

# Support Vector Machines

Adopted from slides by

Jure Leskovec, Anand Rajaraman, Jeff Ullman, <http://www.mmds.org>

and Alexander Ihler

# Linear classifiers

- Classifier model:  $f(x; \theta) = T(\theta x)$
- Loss function:

**Class  $y = \{0, 1\}$**

$$J(\theta) = \frac{1}{m} \sum_i \left( y^{(i)} \phi(\theta x^{(i)}) + (1 - y^{(i)}) \phi(-\theta x^{(i)}) \right)$$

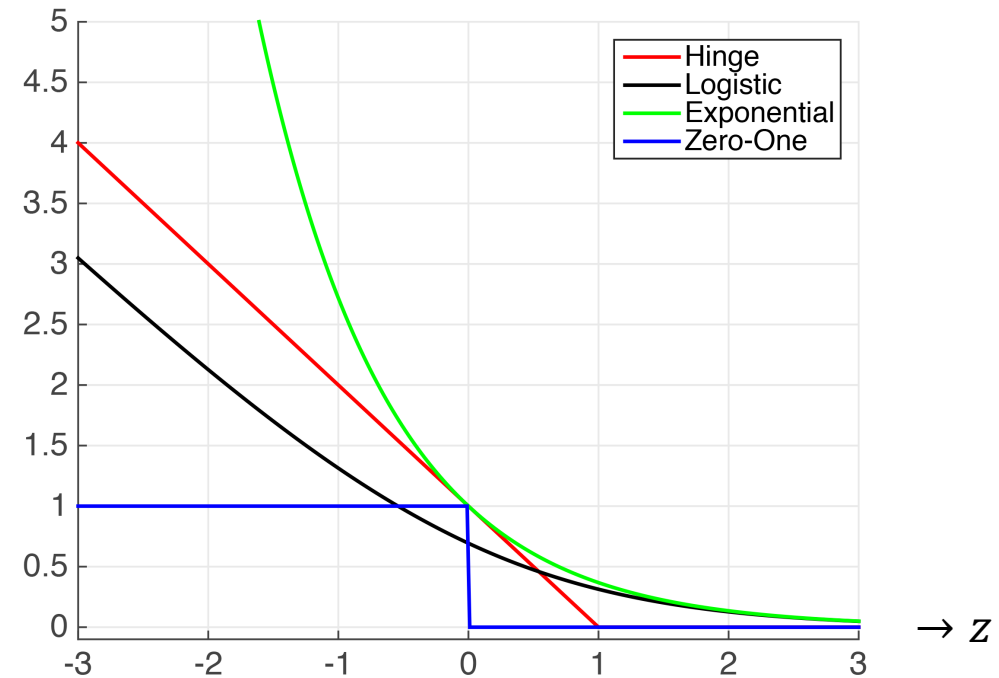
**Class  $y = \{-1, 1\}$**

$$J(\theta) = \frac{1}{m} \sum_i \phi \left( y^{(i)} (\theta x^{(i)}) \right)$$

# Surrogate loss functions

- 0-1:  $\mathcal{L}(z) = \mathbf{1}[z < 0]$
- Logistic:  $\mathcal{L}(z) = -\log \sigma(z)$
- Exponential:  $\mathcal{L}(z) = e^{-\beta z}$
- Hinge:  $\mathcal{L}(z) = \max\{0, 1 - z\}$
- ...

Logistic  
regression

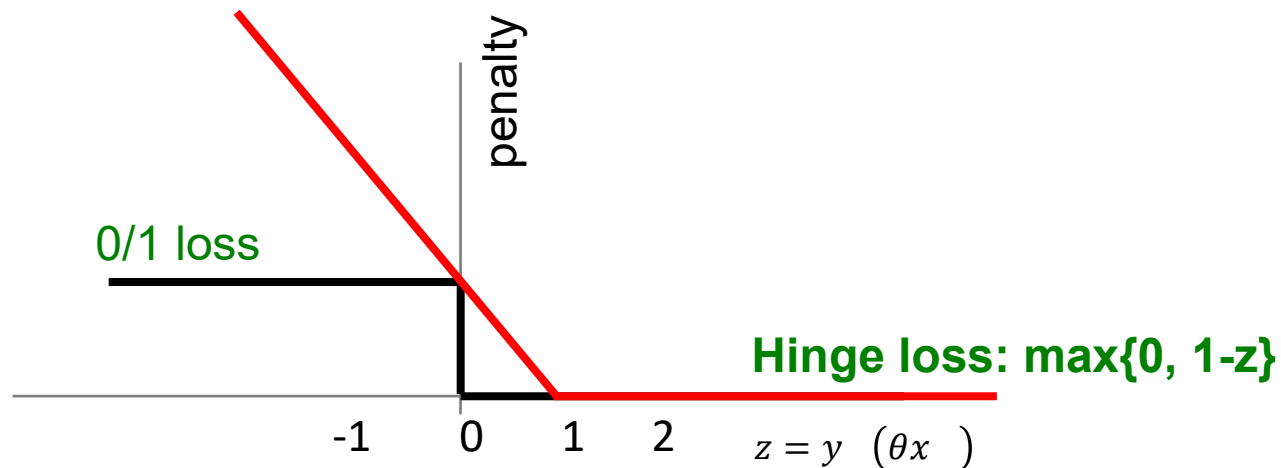


# Support Vector Machines (SVM)

- Classifier:  $f(x; \theta) = \text{sign}(\theta x)$
- Loss function:

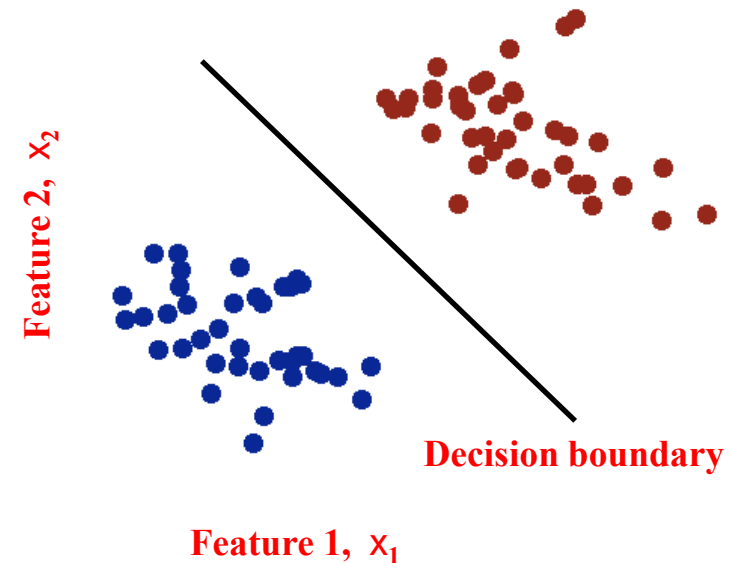
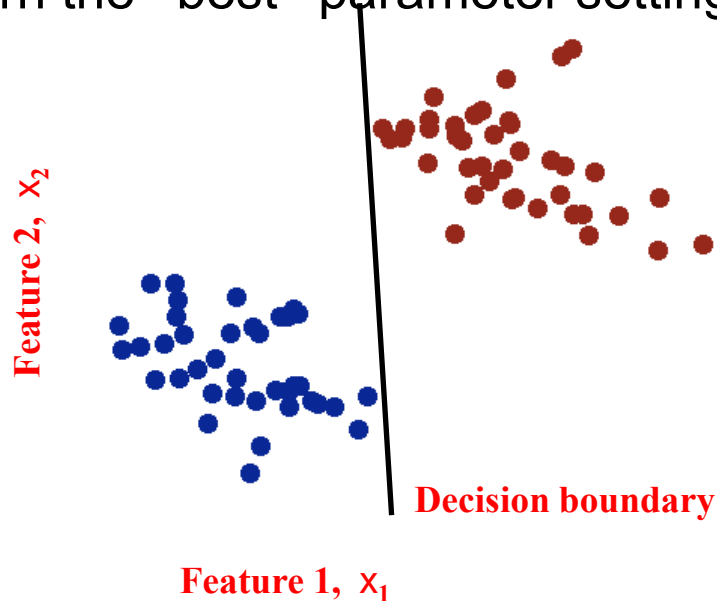
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y^{(i)}(\theta x^{(i)})\} + \frac{\lambda}{2m} \|\theta\|^2$$

