

---

# MACHINE LEARNING

## Linear Classifier: The Perceptron

**Alessandro Moschitti**

Department of information and communication technology

University of Trento

Email: [moschitti@dit.unitn.it](mailto:moschitti@dit.unitn.it)



# Summary

---

- Computational Learning theory
  - Perceptron Learning
  - Margins



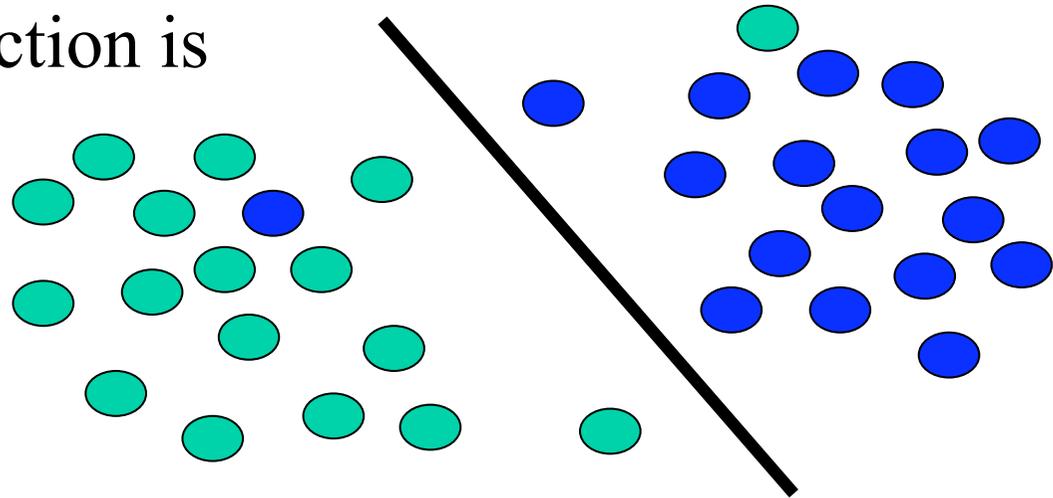
# Linear Classifier (1)

---

- The equation of a hyperplane is

$$f(\vec{x}) = \vec{x} \cdot \vec{w} + b = 0, \quad \vec{x}, \vec{w} \in \mathfrak{R}^n, b \in \mathfrak{R}$$

- $\vec{x}$  is the vector representing the classifying example
- $\vec{w}$  is the gradient to the hyperplane
- The classification function is  
 $h(x) = \text{sign}(f(x))$



