

Discussion 11

王欣奕, wangxy6@shanghaitech.edu.cn

Neural Network

Biological inspiration

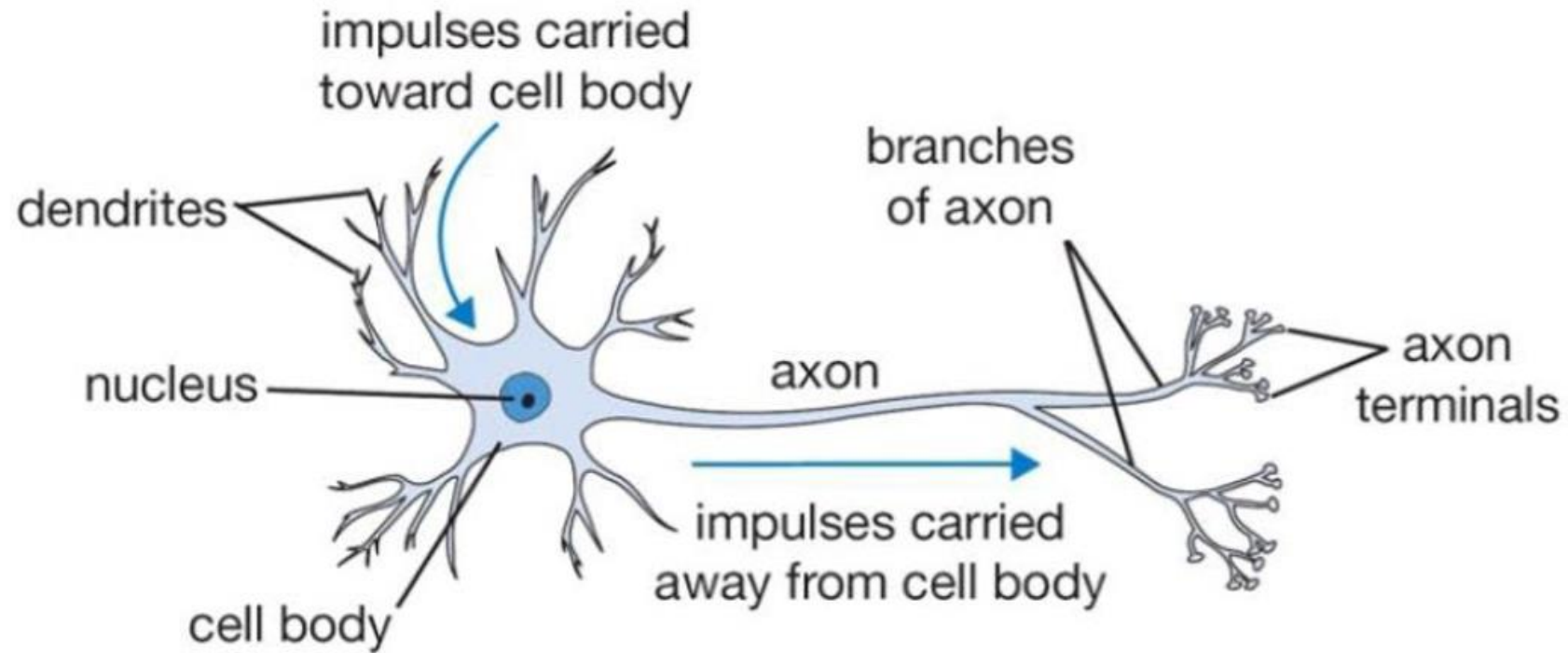
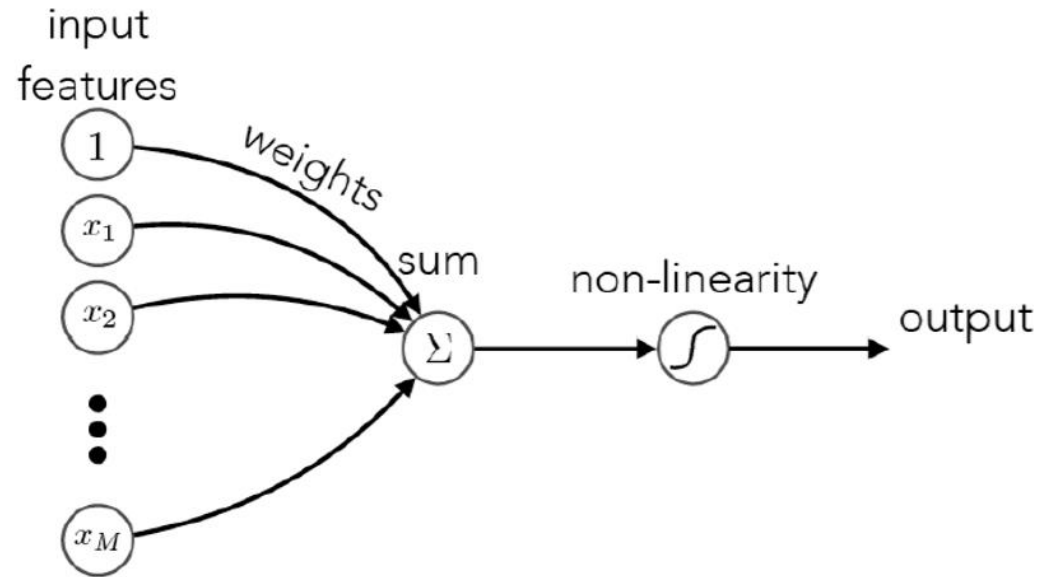


