# Artificial Neural Networks (ANNs)

Biologically inspired computational model:

- (1) <u>Simple</u> computational units (neurons).
- (2) <u>Highly interconnected</u> connectionist view
- (3) Vast <u>parallel</u> computation, consider:
  - Human brain has ~10<sup>11</sup> neurons
  - Slow computational units, switching time ~10<sup>-3</sup> sec (compared to the computer >10<sup>-10</sup> sec)
  - Yet, you can recognize a face in ~10<sup>-1</sup> sec
  - This implies only about 100 sequential, computational neuron steps this seems too low for something as complicated as recognizing a face
  - Parallel processing

ANNs are naturally parallel - each neuron is a self-contained computational unit that depends only on its inputs.



#### Learning

We have seen machine learning with different representations:

- Decision trees -- symbolic representation of various decision rules -- "disjunction of conjunctions"
- (2) Perceptron -- learning of weights that represent alinear decision surface classifying a set of objects into two groups



Different representations give rise to different <u>hypothesis</u> or <u>model spaces</u>. Machine <u>learning algorithms search</u> these model spaces for the <u>best fitting</u> <u>model</u>.

#### The Perceptron

- A simple, single layered neural "network" - only has a single neuron.
- However, even this simple neural network is already powerful enough to perform (linear) classification tasks.



#### The Architecture



#### Computation

A perceptron computes the value,

$$y = \operatorname{sgn}\left(b + \sum_{i=1}^{m} w_i x_i\right)$$

Ignoring the activation function sgn and setting m = 1, we obtain,

$$y' = b + w_1 x_1$$

But this is the equation of a line with slope *w* and offset *b*.

<u>Observation</u>: For the general case the perceptron computes a <u>hyperplane</u> in order to accomplish its classification task,

$$y' = b + \sum_{i=1}^{m} w_i x_i = b + \vec{w} \cdot \vec{x}$$

# Perceptron Learning Revisited





Constructs a line (hyperplane) as a classifier

$$b + \vec{w} \cdot \vec{x}$$

#### Classification



In order for the hyperplane to become a classifier we need to find b and w => learning!

#### What About Non-Linearity?



Can we learn this decision surface? ... Yes! Multi-Layer Perceptrons

# Multi-Layer Perceptrons (ANNs)



#### How do we train?

Perceptron was easy:



Every time we found an error of the predicted value  $f(x_i)$  compared to the label in the training set  $y_i$ , we update w and b.

Not so easy in multi-layer neural networks – the error can occur deep in the network!

#### **Artificial Neural Networks**

#### Feed-forward with Backpropagation

We have to be a bit smarter in the case of ANNs: compute the signal (feed forward) and then use the error at the output to update all the weights by propagating the error back through the network.



## Backpropagation



#### Note: These computations only work if the transfer functions are differentiable.

#### **Neural Network Learning**

 Define the network error in terms of weights w as

 $E = (y' - y)^2$ 

for some training instance *x*.

 Use the gradient (slope) of the error surface to guide the search towards appropriate weights:

$$F_{W_0}^{25}$$

$$\Delta w = -\frac{\partial E}{\partial w}$$

Backpropagation can be understood as a <u>stochastic gradient</u> <u>search</u> on the error surface of the network.

# **Backpropagation Algorithm**

Note: this algorithm is for a NN with a single output node o and a single hidden layer. It can easily be generalized. Initialize the weights in the network (often randomly) Do For each example e in the training set // forward pass y = compute neural net output y' = label for e from training data Calculate error E = (y' - y)<sup>2</sup> at the output units // backward pass Compute error  $\delta_0$  for weights from a hidden node h to the output node o using E Compute error  $\delta_h$  for weights from an input node i to hidden node h using  $\delta_0$ Update the weights in the network Until all examples classified correctly or stopping criterion satisfied Return the network

#### **Representational Power**

- Every bounded continuous function can be approximated with arbitrarily small error by a network with one hidden layer.
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers.



# Hidden Layer Representations

**Target Function:** 



Input		Output
10000000	$\rightarrow$	1000000
01000000	$\rightarrow$	01000000
00100000	$\rightarrow$	00100000
00010000	$\rightarrow$	00010000
00001000	$\rightarrow$	00001000
00000100	$\rightarrow$	00000100
00000010	$\rightarrow$	00000010
00000001	$\rightarrow$	0000001

Can this be learned?

## Hidden Layer Representations

Inputs	Outputs	Input	Hidden	Output	
0	P				
0ff	Æρ	$ 10000000 \rightarrow$	$.89$ $.04$ $.08$ $\rightarrow$	1000000	
O		$01000000 \rightarrow$	.01 .11 .88 $\rightarrow$	01000000	
		$00100000 \rightarrow$	.01 .97 .27 $\rightarrow$	00100000	
		$00010000 \rightarrow$	$.99$ $.97$ $.71$ $\rightarrow$	00010000	111
		$00001000 \rightarrow$	$.03$ $.05$ $.02$ $\rightarrow$	00001000	000
		$00000100 \rightarrow$	$.22$ .99 .99 $\rightarrow$	00000100	011
0	$\checkmark$	$00000010 \rightarrow$	$.80$ $.01$ $.98$ $\rightarrow$	00000010	
		$0000001 \rightarrow$	.60 .94 .01 $\rightarrow$	00000001	110

This neural network architecture is sometimes also called <u>autoencoder</u> because of its ability to invent new representations of the input data and is a popular building block in deep-learning.