$L_2$  regularization



# Linear classifiers

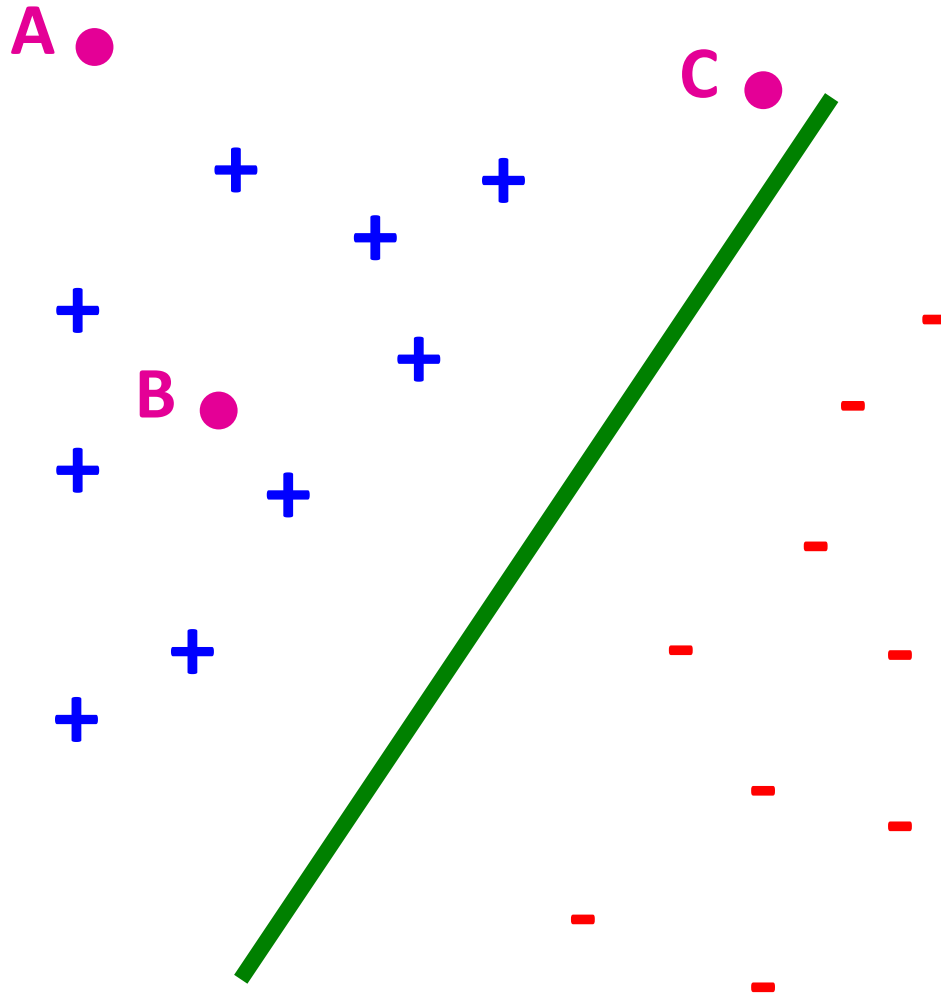
- Which decision boundary is “better”?
  - Both have zero training error (perfect training accuracy)
  - But, one of them seems intuitively better...
- How can we quantify “better”,  
and learn the “best” parameter settings?



# Largest Margin

Margin:

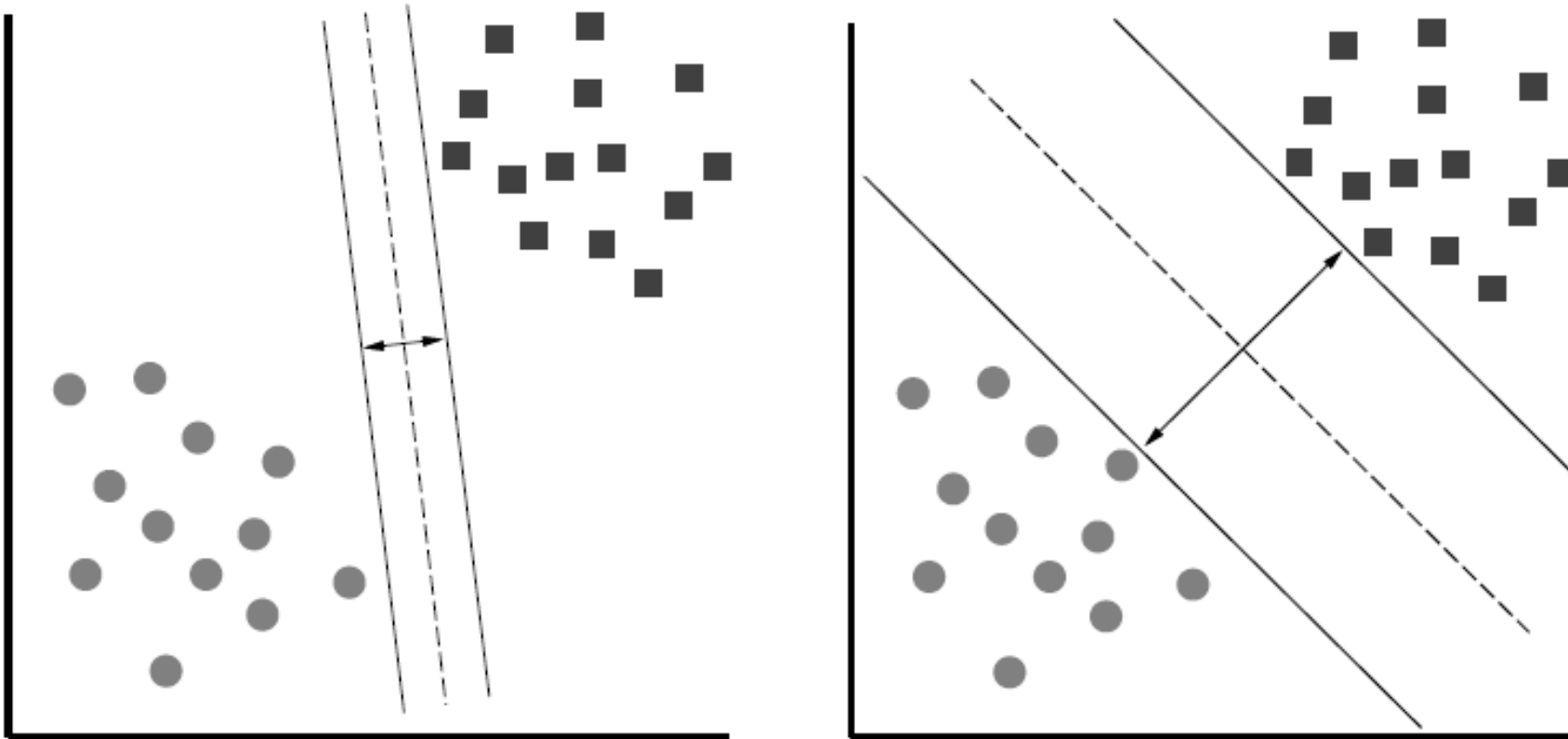
(2 ×) the distance from the decision boundary to the closest example.



- Distance from the separating hyperplane corresponds to the “confidence” of prediction
- Example:
  - We are more sure about the class of **A** and **B** than of **C**

# Largest Margin

- **Margin  $\gamma$ :** Distance of closest example from the decision line/hyperplane



The reason we define margin this way is due to theoretical convenience and existence of generalization error bounds that depend on the value of margin.

# Notations

- Parameters:

$$\theta_0, \theta_1, \dots, \theta_n \quad \Rightarrow \quad b, \overbrace{w_1, w_2, \dots, w_n}^w$$

- Data:

- Training examples:

$$(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$$

- Each example  $i$ :

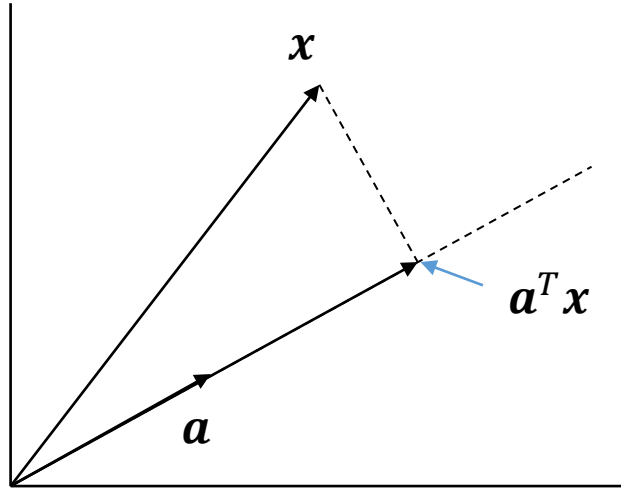
$$x^{(i)} = (x_1^{(i)}, \dots, x_n^{(i)})$$
$$y^{(i)} \in \{-1, +1\}$$