# Linear classifiers (2)

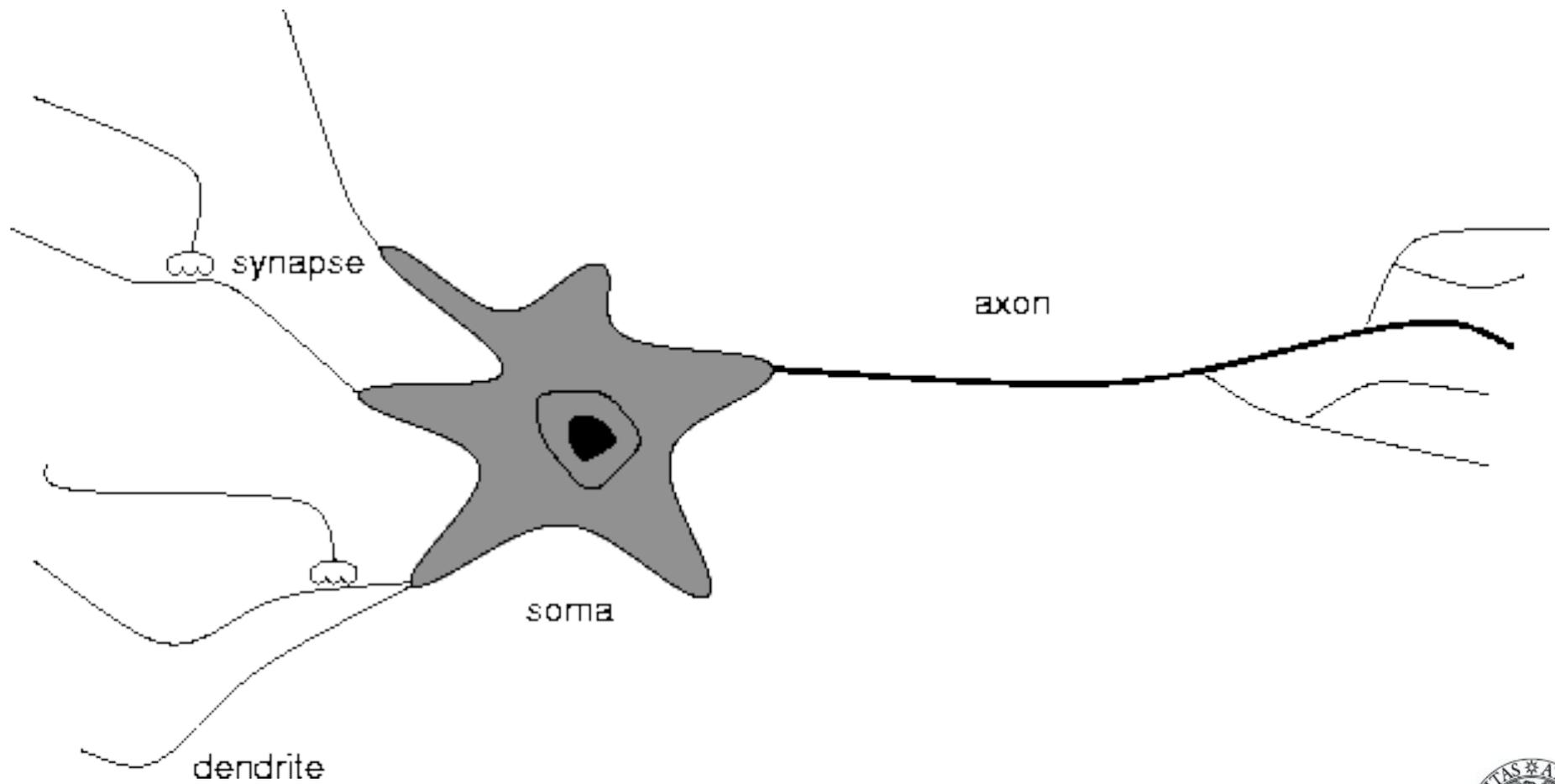
---

- Linear Functions are the simplest ones from an analytical point of view.
- The basic idea is to select a hypothesis with null error on the training-set.
- To learn a linear function a simple neural network of only one neuron is enough (Perceptron)



# An animal neuron

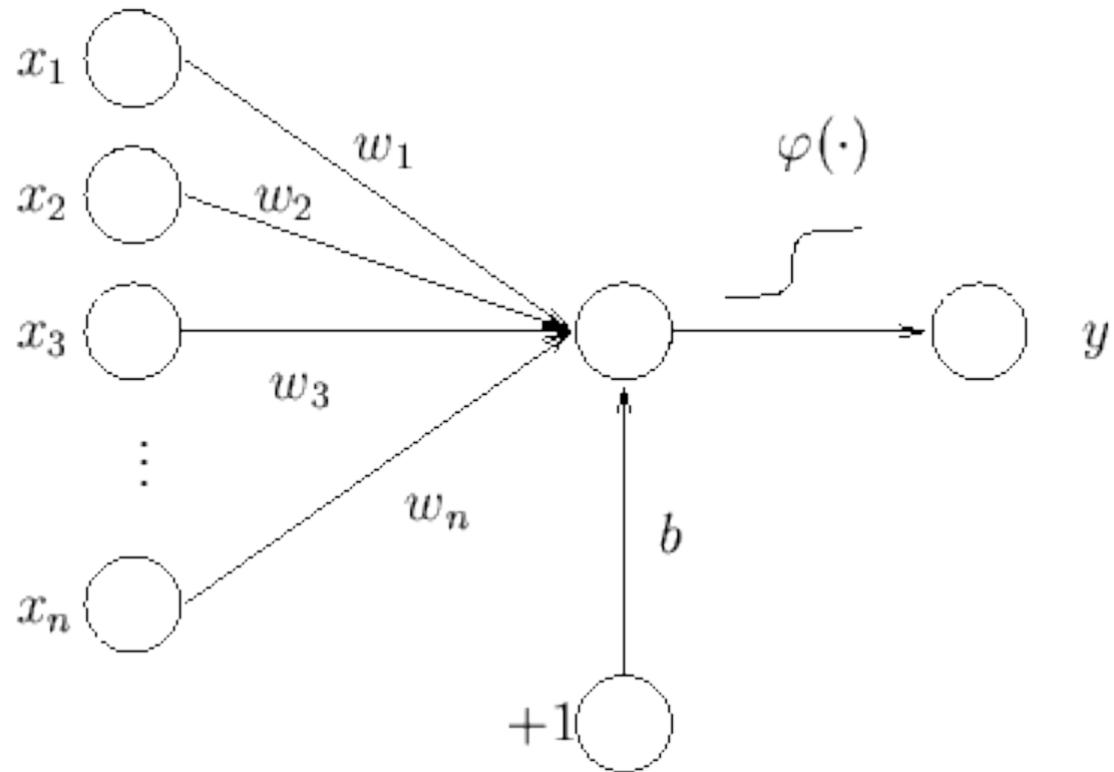
---



# The Perceptron

---

$$\varphi(\vec{x}) = \text{sgn}\left(\sum_{i=1..n} w_i \times x_i + b\right)$$



# Useful Concepts

---

- ***Functional Margin*** of an example with respect to a hyperplane:  $\gamma_i = y_i(\vec{w} \cdot \vec{x}_i + b)$
- ***The distribution of functional margins*** of a hyperplane with respect to a training set  $S$  is the distribution of the margins of the examples in  $S$  wrt the hyperplane  $(\vec{w}, b)$ .
- ***The functional margin of a hyperplane*** is the minimum margin of the distribution



# Notations (con'td)

---

- If we normalize the hyperplane equation, i.e.

$$\left( \frac{\vec{w}}{\|\vec{w}\|}, \frac{b}{\|\vec{w}\|} \right), \text{ we obtain the } \textit{geometric margin}$$

- The *geometric margin* measure the Euclidean distance between the target point and the hyperplane.
- *The training set Margin* is the maximum geometric (functional) margin among all hyperplanes which separates the examples in S.
- The hyperplane associated with the above quantity is called *maximal margin hyperplane*



# Basic Concepts

---

- From  $\cos(\vec{x}, \vec{w}) = \frac{\vec{x} \cdot \vec{w}}{\|\vec{x}\| \cdot \|\vec{w}\|}$

