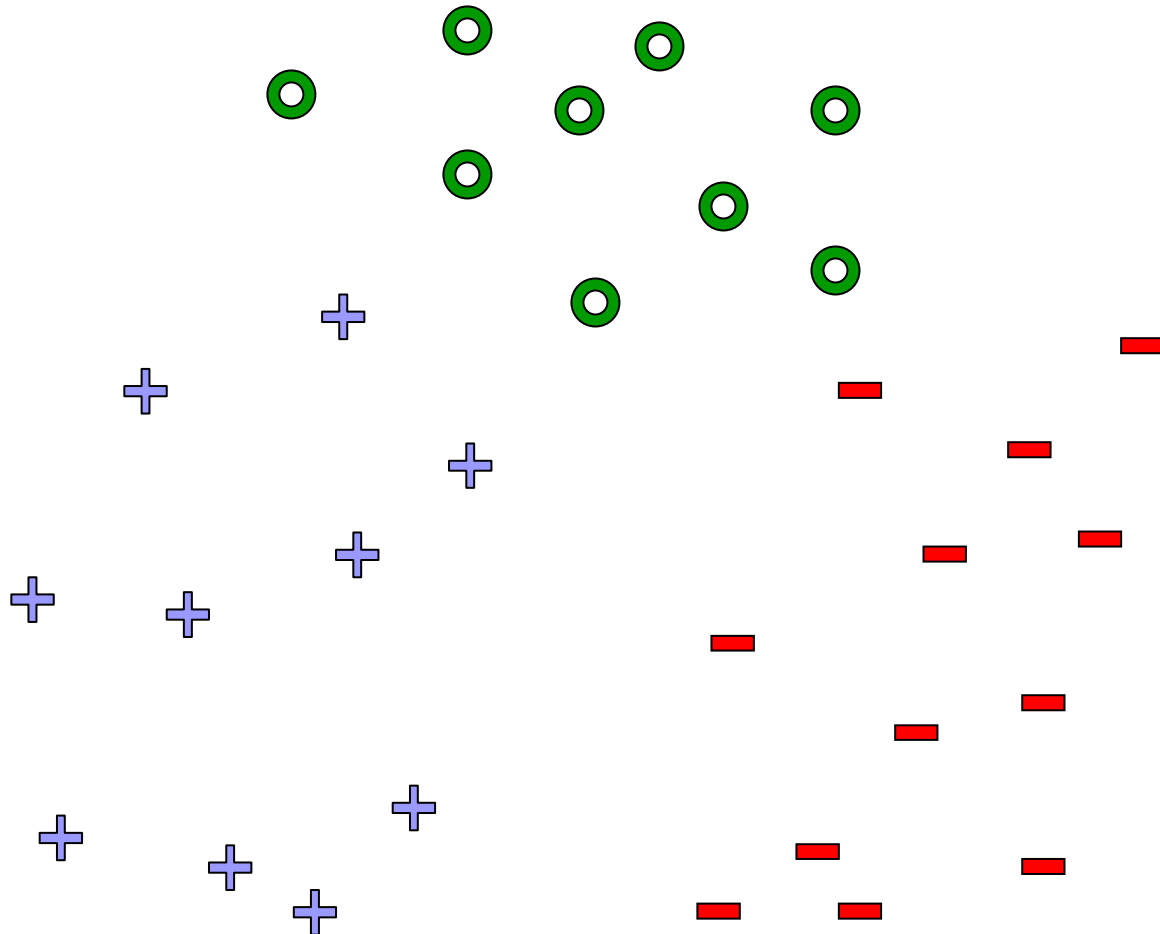
$$P(Y = 1 \mid x, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w}.\mathbf{x}+b)}}$$

**Log loss:**

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

Hinge Loss:

c



log loss

hinge loss

C

0/1 loss

$wx+b=0$

# What about multiple classes?

# One against All



Learn 3 classifiers:

$+$ versus $\{o, -\}$

$-$ versus $\{o, +\}$

$o$ versus $\{-, +\}$

classify x
classifier with
highest confidence

# Learn 1 classifier: Multiclass SVM

$$\begin{cases} w^{(+)}, b^{(+)} \\ w^{(0)}, b^{(0)} \\ w^{(-)}, b^{(-)} \end{cases}$$

