



UNIVERSITY OF LEEDS

---

**Class: Machine Learning**

**Neural Networks: Perceptron**

**Instructor: Matteo Leonetti**

# Learning outcomes

---



UNIVERSITY OF LEEDS

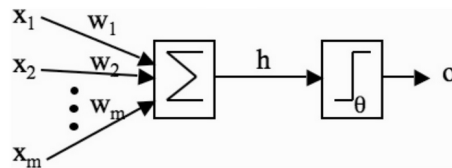
- Define **linear separability**.
- Justify whether a given **error function** is suitable for gradient descent.

# Recap

We want to apply gradient descent:

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

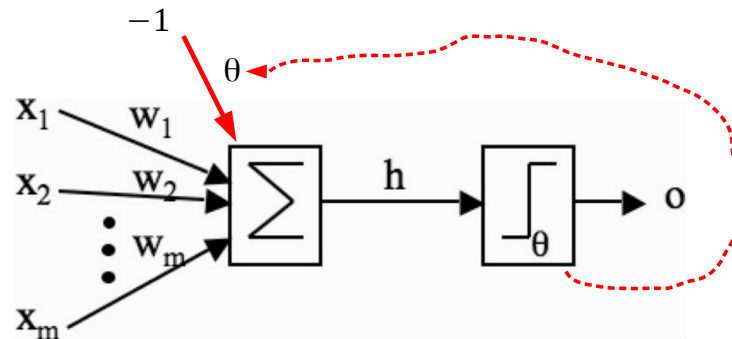
To the parameters of a perceptron:



So as to minimise **an error (or loss) function**, such as:

$$E(\mathbf{X}) = \sum_{\vec{x}_n \in \mathbf{X}} |y_n - t_n|$$

# Bias input



$$\begin{aligned}
 h_w(\mathbf{x}) &= \sum_i w_i x_i = \mathbf{w} \cdot \mathbf{x} \\
 o(h_w) &= \begin{cases} 1 & \text{if } h_w > \theta \\ 0 & \text{if } h_w \leq \theta \end{cases} \\
 &\quad w_1 x_1 + w_2 x_2 + \dots + w_n x_n - 1 \cdot \theta > 0
 \end{aligned}
 \longrightarrow
 \begin{aligned}
 h_w(\mathbf{x}) &= \sum_i w_i x_i - \theta \\
 \mathbf{x}_{new} &= \langle \mathbf{x}, -1 \rangle \quad \mathbf{w}_{new} = \langle \mathbf{w}, \theta \rangle \\
 h_w(\mathbf{x}_{new}) &= \mathbf{w}_{new} \cdot \mathbf{x}_{new} \\
 o(h_w) &= \begin{cases} 1 & \text{if } h_w > 0 \\ 0 & \text{if } h_w \leq 0 \end{cases}
 \end{aligned}$$

We want to optimise the parameters of the perceptron through gradient descent.

For this purpose, we will need to compute **the derivative with respect to the parameters.**

With the weights this is simple (as we will see shortly). However, the threshold of the activation function is in an awkward position.

Nonetheless, we can simply move the threshold on the other side of the inequality of the activation function, and treat it as if it were just a normal weight, associated with an additional input, called **the bias input.**

The bias input must **be fixed at a constant.** The value of the constant is not important, since the parameter theta can adjust accordingly. In books, you will usually find either 1 or -1.

## Question

The decision boundary of the perceptron is the function below:

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} = 0$$

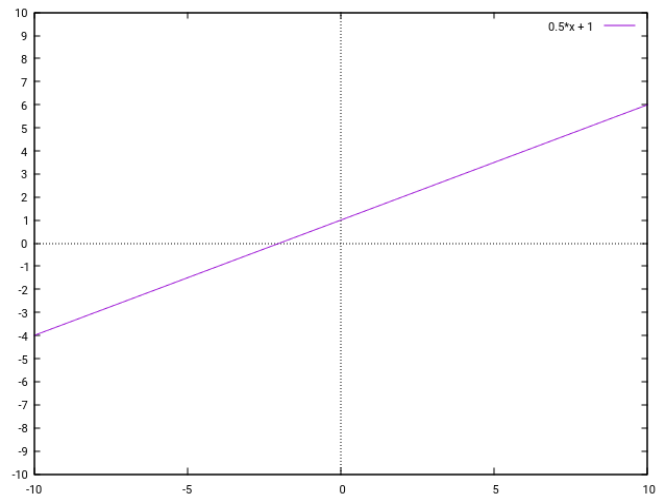
Plot the following function:  $\frac{1}{2}x - y + 1 = 0$

