

PERCEPTRON (F. ROSENBLATT 1957)

①

↳ LINEAR BINARY Classifier

↳ hyperplane H separates 2 classes

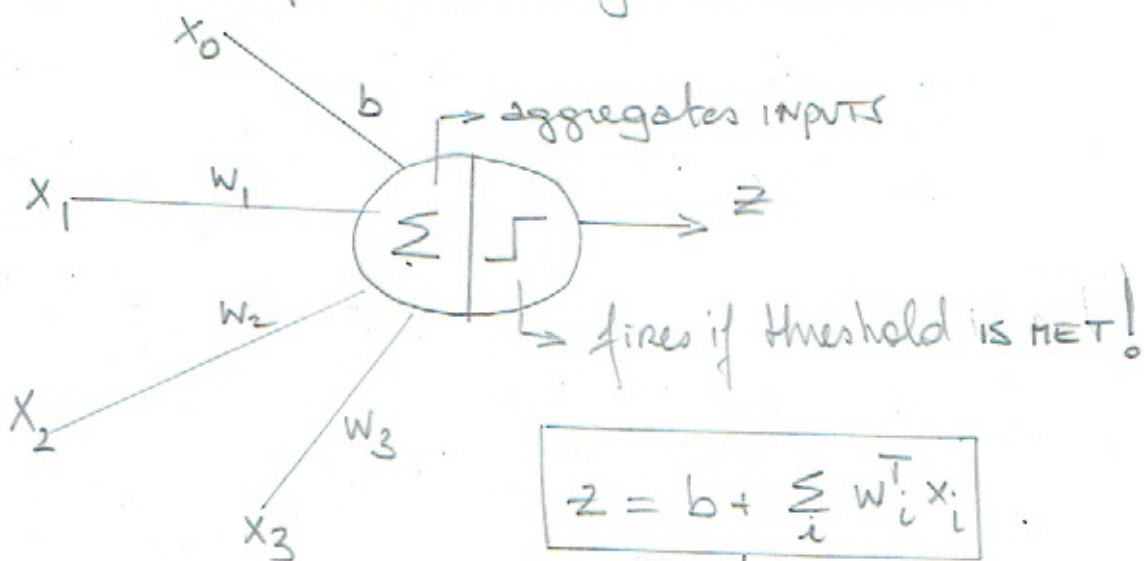
$\oplus y=1$

$\ominus y=0$

↳ if there exists a hyperplane H that separates the 2 classes you will find it!

→ there are infinitely many!

→ Perceptron Convergence Theorem

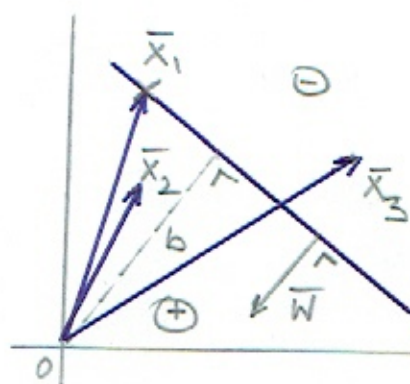


$$z = b + \sum_i w_i^T x_i$$

set $x_0 = 1$
 $b = w_0$

$$z = \sum_i w_i^T x_i$$

↳ - Threshold



↳ distance from ORIGIN.

↳ points to positive class by convention!
unit weight vector

on H : $-w^T x_1 = b \rightarrow w^T x_1 + b = 0 \rightarrow$ per convention belong to \ominus class