- It follows that

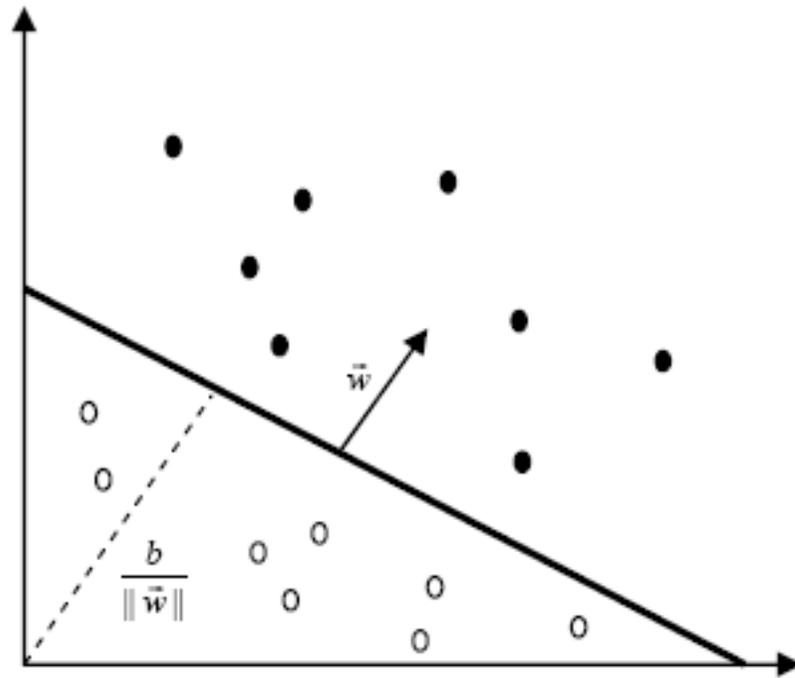
$$\|\vec{x}\| \cos(\vec{x}, \vec{w}) = \frac{\vec{x} \cdot \vec{w}}{\|\vec{w}\|} = \vec{x} \cdot \frac{\vec{w}}{\|\vec{w}\|}$$

- Norm of  $\vec{x}$  times the cosine between  $\vec{x}$  and  $\vec{w}$ , i.e. the projection of  $\vec{x}$  on  $\vec{w}$