- Decision boundary:

$$\theta x = 0 \quad \Rightarrow \quad w \cdot x + b = 0$$



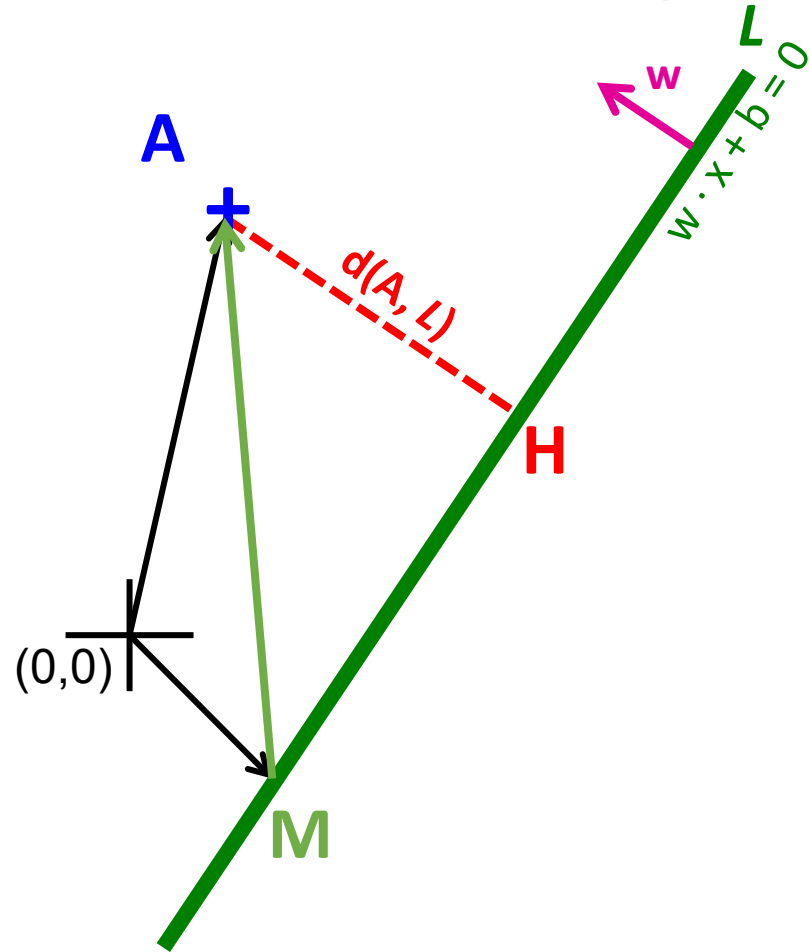
# Projection: Using Inner Products



Projection of  $\mathbf{x}$  along the direction  $\mathbf{a}$  ( $\|\mathbf{a}\| = 1$ )

$$\mathbf{p} = \mathbf{a} (\mathbf{a}^T \mathbf{x})$$
$$\|\mathbf{a}\| = \mathbf{a}^T \mathbf{a} = 1$$

# What is the margin?



## Distance from a point to a line

- **Let:**

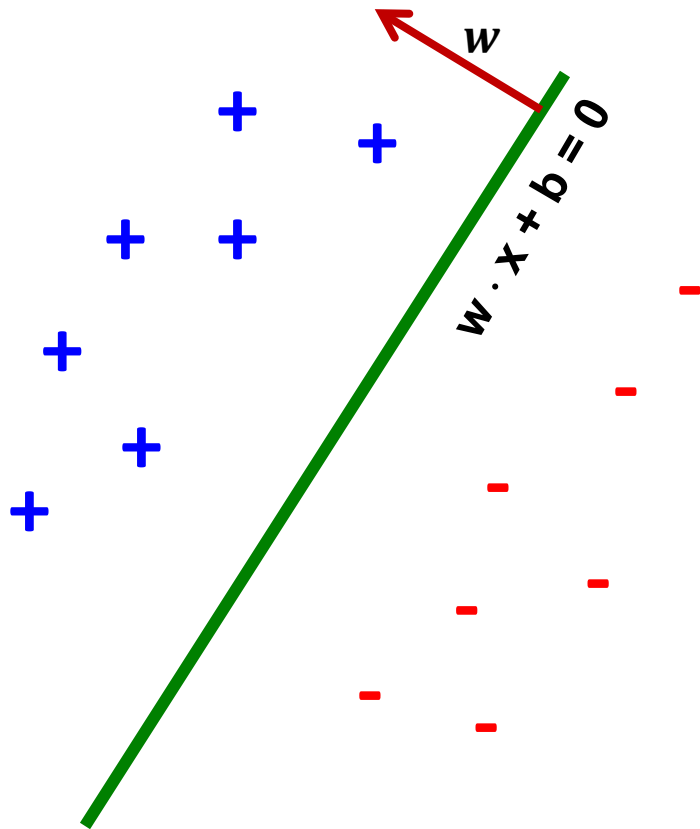
- **Line L:**  $w \cdot x + b = 0$
- **Point A**
- **Point M** on line L

Note we assume  
 $\|w\|_2 = 1$

$$\begin{aligned} d(A, L) &= AH \\ &= w \cdot (A - M) \\ &= w \cdot A - w \cdot M \\ &= w \cdot A + b \end{aligned}$$

Remember  $w \cdot M = -b$   
since **M** belongs to line **L**

# Largest Margin



- Prediction =  $\text{sign}(w \cdot x + b)$
- “**Confidence**” =  $(w \cdot x + b) y$
- For  $i$ -th datapoint:
$$\gamma^{(i)} = (w \cdot x^{(i)} + b) y^{(i)}$$
- Want to solve:
$$\max_{w,b} \min_i \gamma^{(i)}$$
- Can rewrite as
$$\begin{aligned} &\max_{w,b} \gamma \\ &s.t. \forall i, y^{(i)} (w \cdot x^{(i)} + b) \geq \gamma \end{aligned}$$