Figure : The basic computational unit of the brain: Neuron

Single neuron



artificial neuron: *weighted sum and non-linearity*

$$s = \underset{\substack{\text{bias} \\ \uparrow}}{b} + \underset{\substack{\text{weights} \\ \uparrow}}{w_1} \underset{\substack{\text{input features} \\ \uparrow}}{x_1} + \underset{\substack{\text{weights} \\ \uparrow}}{w_2} \underset{\substack{\text{input features} \\ \uparrow}}{x_2} + \cdots + \underset{\substack{\text{weights} \\ \uparrow}}{w_M} \underset{\substack{\text{input features} \\ \uparrow}}{x_M} = \mathbf{w}^T \mathbf{x}$$

sum \nearrow

$$h = g(s)$$

output \nearrow non-linearity \nearrow sum

Activation function

Sigmoid: $\sigma(z) = \frac{1}{1+\exp(-z)}$

Tanh: $\tanh(z) = \frac{\exp(z)-\exp(-z)}{\exp(z)+\exp(-z)}$

ReLU (Rectified Linear Unit): $\text{ReLU}(z) = \max(0, z)$

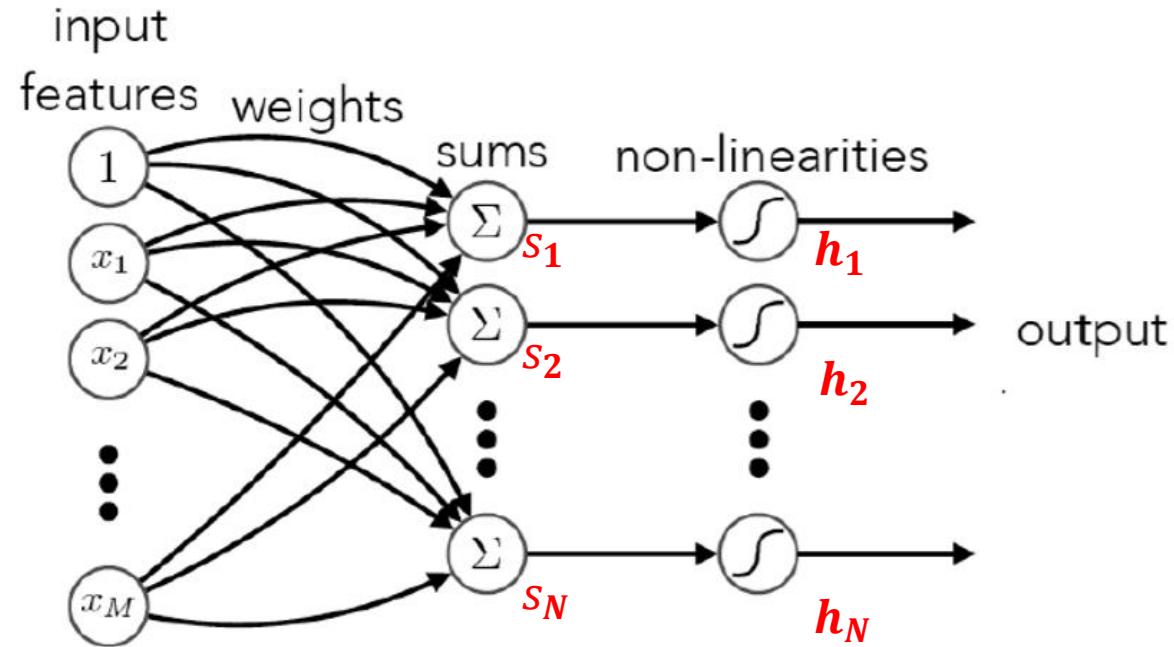
Nice property: $\sigma' = \sigma(1 - \sigma)$

