

1 Neural networks (for classification or regression)

Used to estimate unknown functions (complex, non-linear hypothesis) that are based on a large number of inputs, through the back-propagation algorithm.

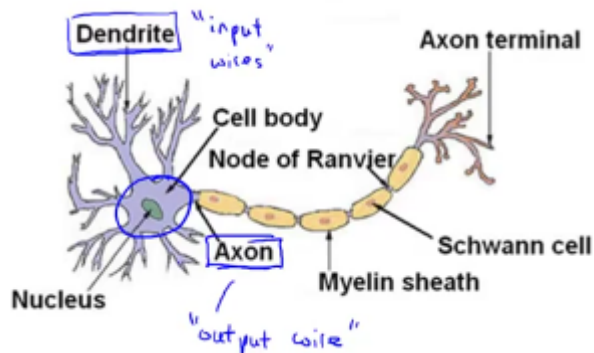
Generally more complex and computationally expensive than other methods, but powerful for certain problems.

The basis of many deep learning methods.

Today is the state of the art technique for many different machine learning problems.

1.1 Description

Neural networks were developed to vaguely simulate the neurons in the brain.



A neuron is a computational unit that receives a number of inputs through its input wires (*dendrites*), does some computation and then sends signals through its output wire (*axon*) to other neurons in the brain.

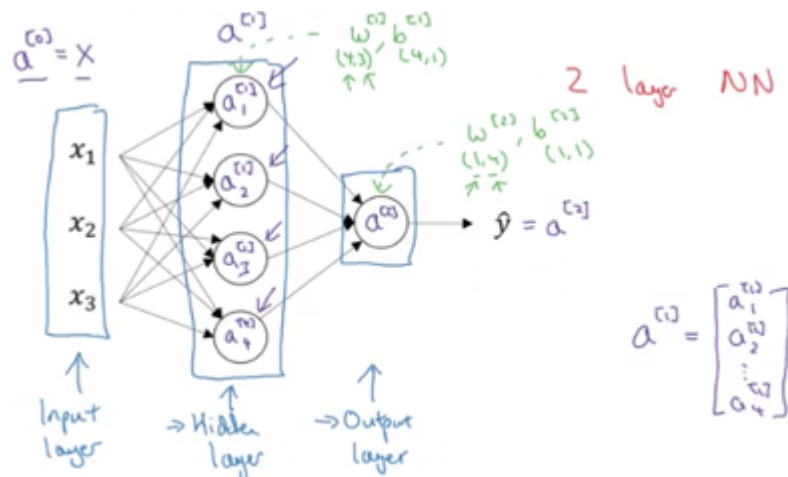
A neural network is a group of neurons.

The inputs are grouped in an *input layer*.

The outputs are grouped in a final *output layer*.

The layers in between are called the *hidden layers*.

Adding more layers helps computing even more complex functions on the input data.



1.2 Notation

- training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$ with m training samples
- each input variable has n features: $x^{(i)} = \{x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}\}$
- output variable $y^{(i)}$, represented by either a single value (in case of regression or binary classification) or by an identity vector (in case of multi-class classification)

- L : the total number of layers in the network (comprising the hidden layers and the output layer)
- $n^{[l]}$: the number of units (neurons) in layer l
- $w^{[l]}$: weights matrix $[n^{[l]} \times n^{[l-1]}]$ controlling function mapping from layer $l - 1$ to layer l ; $w_{ij}^{[l]}$: weight to unit i in layer l from unit j in layer $l - 1$
- $b^{[l]}$: bias vector $[n^{[l]}, 1]$ on layer l ; $b_i^{[l]}$: bias on unit i in layer l
- activation values outputted from layer l

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

- in detail for each unit in the layer

$$z_j^{[l]} = w_{j0}^{[l]} a_0^{[l-1]} + w_{j1}^{[l]} a_1^{[l-1]} + \dots + w_{jn^{[l-1]}}^{[l]} a_{n^{[l-1]}}^{[l-1]} + b_j^{[l]}$$

$$a_j^{[l]} = g^{[l]}(z_j^{[l]})$$

- output $\hat{y} = a^{[L]} = h(x)$

1.3 Neural networks for multi-class classification

- Model's architecture
 - fully connected network
 - L layers
 - input layer with m units (a training sample x)
 - output layer with K units ($\hat{y} = h_{\Theta}(x) \in \mathbb{R}^K$)
 - $L - 1$ hidden layers with s_l units (for $l \in \{1, \dots, L - 1\}$)

- Model's parameters

$$w = \{w^{[1]}, w^{[2]}, \dots, w^{[L]}\}$$

$$b = \{b^{[1]}, b^{[2]}, \dots, b^{[L]}\}$$

- Cost function (depends on the output activation function)

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}, y) = \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left(-y_k^{(i)} \log \hat{y}_k^{(i)} - (1 - y_k^{(i)}) \log (1 - \hat{y}_k^{(i)}) \right)$$

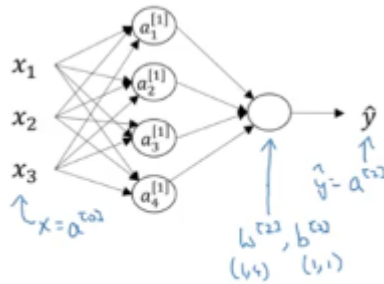
- Goal

$$\min_{w, b} J(w, b)$$

- Algorithm

- start with some initial values for w and b (usually random values in $[-\epsilon, \epsilon]$)
- for each epoch
 - set $a^{[0]} = x$
 - perform forward-propagation to compute $a^{[l]}$ for $l = \{1, 2, \dots, L\}$

- $z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$
- $a^{[l]} = g^{[l]}(z^{[l]})$



Given input x :

$$\begin{aligned} \rightarrow z^{[1]} &= W^{[1]} \underset{(4,1)}{a^{[0]}} + \underset{(1,1)}{b^{[1]}} \\ \rightarrow a^{[1]} &= \sigma(z^{[1]}) \\ \rightarrow z^{[2]} &= W^{[2]} \underset{(1,4)}{a^{[1]}} + \underset{(1,1)}{b^{[2]}} \\ \rightarrow a^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

- perform back-propagation: back propagate the error through each layer

- last layer

$$dz^{[L]} = \frac{\partial J}{\partial z^{[L]}} = a^{[L]} - y$$

$$dw^{[L]} = \frac{\partial J}{\partial w^{[L]}} = \frac{1}{m} dz^{[L]} a^{[L-1]}$$

$$db^{[L]} = \frac{\partial J}{\partial b^{[L]}} = \frac{1}{m} dz^{[L]}$$

- previous layers

$$dz^{[l]} = w^{[l+1]^T} dz^{[l+1]} * g'^{[l]}(z^{[l]})$$

$$dw^{[l]} = \frac{1}{m} dz^{[l+1]} a^{[l-1]}$$

$$db^{[l]} = \frac{1}{m} dz^{[l+1]}$$

- update the weights and biases for every layer

$$w^{[l]} = w^{[l]} - \alpha dw^{[l]}$$

$$b^{[l]} = b^{[l]} - \alpha db^{[l]}$$