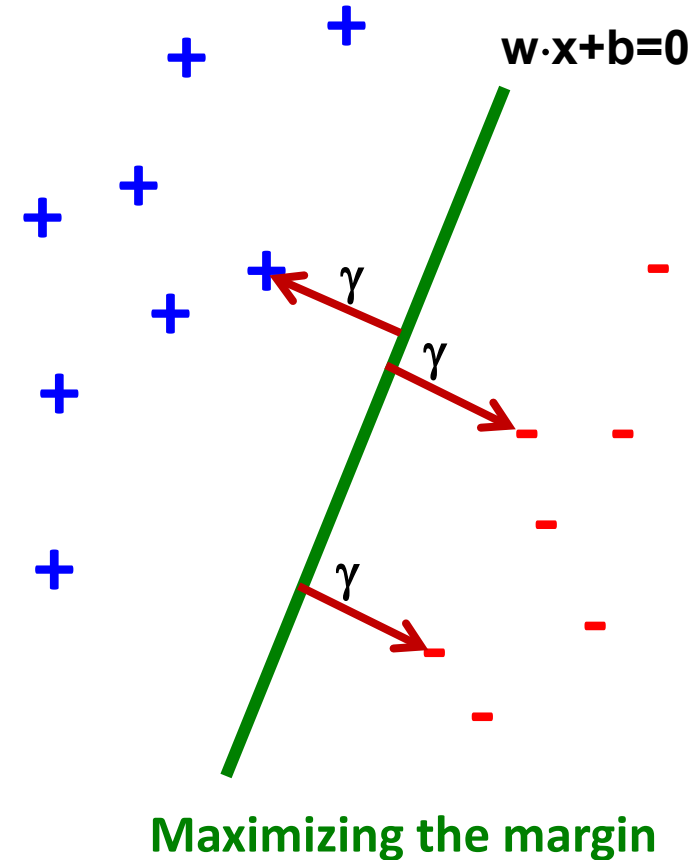
# Support Vector Machine

- **Maximize the margin:**

- Good according to intuition, theory (c.f. “VC dimension”) and practice

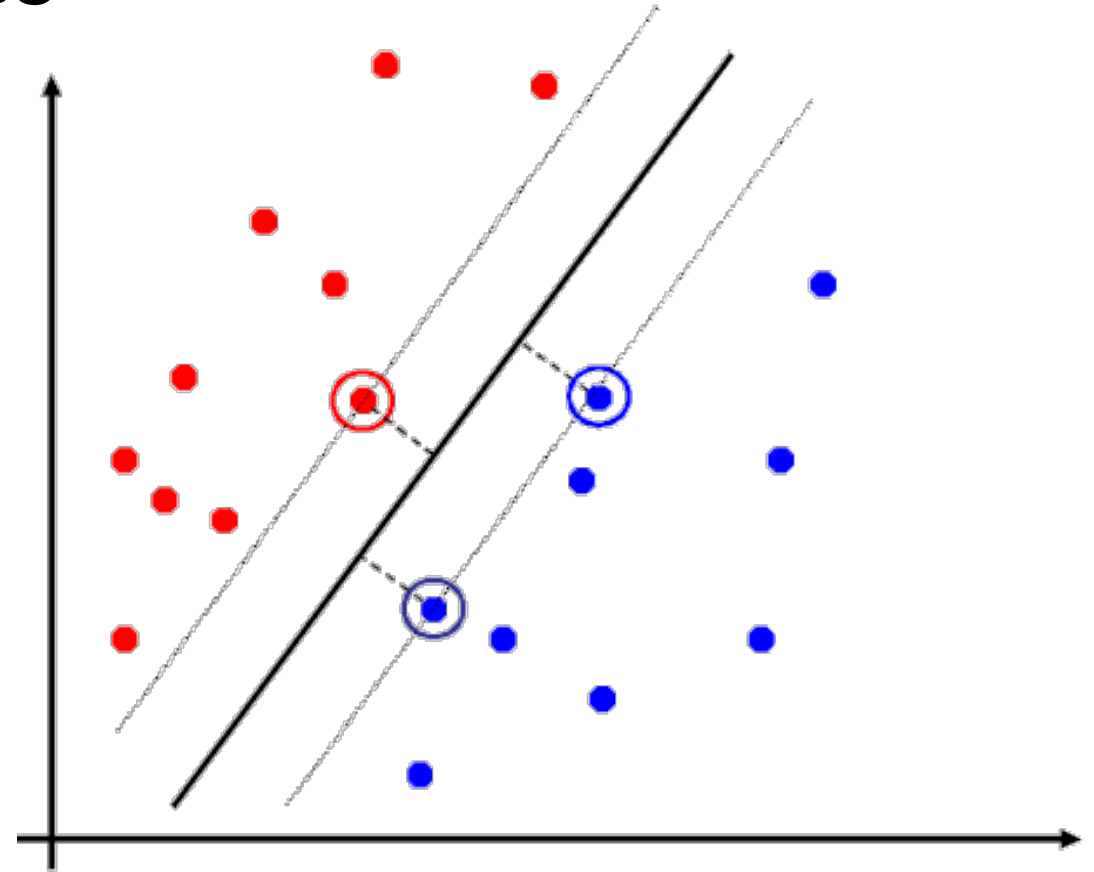
$$\begin{aligned} \max_{w,b} \gamma \\ \text{s.t. } \forall i, y^{(i)} (w \cdot x^{(i)} + b) \geq \gamma \end{aligned}$$

- $\gamma$  is margin ... distance from the separating hyperplane



# Support Vector Machines

- Separating hyperplane is defined by the support vectors
  - Points on  $\pm$  planes from the solution
  - If you knew these points, you could ignore the rest
  - Generally,  $n+1$  support vectors (for  $n$  dim. data)



# Support Vector Machines

Deriving the margin

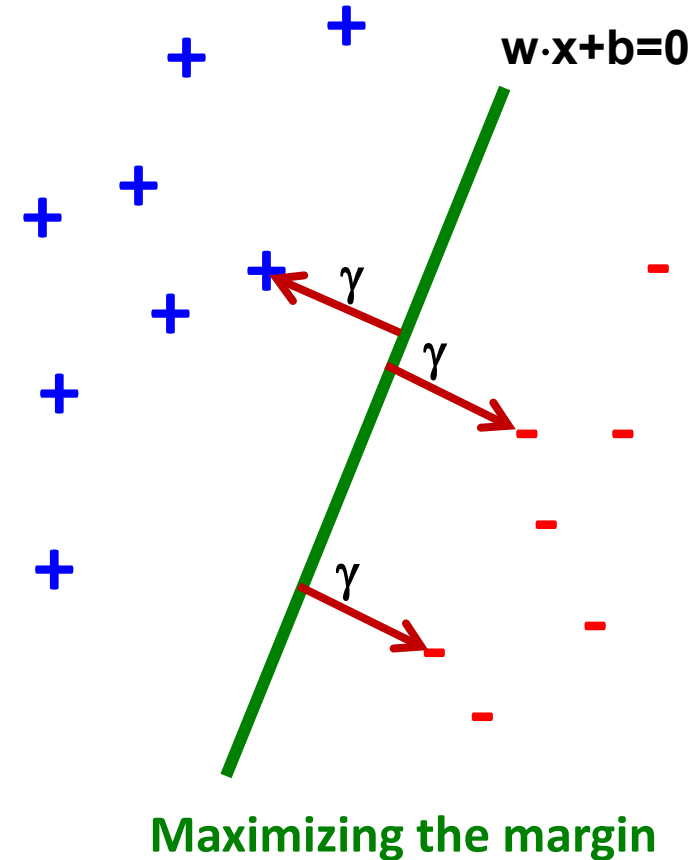
# Support Vector Machine

- **Maximize the margin:**

- Good according to intuition, theory (c.f. “VC dimension”) and practice

$$\begin{aligned} \max_{w,b} \gamma \\ \text{s.t. } \forall i, y^{(i)} (w \cdot x^{(i)} + b) \geq \gamma \end{aligned}$$

- $\gamma$  is margin ... distance from the separating hyperplane



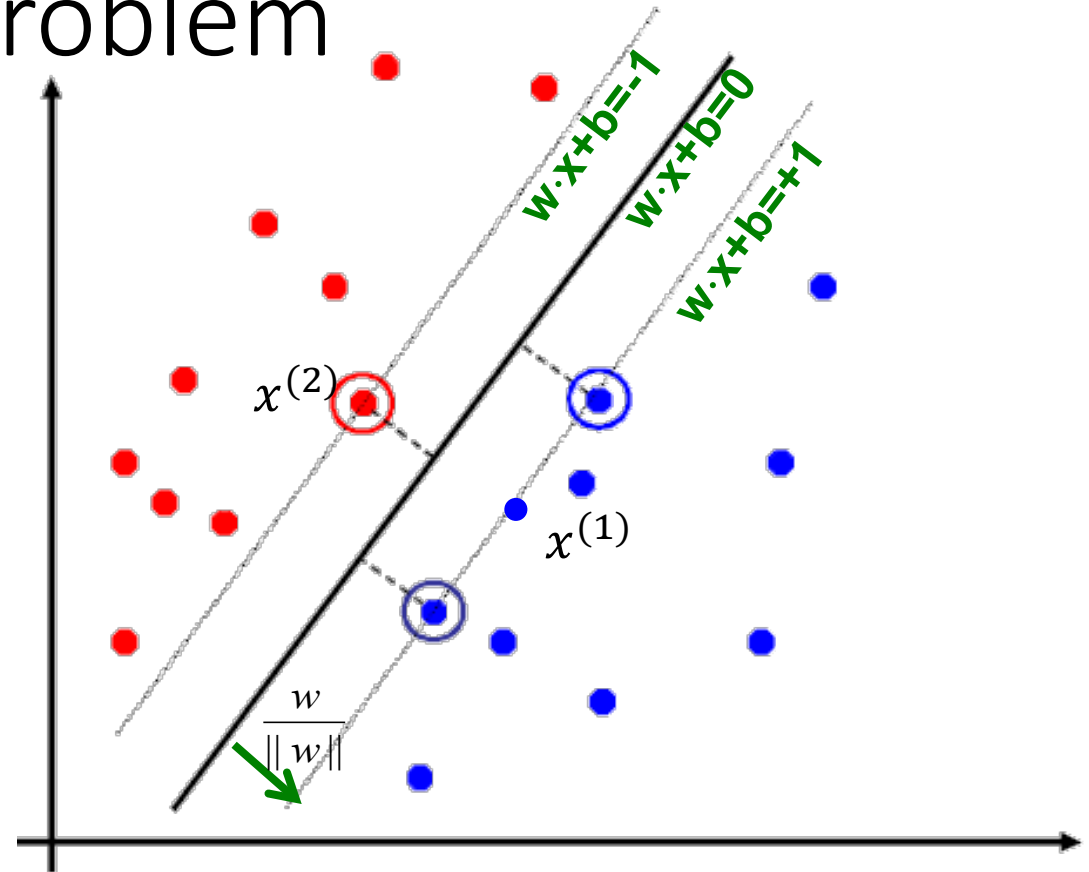
# Canonical Hyperplane: Problem

- **Problem:**

- Let  $(w \cdot x + b)y = \gamma$   
then  $(2w \cdot x + 2b)y = 2\gamma$ 
  - Scaling  $w$  increases margin!