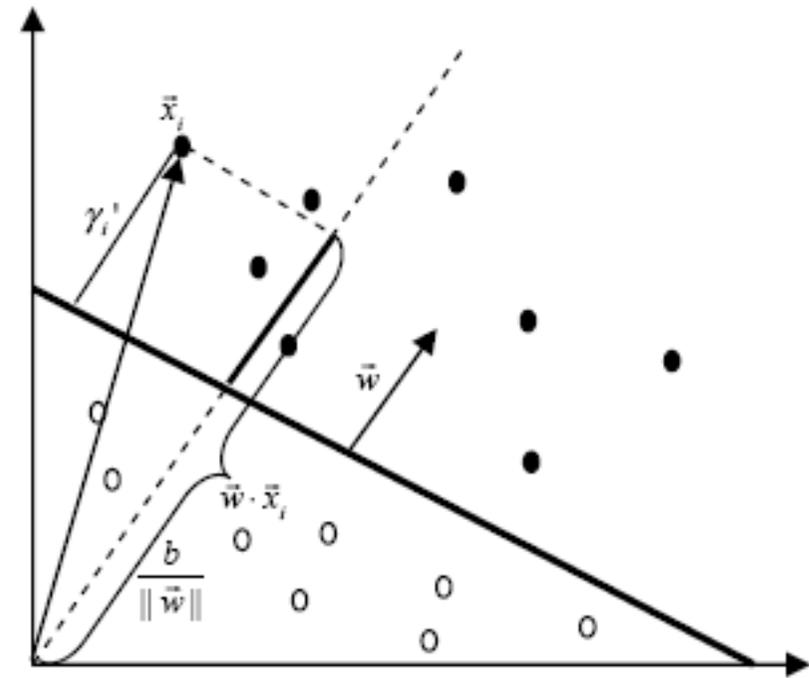


# Geometric Margin

---



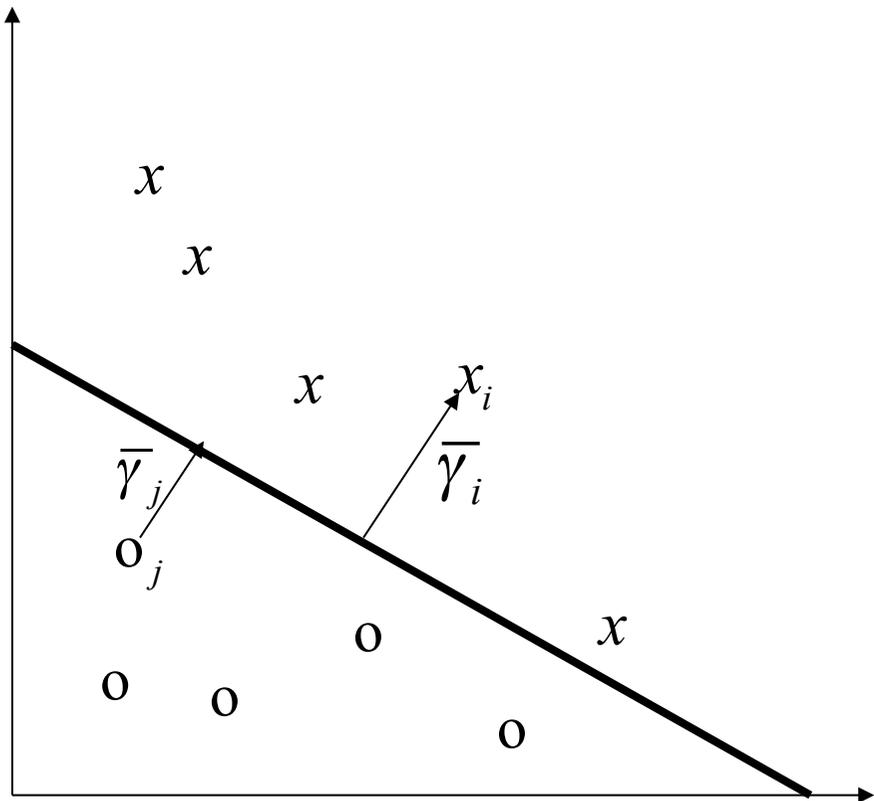
A



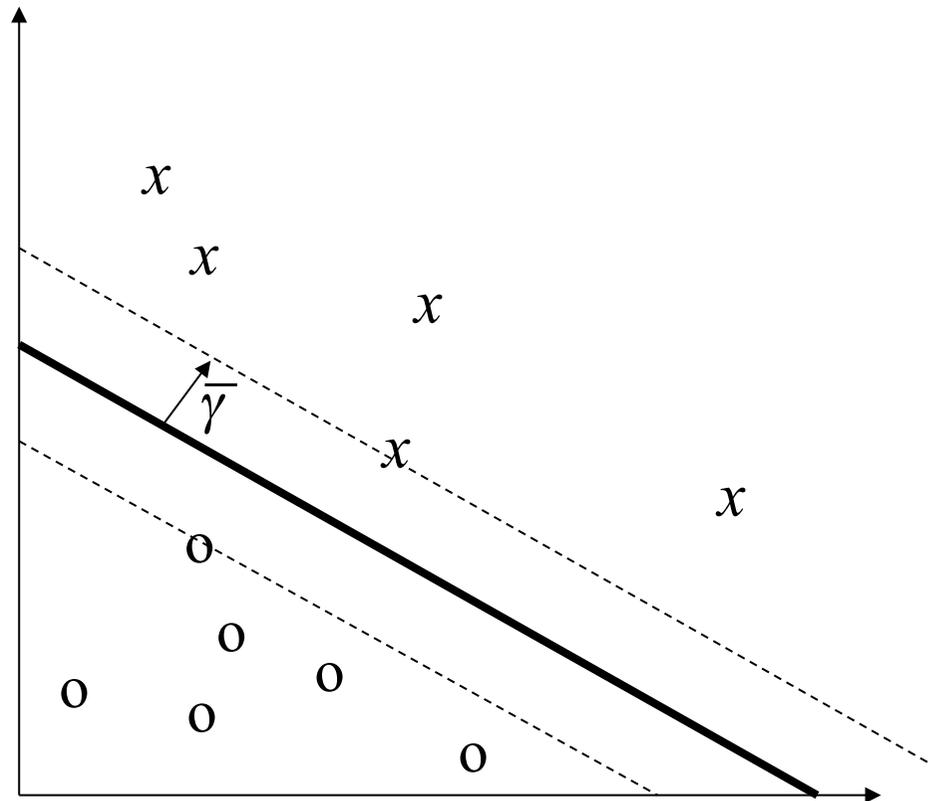
B

# Geometric margins of 2 points and hyperplane margin

---



Geometric Margin

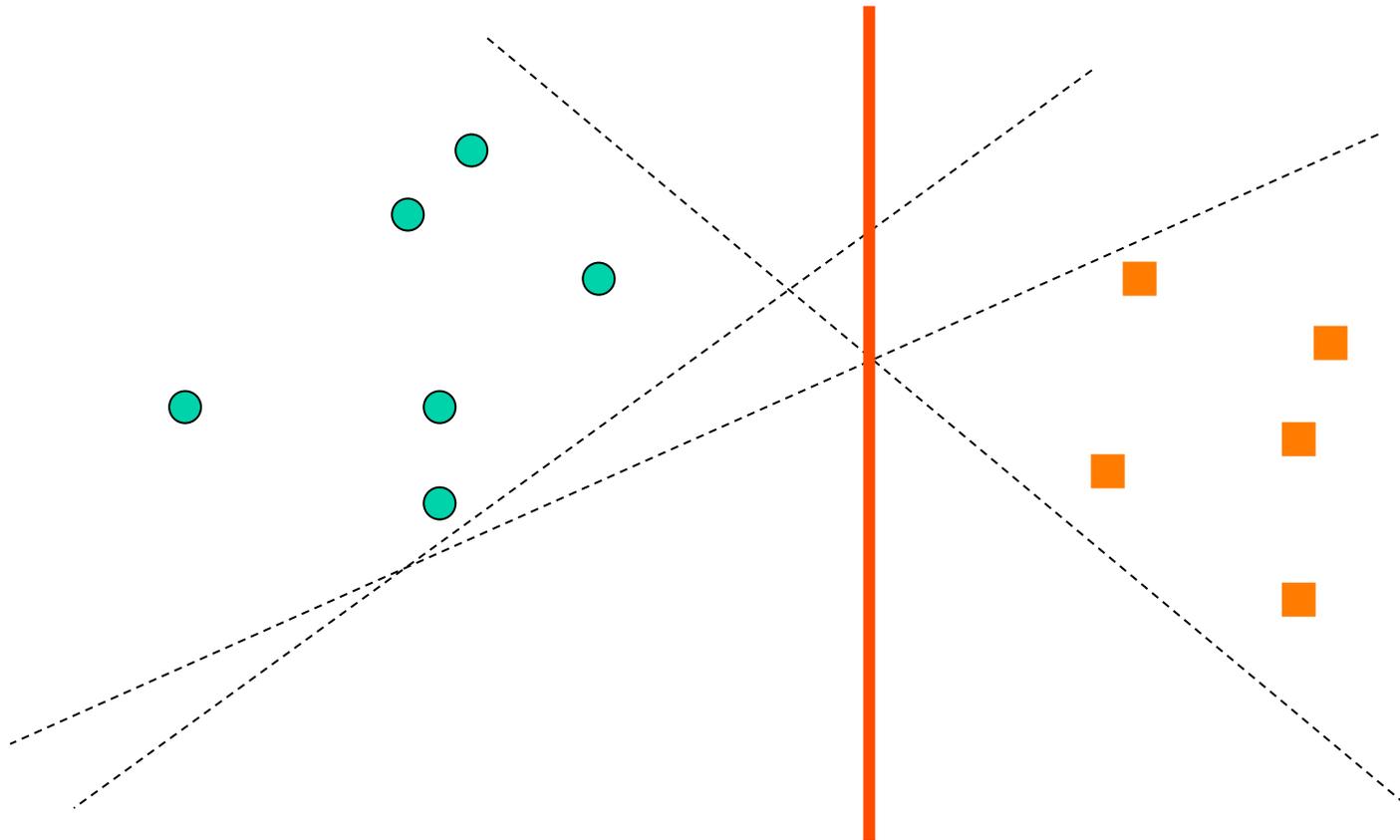


Hyperplane Margin



# Maximal margin vs other margins

---



# Perceptron training on a data set (on-line algorithm)

---

$$\vec{w}_0 \leftarrow \vec{0}; b_0 \leftarrow 0; k \leftarrow 0; R \leftarrow \max_{1 \leq i \leq l} \|\vec{x}_i\|$$

