### 2. PERCEPTRONS

Stephan Robert-Nicoud HEIG-VD/HES-SO



HE

IG



APPRENTISSAGE PAR RESEAUX DE NEURONES ARTIFICIELS

### **Objectives**

- Understand what motivates the inspiration from biological neural networks
- Study the basic models of artificial neurons and their functioning
- Understand what does an artificial neuron do from the mathematical point of view
- Understand how does an artificial neuron "learn" from training data







	von Neumann computer	Brain
processor	Complex High frequency (GHz) one or several (cores)	Simple (neuron) low frequency (Hz) Many!
memory	Separated from the processor Address-based access	Integrated to computing (distributed) Content-addressable
computing	centralized sequential Programmable	Distributed Parallel Learning
reliability	Highly vulnerable	Very robust



APE 2025





Connectionism: modeling of mental or behavioral phenomena as emergent processes of interconnected networks of simple units

#### **Connectionist systems**

We ideally use a massively parallel architecture where each computational element (neuron) performs a sort of a correlation between inputs and stored values called synaptic weights.

Paradigm shift: we replace programming by learning

Réseau de neurones links = weights inputs output APE 2025

#### The artificial neuron

MCulloch-Pítts (1948)



#### **Threshold logic**

We can interpret that excitatory inputs are weighted by 1 and the inhibitory inputs are weighted by -1





NON

if  $(x1 + x2) \ge 2$  then y = 1, otherwise y = 0

### **Rosenblatt's Perceptron (1958)**



If  $(w_1x_1 + w_2x_2 + ... + w_nx_n) \ge \Theta$  then y = 1, otherwise y = -1

If  $\Sigma w_i x_i \ge \Theta$  then y = 1, otherwise y = -1

HE<sup>™</sup> TG

#### Perceptron



#### NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo of Computer Designed to Read and Grow Wiser

WASHINGTON, July 7 (UPI) —The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen.,

The service said it would use this principle to build the first of its Perceptron thinking machines that will be able to read and write. It is expected to be finished in about a year at a cost of \$100.000.

The New York Times, July 8 1958

# What a Perceptron can do ? geometric interpretation

The activation function is discontinuous, what happens when the sum of weighted inputs equals the threshold ?



2D case:  $w_1x_1 + w_2x_2 = \Theta$ , therefore  $x_2 = -(w_1/w_2)x_1 + \Theta/w_2$ 



# The Perceptron as a pattern classifier

A Perceptron is capable of linearly separating the space into two distinct regions. In 2D the separation is a line, in 3D it is a plane and in n-dimensions it is a hyperplane



Learning problem: in a classification problem, it consists on finding the weights and threshold values to correctly separate blue from red data points

HE

#### Bias: a learnable threshold of activation

We would like to get rid of the threshold



If  $\Sigma w_i x_i \ge \Theta$  then y = 1, otherwise y = 0

Let's introduce Wo =  $-\Theta$ :

#### or

If 
$$\Sigma w_i x_i - \Theta \ge 0$$
 then  $y = 1$ ,  
otherwise  $y = 0$ 

If 
$$\Sigma w_i x_i + w_0 \ge 0$$
 then  $y = 1$ ,  
otherwise  $y = 0$ 

Wo is called a "bias" and can be seen as an extra weight that can be learned using a learning algorithm APE 2025

#### **Biasing the output for zero inputs**





HE" IG





### **Perceptron learning algorithm**

- 1. Randomly initialize weights
- 2. Compute the neuron's output  $\mathbf{y}$  for a given input vector  $\mathbf{x}$

 $y = \sum_{j} w_j x_j$ 

3. Update weights:  $W_j(t+1) = W_j(t) + \eta(d-y)x$  //so-called delta rule

d is the desired output and  $\eta$  is the learning rate, 0.0 <  $\eta$  < 1.0

4. Repeat 2 and 3 for a given number of steps or until the error is smaller than a given threshold, i.e., error <  $1/2 \sum_{p} (y_{p-d_p})^2$ 