- **Solution:**

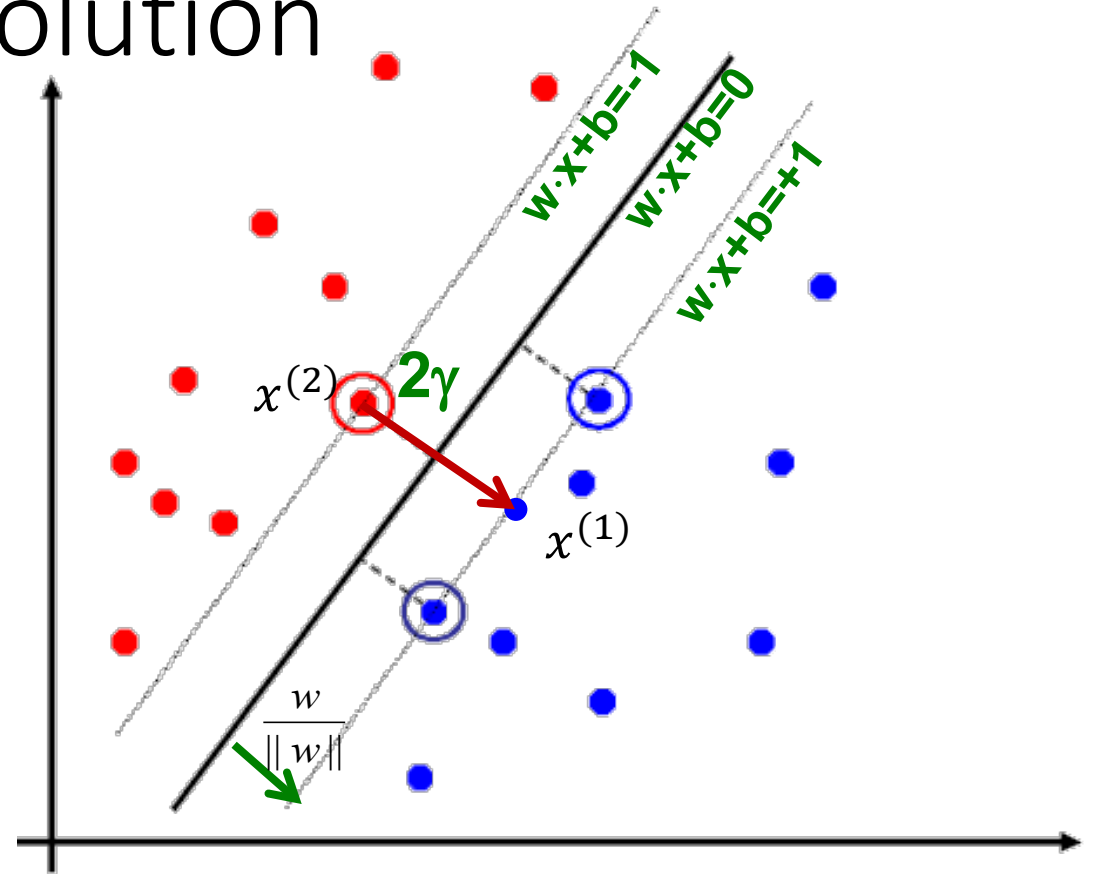
- Let's require **support vectors**  $x^{(j)}$  to be on the plane defined by:  
 $w \cdot x^{(j)} + b = \pm 1$





# Canonical Hyperplane: Solution

- **Want to maximize margin!**
- **What is the relation between  $x_1$  and  $x_2$ ?**
  - $x^{(1)} = x^{(2)} + 2\gamma \frac{w}{\|w\|}$
  - **We also know:**
    - $w \cdot x^{(1)} + b = +1$
    - $w \cdot x^{(2)} + b = -1$
- **So:**
  - $w \cdot x^{(1)} + b = +1$
  - $w \left( x^{(2)} + 2\gamma \frac{w}{\|w\|} \right) + b = +1$
  - $\underbrace{w \cdot x^{(2)} + b}_{-1} + 2\gamma \frac{w \cdot w}{\|w\|} = +1$



$$\Rightarrow \gamma = \frac{\|w\|}{w \cdot w} = \frac{1}{\|w\|}$$

**Note:**  
 $w \cdot w = \|w\|^2$

# Maximizing the Margin

- We started with

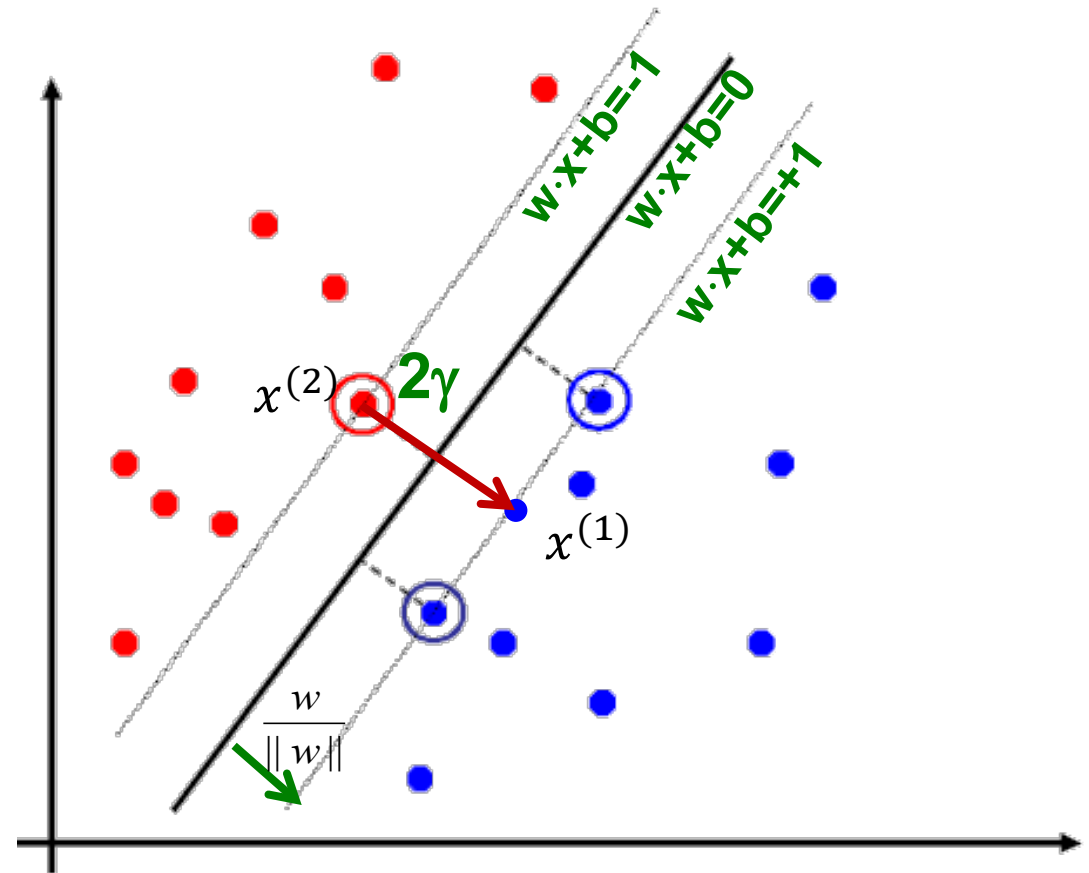
$$\begin{aligned} & \max_{w,b} \gamma \\ & s.t. \forall i, y^{(i)}(w \cdot x^{(i)} + b) \geq \gamma \end{aligned}$$

But  $w$  can be arbitrarily large!

- We normalized and...

$$\arg \max \gamma = \arg \max \frac{1}{\|w\|} = \arg \min \|w\| = \arg \min \frac{1}{2} \|w\|^2$$