Repeat

for  $i = 1$  to  $m$

if  $y_i(\vec{w}_k \cdot \vec{x}_i + b_k) \leq 0$  then

$$\vec{w}_{k+1} = \vec{w}_k + \eta y_i \vec{x}_i$$

$$b_{k+1} = b_k + \eta y_i R^2$$

$$k = k + 1$$

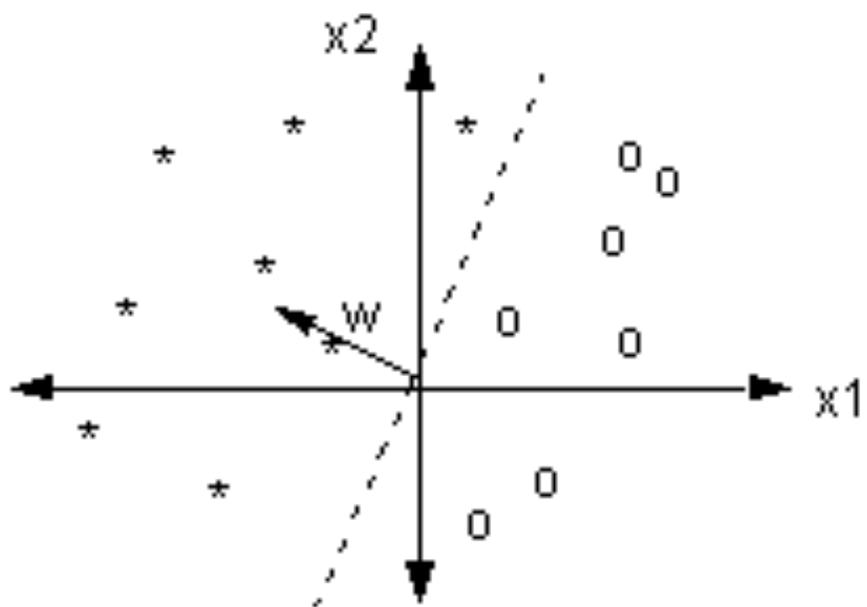
endif

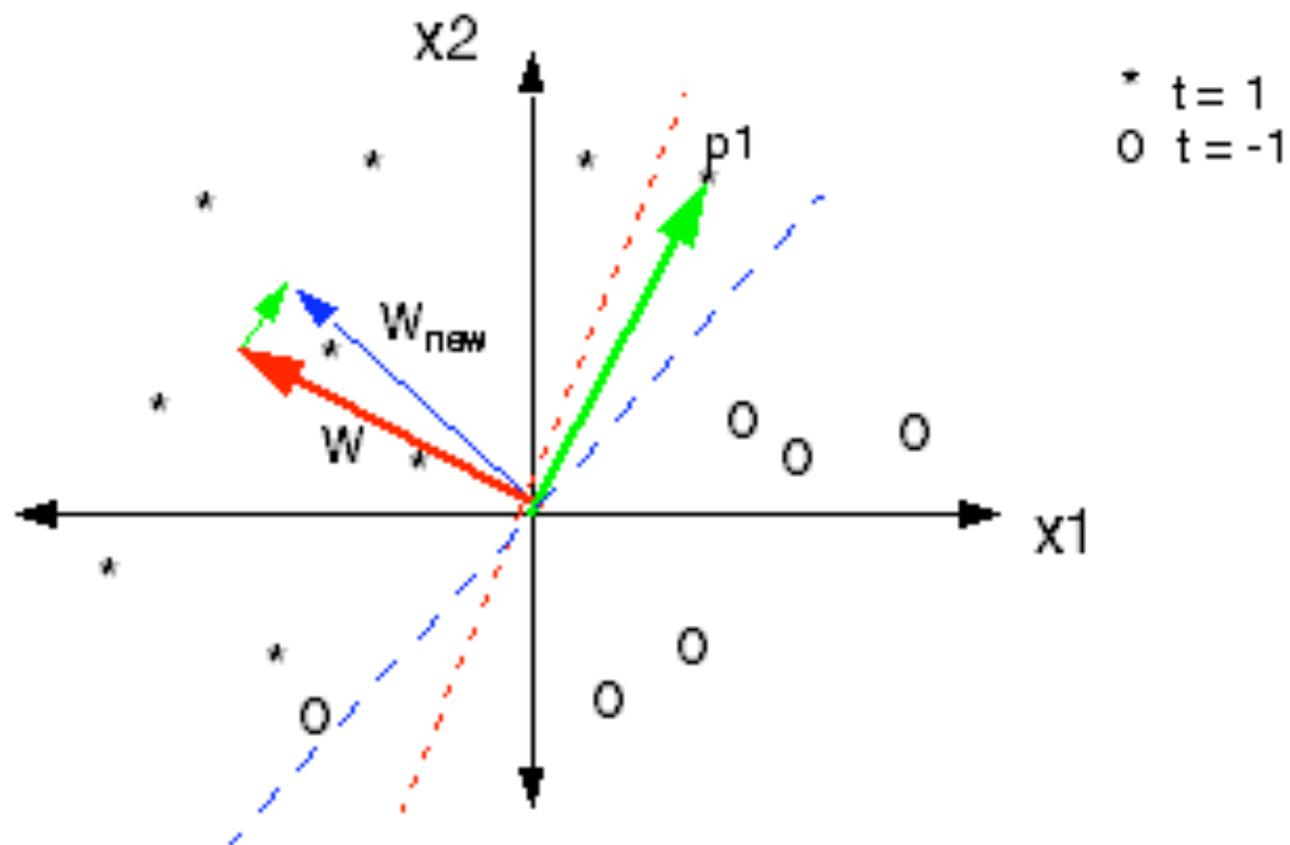
endfor

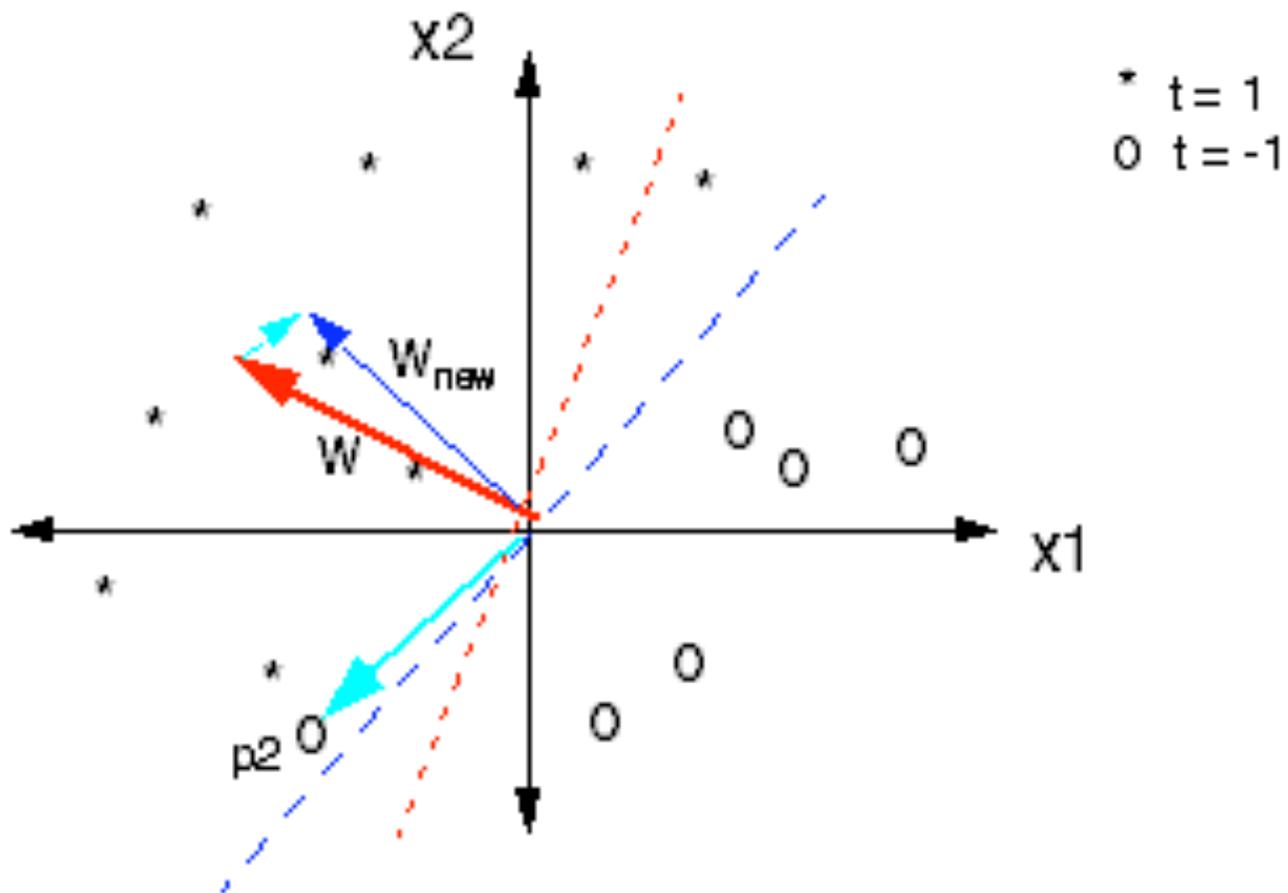
until no error is found

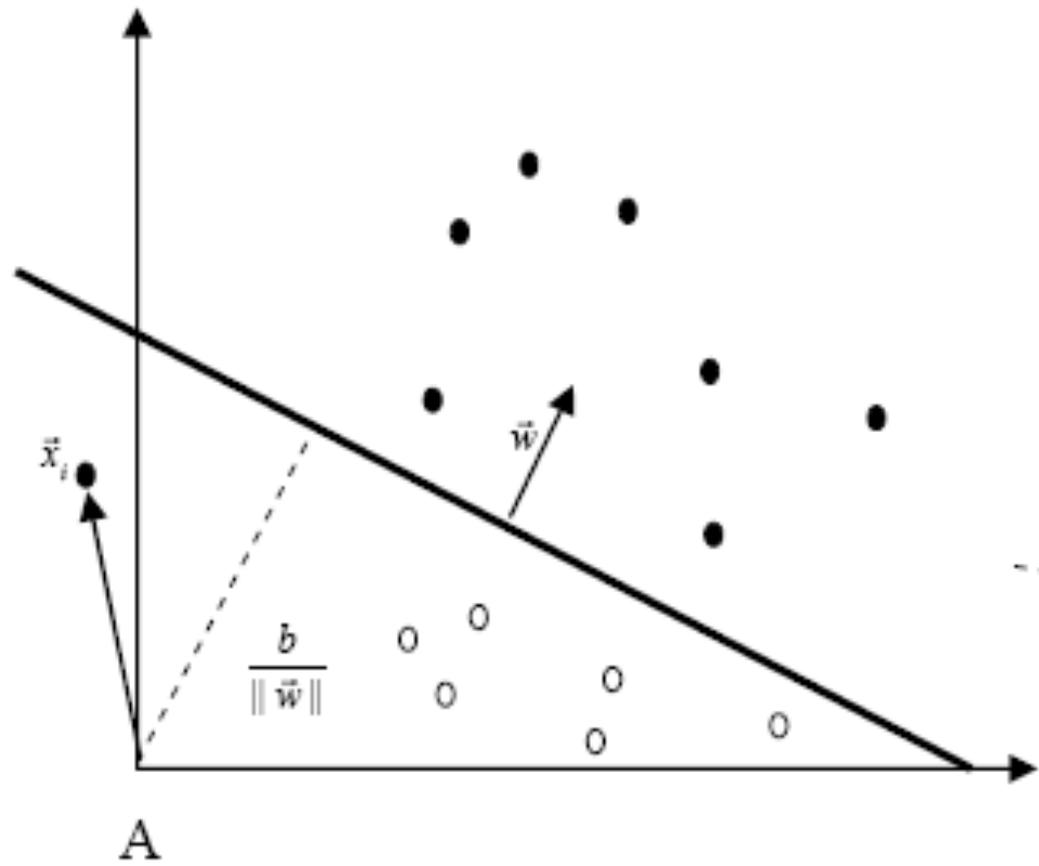
return  $k, (\vec{w}_k, b_k)$

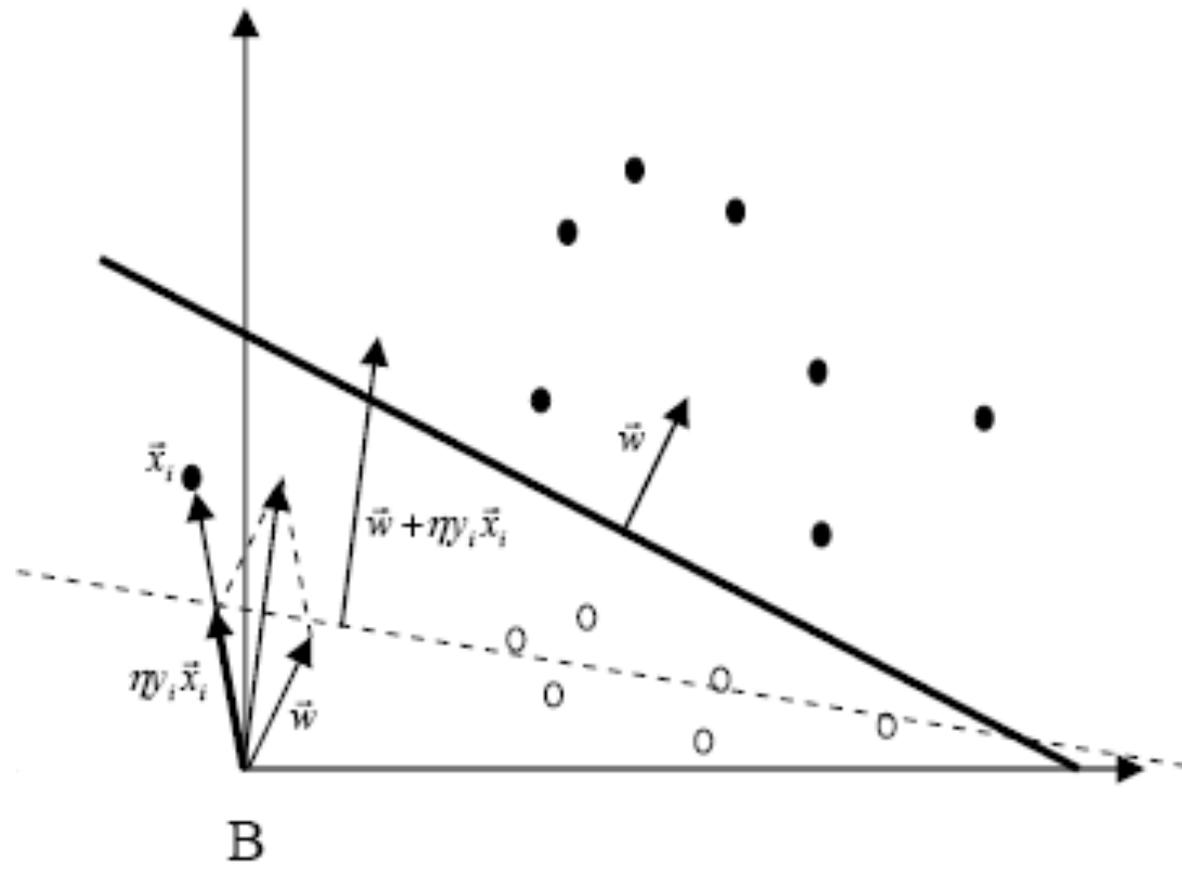


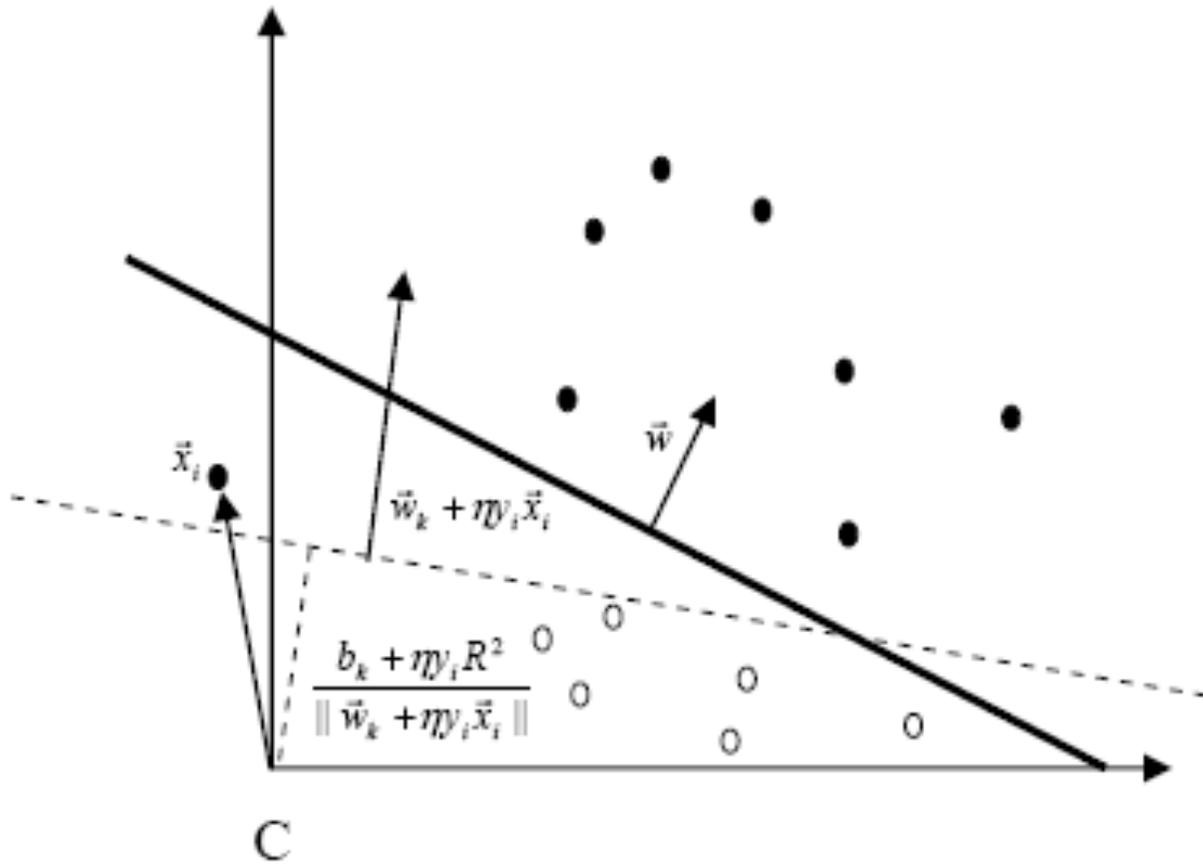












# Novikoff's Theorem

---

Let  $S$  be a non-trivial training-set and let

$$R = \max_{i=1,\dots,m} \|x_i\|.$$

Let us suppose there is a vector  $\mathbf{w}^*$ ,  $\|\mathbf{w}^*\| = 1$  and

$$y_i(\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*) \geq \gamma, \quad i = 1, \dots, m,$$

with  $\gamma > 0$ . Then the maximum number of errors of the perceptron is:

$$t^* = \left( \frac{2R}{\gamma} \right)^2,$$