### Learning in action (weight updating)

1) No classification error 2) Error correction after a weight update







 $W_j(t+1) = W_j(t) + \eta(d-y)\mathbf{x}$ 





#### **Gradient descent**

#### Widrow-Hoff algorithm / Delta rule

The activation function of the Perceptron has a discontinuity, thus it is not derivable at that point. By using a linear activation function, Widrow and Hoff introduced a new model called ADALINE or ADAptive LINear Element



So, what is learning ? : the idea is to find the weights Wi that minimize an error function E, such as this one:

$$E = 1/2\sum_{i=1}^{p} ||y_i - d_i||^2, y = \sum_{j=1}^{p} w_j x_j$$

A way to do this exploits the first derivative of E with respect to the weights. HE<sup>\*\*</sup> This method is called gradient descent, and dates back to Cauchy (1847). APE 2025

# **Gradient descent (2)** Widrow-Hoff algorithm / Delta rule



ΗE

Set of weights minimizing the error

$$E = \frac{1}{2} \sum_{i=1}^{P} ||y_i - d_i||^2, y = \sum_j w_j x_j$$
$$\nabla E = \frac{\partial E}{\partial w} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial w}$$
$$\frac{\partial E}{\partial y} = y - d$$
$$\frac{\partial W_j}{\partial w_j} = \frac{\partial}{\partial w_j} \sum_j w_j x_j = x_j$$
$$\Delta w_j = -\nabla E = \eta (d - y) x_j$$

#### **Generalized Delta rule**



 $y_{j} = \frac{1}{1 + e^{-S_{j}}}, \quad S_{j} = \sum w_{i}x_{i}$  $\Delta w_{j} = \eta(d - y)y_{j}'x_{j}$  $y_{j}' = y_{j}(1 - y_{j})$ 

The tanh and ReLu are other widely used activation functions

APE 2025

# sigmoid function and derivative
def sigmoid(neta):
 output = 1 / (1 + np.exp(-neta))
 d\_output = output \* (1 - output)
 return (output, d\_output)

# input dataset X = np.array([ [0.2,0.2], [0.2,0.5], ..... [0.8,0.7], [0.8,0.9] ])

# output classes y = np.array([[0,0,0,0, ..... 1,1,1,1]]).T

# seed random numbers to make calculation
# deterministic (just a good practice)
np.random.seed(1)

# initialize weights randomly with mean 0
weights = np.random.normal(size=2)
bias = np.random.normal(size=1)

HE IG

#### inputs = Xtarget = y

for iter in xrange(100):

# forward propagation
neta = np.dot(inputs, weights) + bias
output, d\_output = sigmoid(neta)

# how much did we miss?
error = target - output

# learning rate alpha = 0.1

d\_w\_x = alpha \* error \* d\_output \* inputs[:,0] d\_w\_y = alpha \* error \* d\_output \* inputs[:,1] d\_b = alpha \* error \* d\_output

# update weights
weights += np.array([np.sum(d\_w\_x), np.sum(d\_w\_y)])
bias += np.sum(d\_b)

print "Network After Training:" print weights, bias

### Activation functions (1)



# **Activation functions (2)**

~

0.0

Input value

 $\mathbf{v}$ 

0.0

Input value

2.5

5.0

0.51

2.5

- output

7.5

— output

7.5

first derivative

10.0

5.0

first derivative

10.0

0.61







APE 2025

# A two-class classifier example (1)





## A two-class classifier example (2)



### A two-class classifier example (3)



#### **Two-class classification threshold**

- □ A threshold has to be fixed by the engineer
- If we have a validation dataset, we can evaluate the effect of such a threshold and decide the best value
- In general, there is a trade-off between the true positive rate (sensitivity) and the false positive rate (1 - specificity) of such classifier.

# A linearly separable three-class problem





HE<sup>™</sup> IG



□ It is not possible to draw a line to separate these two classes, e.g., red points from blue points.



APE 2025