## Simultaneously learn 3 sets of weights



For $+$ examples:

$$(w^{(+)} x_j + b^{(+)}) \geq 1 + w^{(-)} x_j + b^{(-)}$$

$$w^{(+)} x_j + b^{(+)} \geq 1 + w^{(0)} x_j + b^{(0)}$$

For $-$ examples
$w^{(-)}, b^{(-)}$ win

For $0$ examples
$w^{(0)}, b^{(0)}$ win

$$\mathbf{w}^{(y_j)}.\mathbf{x}_j + b^{(y_j)} \geq \mathbf{w}^{(y')}.\mathbf{x}_j + b^{(y')} + 1, \ \forall y' \neq y_j, \ \forall j$$
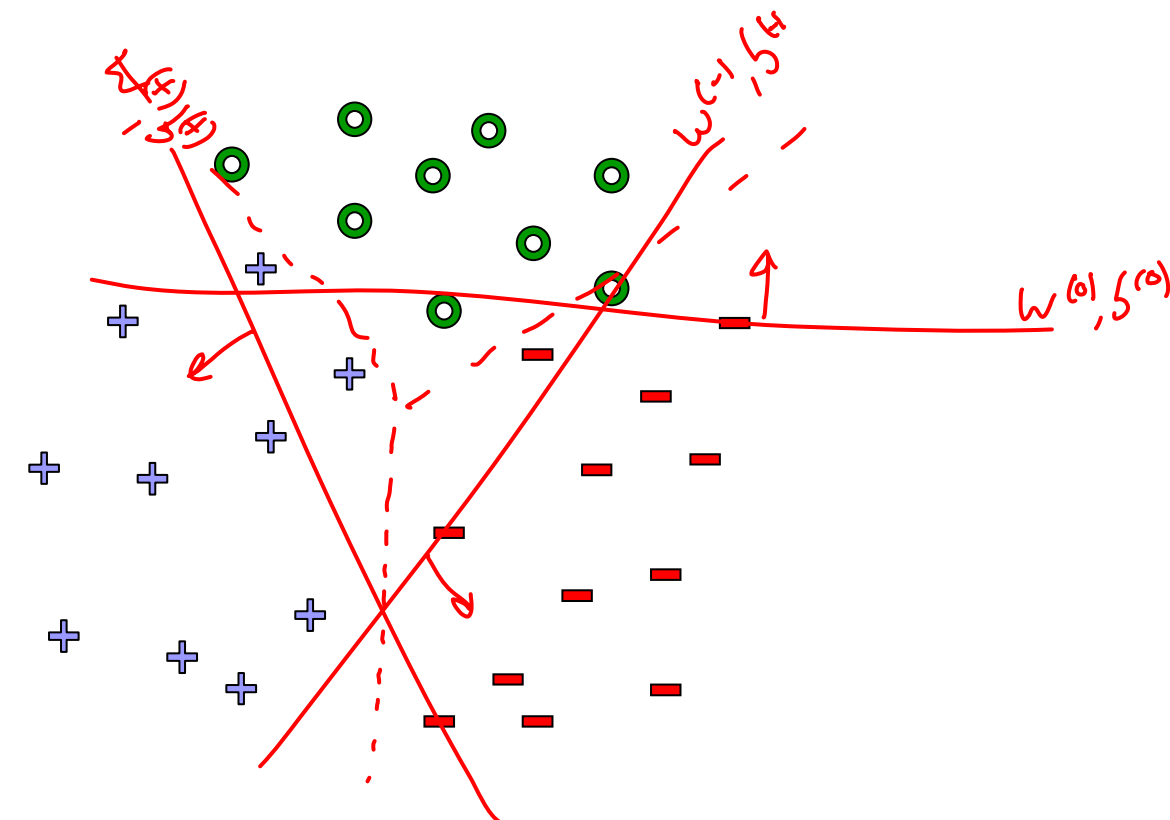
# Learn 1 classifier: Multiclass SVM

possible classes

$$\text{minimize}_{\mathbf{w}} \quad \sum_{y} \mathbf{w}^{(y)} . \mathbf{w}^{(y)} + C \sum_{j} \xi_{j}$$

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

$$\xi_j \geq 0, \ \forall j$$

# 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

# Acknowledgment

■ SVM applet:

☐ http://www.site.uottawa.ca/~gcaron/applets.htm