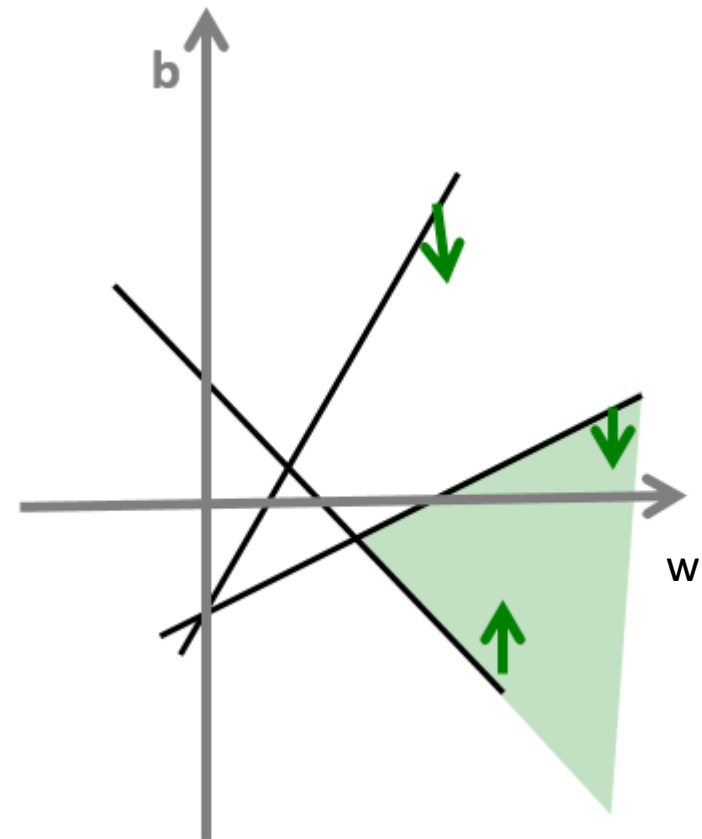
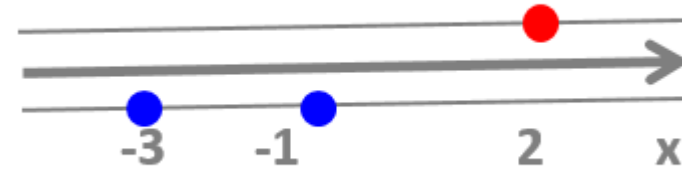
- Then: 
$$\begin{aligned} & \min_{w,b} \frac{1}{2} \|w\|^2 \\ & s.t. \forall i, y^{(i)}(w \cdot x^{(i)} + b) \geq 1 \end{aligned}$$



← Quadratic Programming

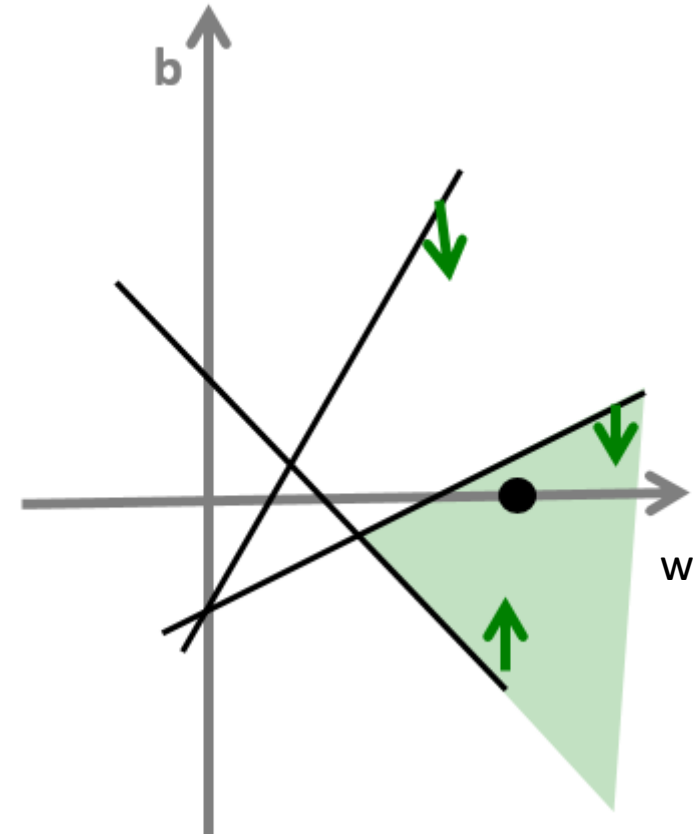
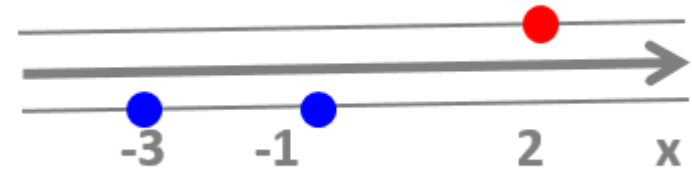
# A 1D Example

- Suppose we have three data points
  - $x = -3, y = -1$
  - $x = -1, y = -1$
  - $x = 2, y = 1$
- Many separating perceptrons,  $\text{sign}[wx + b]$ 
  - Anything with  $w x + b = 0$  between -1 and 2
- We can write the margin constraints
  - $(-1) * (w(-3) + b) > 1 \quad \Rightarrow b < 3w - 1$
  - $(-1) * (w(-1) + b) > 1 \quad \Rightarrow b < w - 1$
  - $w(2) + b > 1 \quad \Rightarrow b > -2w + 1$



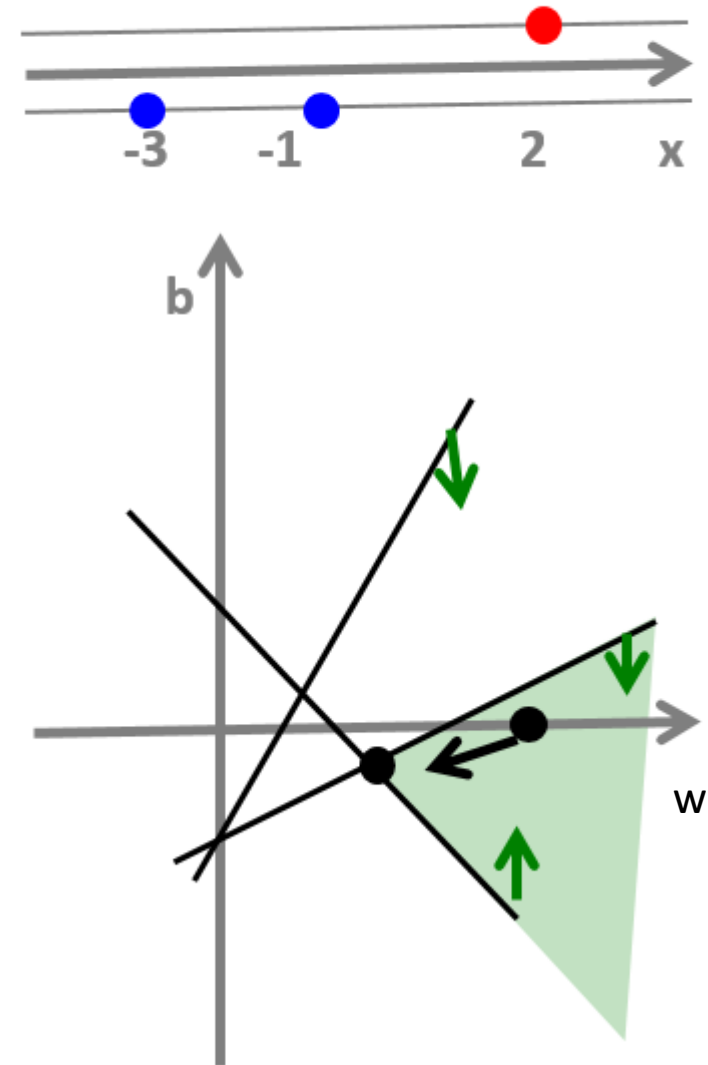
# A 1D Example

- Suppose we have three data points
  - $x = -3, y = -1$
  - $x = -1, y = -1$
  - $x = 2, y = 1$
- Many separating perceptrons,  $\text{sign}[wx + b]$ 
  - Anything with  $w x + b = 0$  between -1 and 2
- We can write the margin constraints
  - $(-1) * (w(-3) + b) > 1 \quad \Rightarrow b < 3w - 1$
  - $(-1) * (w(-1) + b) > 1 \quad \Rightarrow b < w - 1$
  - $w(2) + b > 1 \quad \Rightarrow b > -2w + 1$
- Ex:  $w = 1, b = 0$



# A 1D Example

- Suppose we have three data points
  - $x = -3, y = -1$
  - $x = -1, y = -1$
  - $x = 2, y = 1$
- Many separating perceptrons,  $\text{sign}[wx + b]$ 
  - Anything with  $w x + b = 0$  between -1 and 2
- We can write the margin constraints
  - $(-1)*(w(-3) + b) > 1 \quad \Rightarrow b < 3w - 1$
  - $(-1)*(w(-1) + b) > 1 \quad \Rightarrow b < w - 1$
  - $w(2) + b > 1 \quad \Rightarrow b > -2w + 1$
- Ex:  $w = 1, b = 0$
- Minimize  $||w|| \Rightarrow w = .66, b = -.33$ 
  - Two data on the margin; constraints “tight”