## Solution



UNIVERSITY OF LEEDS

Plot the following function:  $\frac{1}{2}x - y + 1 = 0$



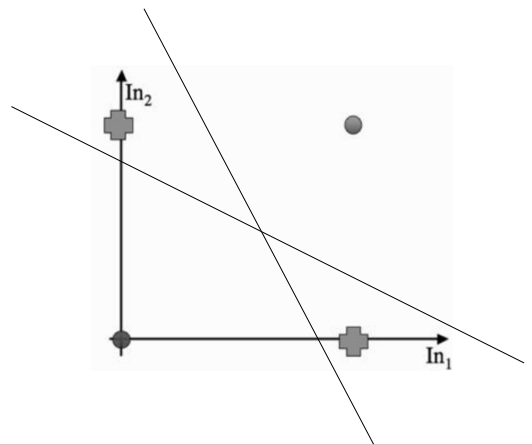
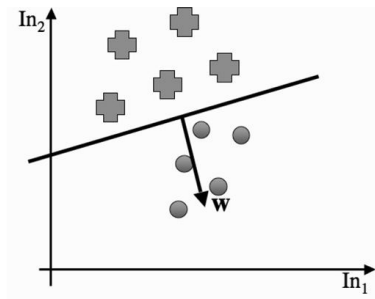
# Linear separability

We have established that the decision boundary is a hyperplane.

$$h_w(x) = w^T x + w_0 = 0$$

XOR

Not linearly separable!



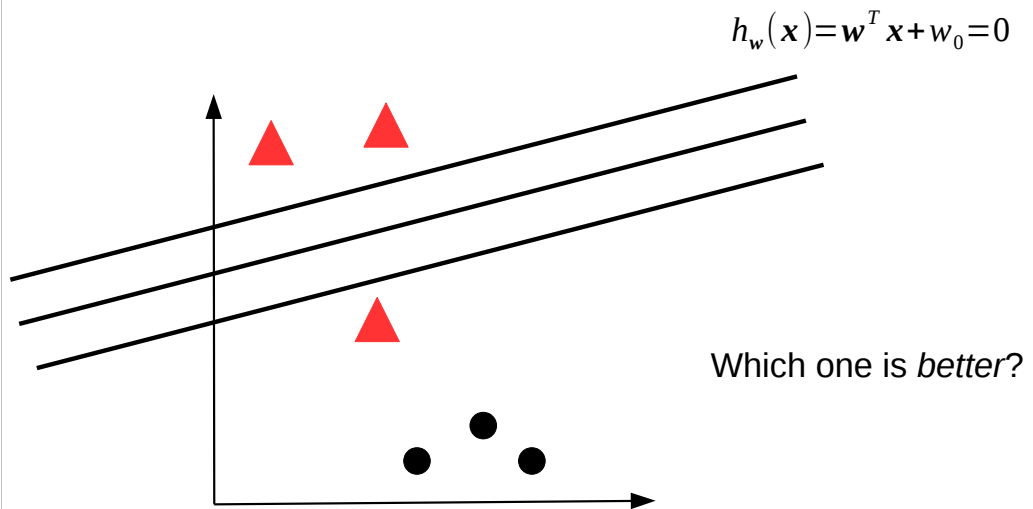
The decision boundary of the perceptron is a hyperplane (in 2D, a straight line).

The vector of weights (excluding the bias weight) is a vector orthogonal to the hyperplane, and plays an important role.

If the dataset is such that an hyperplane that separates the two classes exists, it is said to be linearly separable.

The logic function XOR over points in  $\{0,1\} \times \{0,1\}$  is a well-known example of a non-linearly separable dataset.

## Number of mistakes as Error



If we use the number of mistakes as an error function, all the lines drawn here will have the same error (since they all make the same number of mistakes).

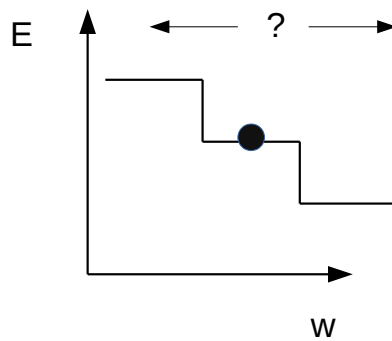
Therefore, it is not possible to use the number of mistakes to inform the search with a local optimization method. Most local steps achieve no difference in the objective (in our case the error) function, and it is impossible to know whether the algorithm is progressing towards a minimum or not.



# Number of mistakes

$$E(\mathbf{X}) = \sum_{\vec{x}_n \in \mathbf{X}} |y_n - t_n|$$

Number of mistakes on the dataset. Piecewise constant  $\rightarrow$  no gradient.

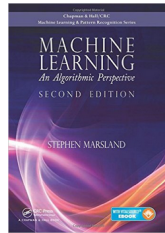


There is no local information on the direction of improvement

Rather than this function, which is piecewise constant, we would like an error function in which the error decreases as the current solution gets closer to a misclassified point.



## Conclusion



## Section 3.3 & 3.4