$-w^T x_2 < b \rightarrow w^T x_2 + b > 0$

$-w^T x_3 > b \rightarrow w^T x_3 + b < 0$

if H through origin

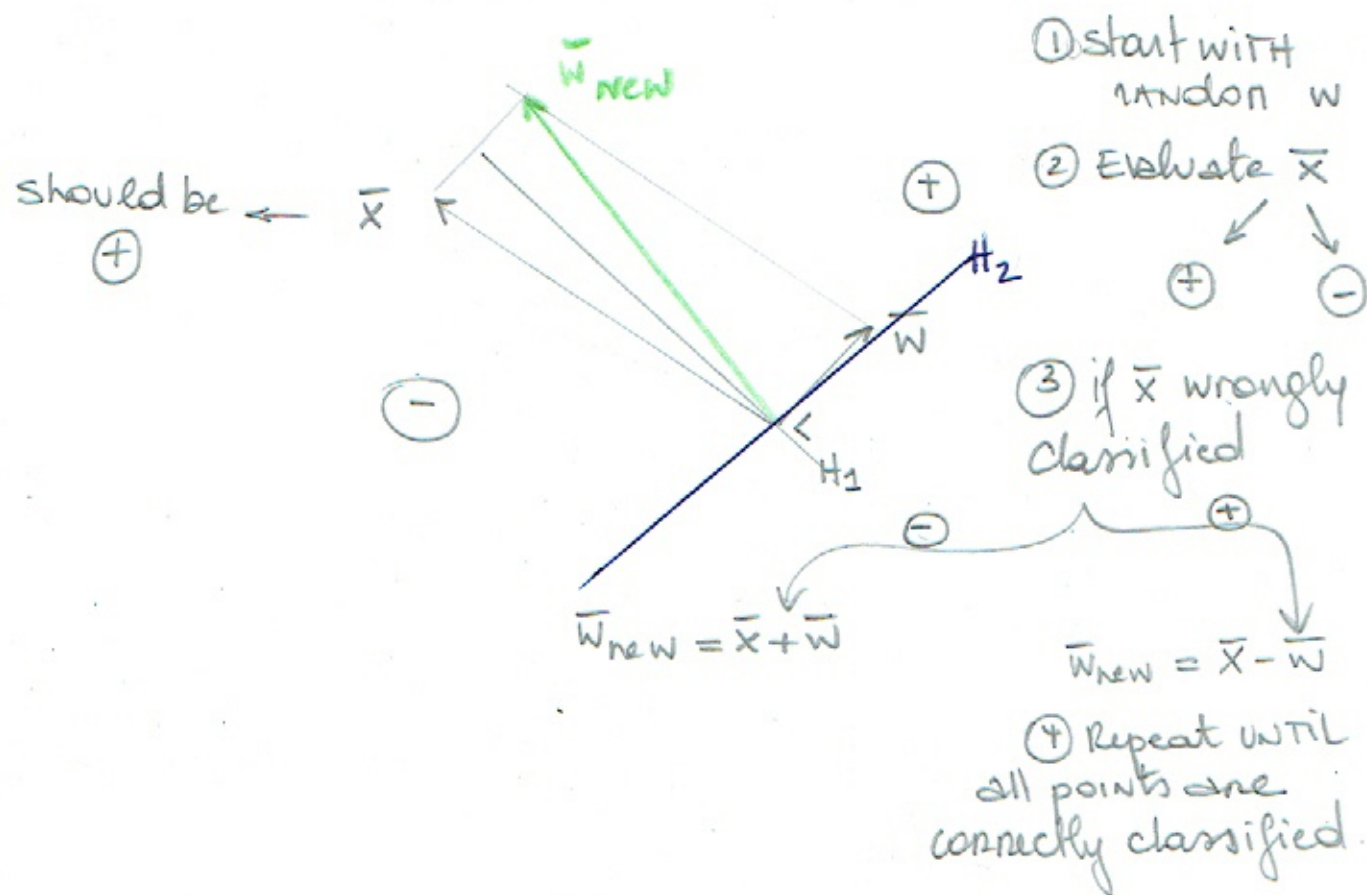
$$\text{ON } H: w^T x = 0$$

$$w^T x > 0 \quad \oplus \text{ class}$$

$$w^T x \leq 0 \quad \ominus \text{ class}$$

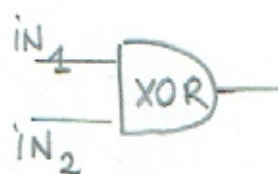
perceptron learning Algorithm

(2)