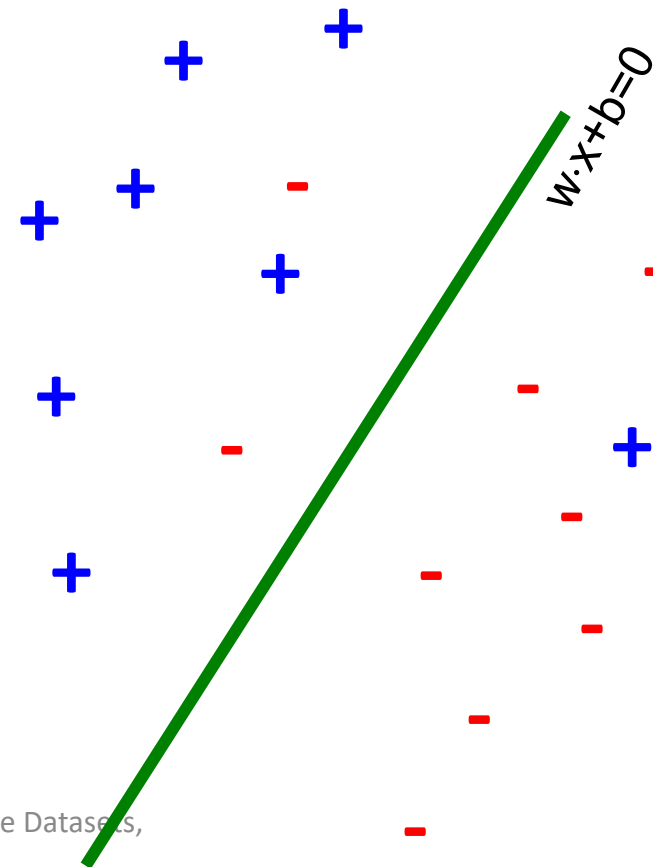


# Non-linearly Separable Data

- If data is **not separable** introduce **penalty**:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \cdot (\text{\#number of mistakes})$$
$$s.t. \forall i, y^{(i)}(w \cdot x^{(i)} + b) \geq 1$$

- Minimize  $\|w\|^2$  plus the number of training mistakes
  - Set  $C$  using cross validation
- **How to penalize mistakes?**
    - All mistakes are not equally bad!



# Support Vector Machines

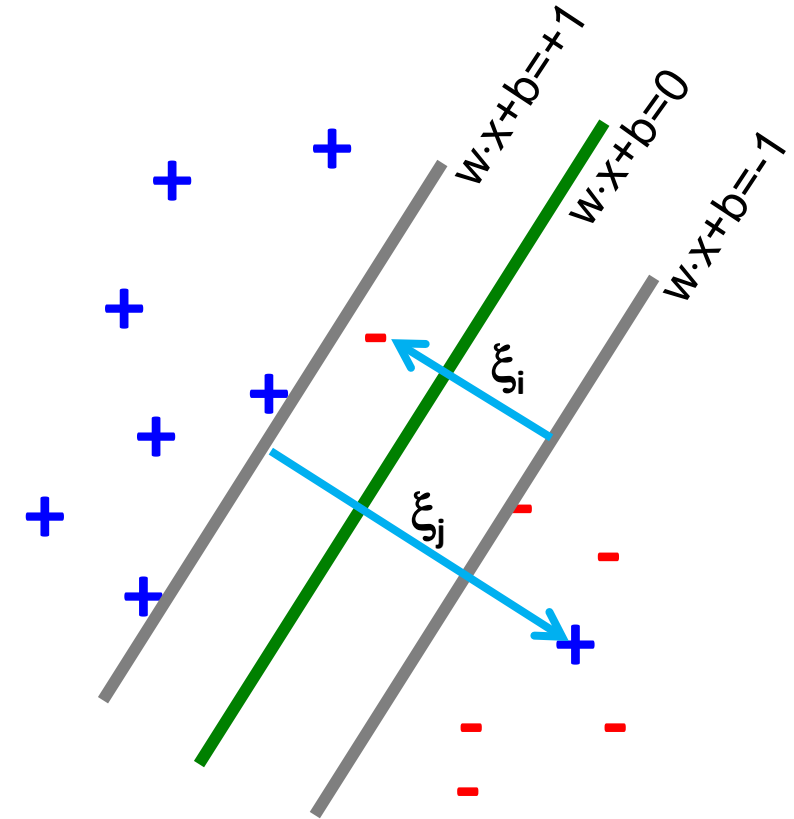
- Introduce **slack variables**  $\xi_i$

$$\min_{w, b, \xi_i \geq 0} \frac{1}{2} \|w\|^2 + C \cdot \sum_{i=1}^n \xi_i$$

$$s.t. \forall i, y^{(i)}(w \cdot x^{(i)} + b) \geq 1 - \xi_i$$

- If point  $x^{(i)}$  is on the wrong side of the margin then get penalty  $\xi_i$

**SVM with “soft” constraints**



**For each data point:**

If margin  $\geq 1$ , don't care

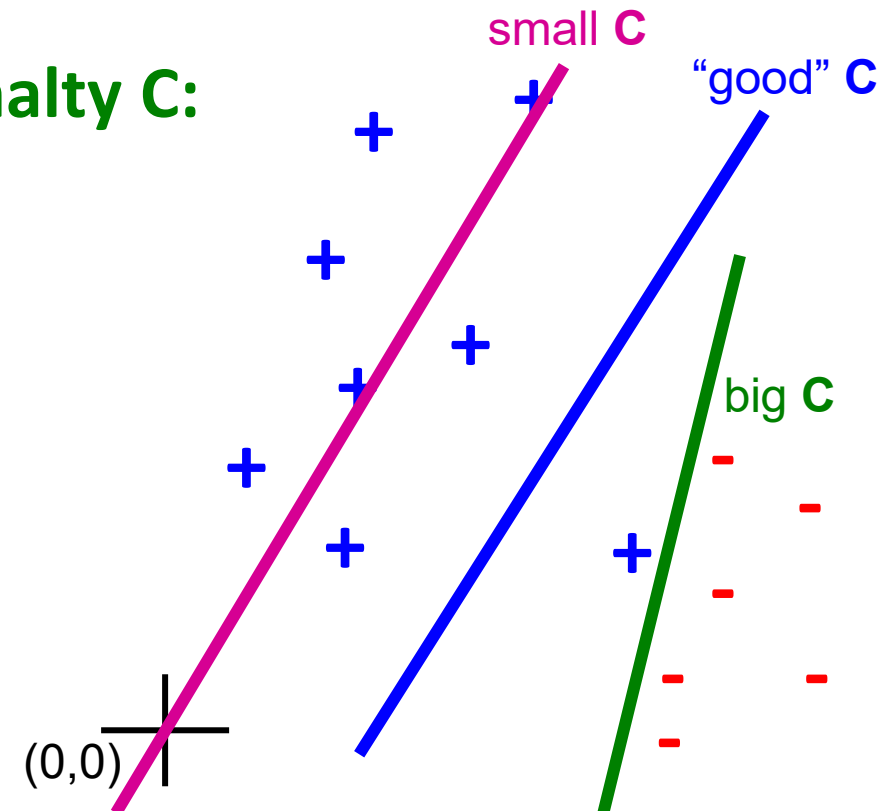
If margin  $< 1$ , pay linear penalty

# Slack Penalty $C$

$$\min_{w, b, \xi_i \geq 0} \frac{1}{2} \|w\|^2 + C \cdot \sum_{i=1}^n \xi_i$$
$$s.t. \forall i, y^{(i)}(w \cdot x^{(i)} + b) \geq 1 - \xi_i$$

- **What is the role of slack penalty  $C$ :**

- $C=\infty$ : Only want to  $w, b$  that separate the data
- $C=0$ : Can set  $\xi_i$  to anything, then  $w=0$  (basically ignores the data)

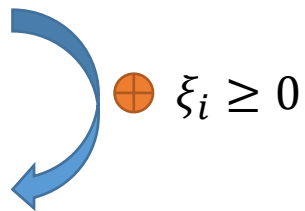





# Support Vector Machines

- Combine the constraints and the objective function

$$\begin{array}{ll} \min_{w, b, \xi_i \geq 0} & \frac{1}{2} \|w\|^2 + C \cdot \sum_{i=1}^n \xi_i \\ \text{s.t. } \forall i, & y^{(i)}(w \cdot x^{(i)} + b) \geq 1 - \xi_i \end{array}$$

$$\Rightarrow \xi_i \geq 1 - y^{(i)}(w \cdot x^{(i)} + b)$$


$$\xi_i = \begin{cases} 0, & \text{if } 1 - y^{(i)}(w \cdot x^{(i)} + b) \leq 0 \\ 1 - y^{(i)}(w \cdot x^{(i)} + b), & \text{if } 1 - y^{(i)}(w \cdot x^{(i)} + b) > 0 \end{cases}$$


$$\xi_i = \max\{0, 1 - y^{(i)}(w \cdot x^{(i)} + b)\}$$

# Support Vector Machines

- SVM in the “natural” form

$$\operatorname{argmin}_{w,b} \underbrace{\frac{1}{2} w \cdot w}_{\text{Margin}} + C \cdot \underbrace{\sum_{i=1}^n \max\{0, 1 - y^{(i)}(w \cdot x^{(i)} + b)\}}_{\text{Empirical loss } \mathbf{L} \text{ (how well we fit training data)}}$$

# Support Vector Machines

How to compute the margin?