# Observations

---

- The theorem states that independently of the margin size, if data is linearly separable the perceptron algorithm finds the solution in a finite amount of steps.
- This number is inversely proportional to the square of the margin.
- The bound is invariant with respect to the scale of the *patterns* (i.e. only the relative distances count).
- The learning rate is not essential for the convergence.



# Dual Representation

---

- The decision function can be rewritten as:

$$h(x) = \text{sgn}(\vec{w} \cdot \vec{x} + b) = \text{sgn}\left(\sum_{j=1..m} \alpha_j y_j \vec{x}_j \cdot \vec{x} + b\right) =$$

$$\text{sgn}\left(\sum_{i=1..m} \alpha_j y_j \vec{x}_j \cdot \vec{x} + b\right)$$

- as well as the updating function

$$\text{if } y_i \left( \sum_{j=1..m} \alpha_j y_j \vec{x}_j \cdot \vec{x}_i + b \right) \leq 0 \text{ then } \alpha_i = \alpha_i + \eta$$

- The learning rate  $\eta$  only affects the re-scaling of the hyperplane, it does not affect the algorithm, so we can fix  $\eta = 1$ .



# First properties of SVMs

---

- **DUALITY** is the first feature of Support Vector Machines
- SVMs are learning machines using the following function:

$$f(x) = \text{sgn}(\vec{w} \cdot \vec{x} + b) = \text{sgn}\left(\sum_{j=1..m} \alpha_j y_j \vec{x}_j \cdot \vec{x} + b\right)$$

- Note that data appears only as scalar product (for both testing and learning phases)
- The Matrix  $G = \left(\vec{x}_i \cdot \vec{x}_j\right)_{i,j=1}^m$  is called Gram matrix



# Limits of Linear Classifiers

---

- Data must be linearly separable
- Noise (almost all classifier types)
- Data must be in vectorial format



# Solutions

---

- **Multi-Layers Neural Network:** back-propagation learning algorithm.
- **SVMs:** kernel methods.

The learning algorithm is decoupled by the application domain which is encoded by a kernel function