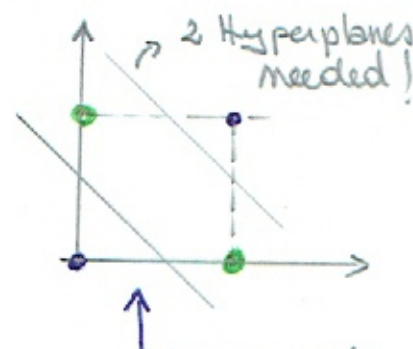


Limitations of perceptron → Minsky & Papert 1969

XOR problem XOR = exclusive OR port



IN_1	IN_2	OUT
0	0	0
1	1	0
0	1	1
1	0	1



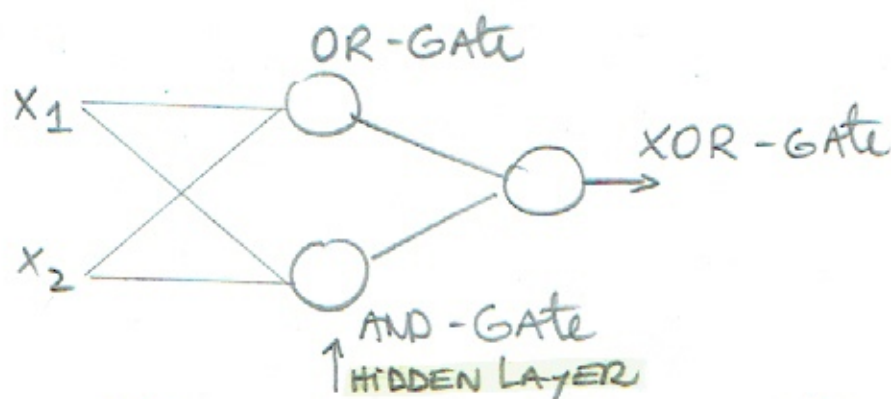
NEED FOR NON-LINEARITY → you cannot separate classes with 1 linear hyperplane!

Solution →

Multi-layer Perceptron

↳ provides non-linearity through introduction of HIDDEN LAYERS + Activation functions

or Hebb's activation function



AND, OR Gates can be trained by perceptron learning rule!
 HOWEVER whole network CANNOT!

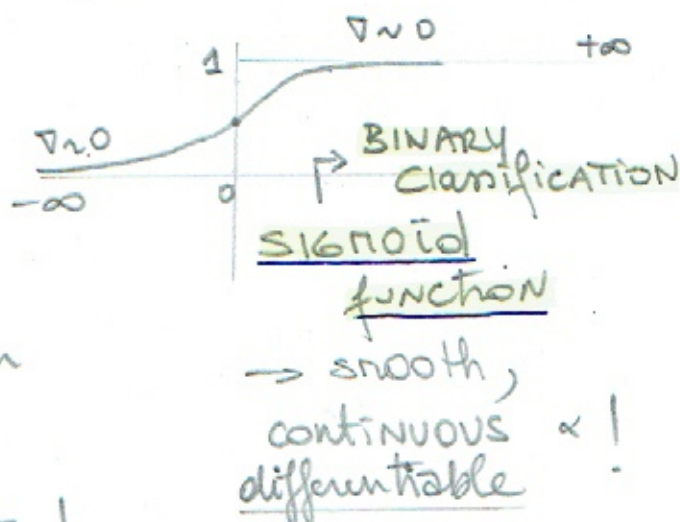
Activation Function

①

$$f(x) = \frac{1}{1 + e^{-x}}$$

↓
 at output hidden layer

↓
 NON-LINEARITY!



but: $e^{-x} \rightarrow$ expensive

outputs close to 0 and 1 \rightarrow kills gradient descent!
 linear score \rightarrow probability!

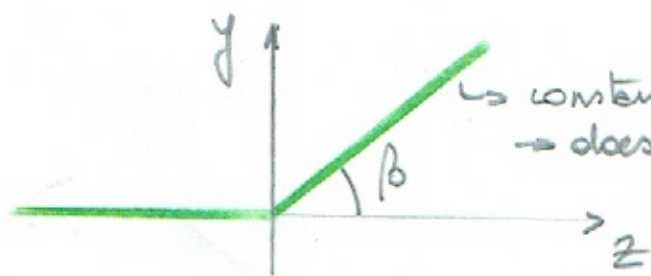
② Softmax \rightarrow multi-class classification

\hookrightarrow VECTOR of K real numbers \rightarrow probability distribution of K probabilities

③ Relu \rightarrow rectified linear unit

VERY POPULAR!

for hidden layers!



\hookrightarrow constant ∇

\rightarrow does not saturate in positive region!

$$y = \max(0, z) \quad \begin{cases} z < 0 \rightarrow 0 \\ z > 0 \rightarrow y \end{cases}$$

$$y = \max(0, w^T x + b)$$

less computationally expensive
 Faster because constant ∇