# SVM: How to estimate $w$ ?

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} w \cdot w + C \cdot \sum_{i=1}^n \xi_i \\ \text{s.t. } & \forall i, y^{(i)} (w \cdot x^{(i)} + b) \geq 1 - \xi_i \end{aligned}$$

- **Want to estimate  $w$  and  $b$ !**
  - **Standard way:** Use a solver!
    - **Solver:** software for finding solutions to “common” optimization problems
- **Use a quadratic solver:**
  - Minimize quadratic function
  - Subject to linear constraints
- **Problem:** Solvers are inefficient for big data!

# SVM: How to estimate $w$ ?

- **Want to minimize  $J(w, b)$ :**

$$J(w, b) = \frac{1}{2} \sum_{j=1}^d (w_j)^2 + C \sum_{i=1}^n \max \left\{ 0, 1 - y^{(i)} \left( \sum_{j=1}^d w_j x_j^{(i)} + b \right) \right\}$$

**Empirical loss  $L(x^{(i)}, y^{(i)})$**

- **Compute the gradient  $\nabla_j$  w.r.t.  $w_j$**

$$\nabla J_j = \frac{\partial J(w, b)}{\partial w_j} = w_j + C \sum_{i=1}^n \frac{\partial L(x^{(i)}, y^{(i)})}{\partial w_j}$$

$$\begin{aligned} \frac{\partial L(x^{(i)}, y^{(i)})}{\partial w_j} &= 0 && \text{if } y^{(i)}(w \cdot x^{(i)} + b) \geq 1 \\ &= -y^{(i)} x_j^{(i)} && \text{else} \end{aligned}$$

# SVM: How to estimate $w$ ?

- **Batch Gradient Descent:**

**Iterate until convergence:**

- **For  $j = 1 \dots d$**

- **Evaluate:**  $\nabla J_j = \frac{\partial J(w, b)}{\partial w_j} = w_j + C \sum_{i=1}^n \frac{\partial L(x^{(i)}, y^{(i)})}{\partial w_j}$
  - **Update:**

$$w'_j \leftarrow w_j - \eta \nabla J_j$$

- **$w \leftarrow w'$**

$\eta$ ...learning rate parameter  
 $C$ ... regularization parameter

# SVM: How to estimate $w$ ?

- **Stochastic Gradient Descent**

- Instead of evaluating gradient over all examples evaluate it for each **individual** training example

$$\nabla J_j(x^{(i)}) = w_j + C \cdot \frac{\partial L(x^{(i)}, y^{(i)})}{\partial w_j}$$

← Notice: no summation over  $i$  anymore

- **Stochastic gradient descent:**

**Iterate until convergence:**

- **For**  $i = 1 \dots n$ 
  - **For**  $j = 1 \dots d$ 
    - **Compute:**  $\nabla J_j(x^{(i)})$
    - **Update:**  $w'_j \leftarrow w_j - \eta \nabla J_j(x^{(i)})$

# Support Vector Machines

Example



# Example: Text categorization

- **Example by Leon Bottou:**
  - **Reuters RCV1** document corpus
    - Predict a category of a document
      - One **vs.** the rest classification
  - **$n = 781,000$**  training examples (documents)
  - 23,000 test examples
  - **$d = 50,000$**  features
    - One feature per word
    - Remove stop-words
    - Remove low frequency words

# Example: Text categorization

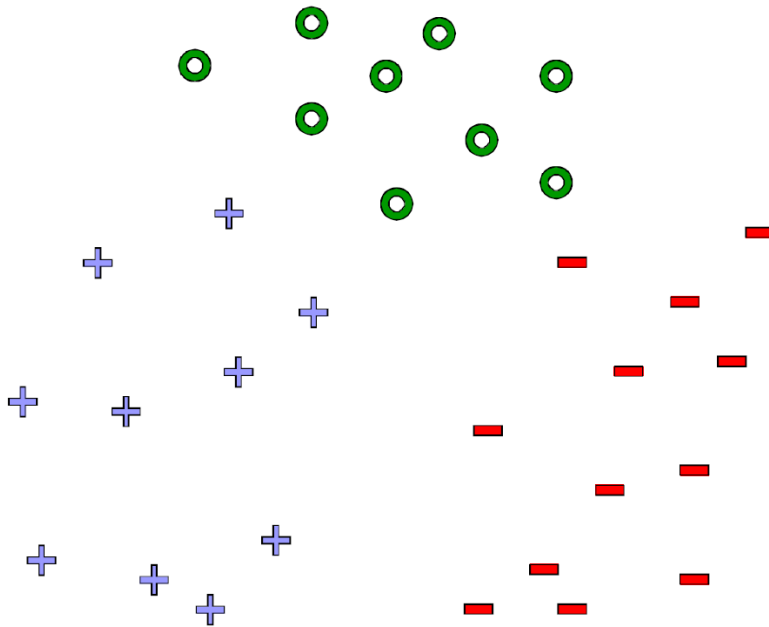
- **Questions:**

- (1) Is **SGD** successful at minimizing  $J(\mathbf{w}, \mathbf{b})$ ?
- (2) How quickly does **SGD** find the min of  $J(\mathbf{w}, \mathbf{b})$ ?
- (3) What is the error on a test set?

	<i>Training time</i>	<i>Value of <math>J(\mathbf{w}, \mathbf{b})</math></i>	<i>Test error</i>
Standard SVM	23,642 secs	0.2275	6.02%
“Fast SVM”	66 secs	0.2278	6.03%
<b>SGD-SVM</b>	1.4 secs	0.2275	6.02%

- (1) SGD-SVM is successful at minimizing the value of  $J(\mathbf{w}, \mathbf{b})$
- (2) SGD-SVM is super fast
- (3) SGD-SVM test set error is comparable

# What about multiple classes?



- **Idea 1:**

- One against all**

Learn 3 classifiers

- + vs. {o, -}
- - vs. {o, +}
- o vs. {+, -}

Obtain:

$$w_+ b_+, w_- b_-, w_o b_o$$

- **How to classify?**

- Return class  $c$

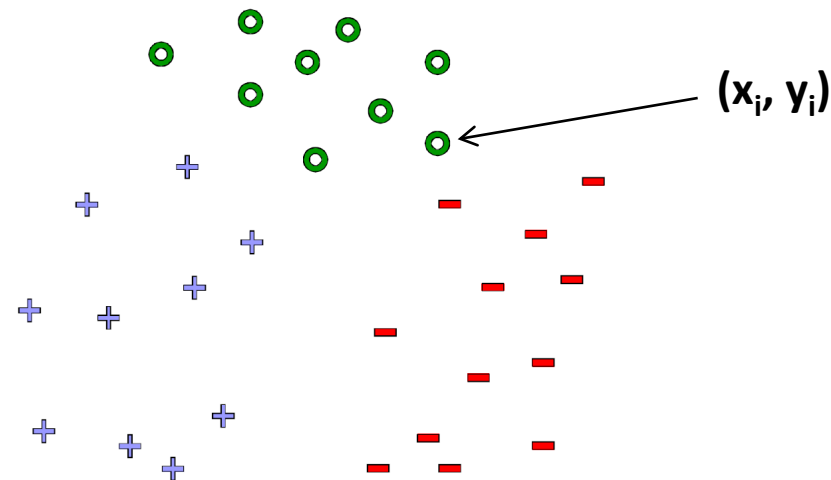
$$\arg \max_c w_c x + b_c$$

# Multiclass SVM

- **Idea 2: Learn 3 sets of weights simultaneously!**

- For each class  $c$  estimate  $w_c, b_c$
- **Want the correct class  $y_i$  to have highest margin:**

$$w_{y_i} x_i + b_{y_i} \geq 1 + w_c x_i + b_c \quad \forall c \neq y_i, \forall i$$



# Multiclass SVM

- **Optimization problem:**

$$\min_{w,b} \frac{1}{2} \sum_c \|w_c\|^2 + C \sum_{i=1}^n \xi_i$$
$$w_{y_i} \cdot x_i + b_{y_i} \geq w_c \cdot x_i + b_c + 1 - \xi_i \quad \forall c \neq y_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

- To obtain parameters  $w_c, b_c$  (for each class  $c$ ) we can use similar techniques as for 2 class **SVM**

- **SVM is widely perceived a very powerful learning algorithm**