Loss function

$$\text{MLE : } L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)$$

$$\text{MAP : } L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

Single layer network



$$s_1 = w_1^T x$$

$$s_2 = w_2^T x$$

$$\vdots$$

$$s_N = w_N^T x$$

$$h_1 = \sigma(s_1)$$

$$h_2 = \sigma(s_2)$$

$$\vdots$$

$$h_N = \sigma(s_N)$$



$$s = W^T x$$

$$h = \sigma(s)$$

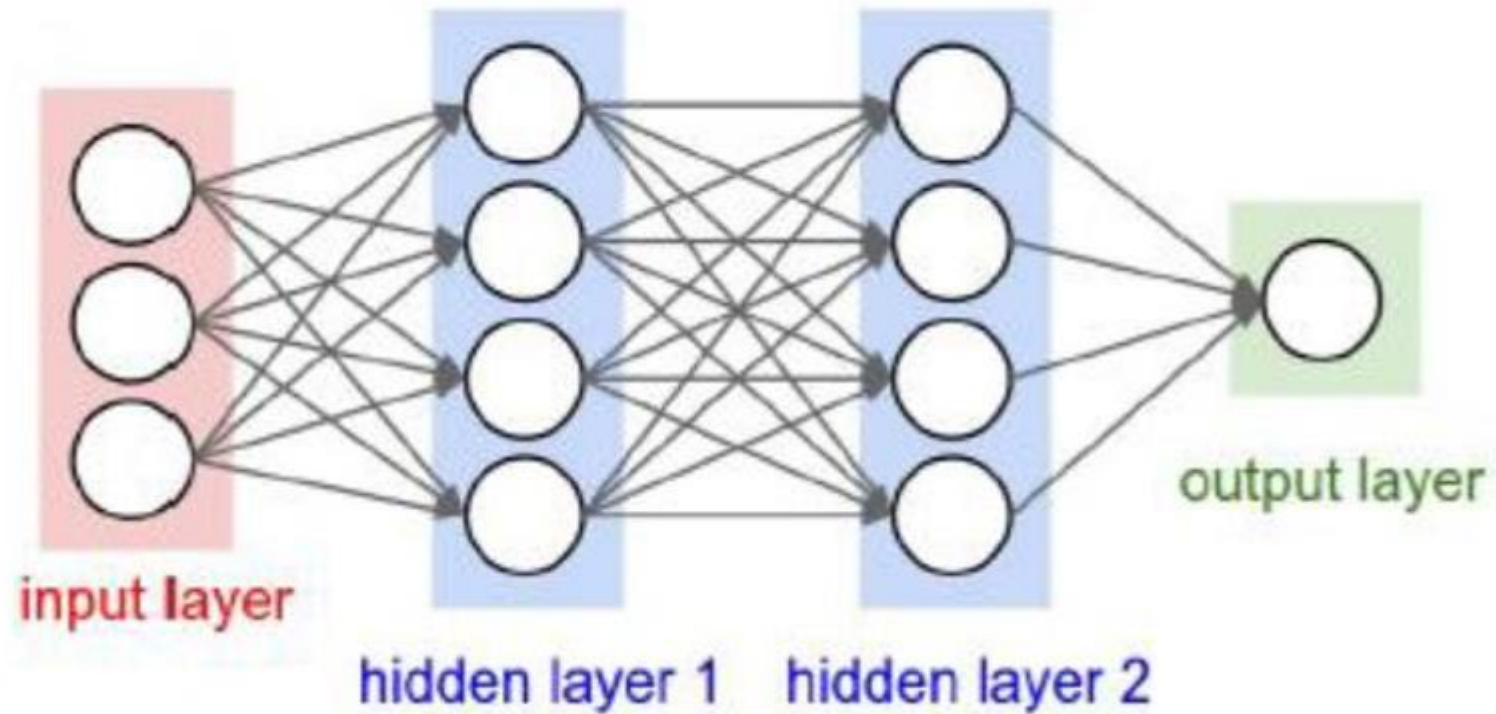
Single layer network

$$\mathbf{s} = \mathbf{W}^T \mathbf{x} \quad \mathbf{h} = \sigma(\mathbf{s})$$

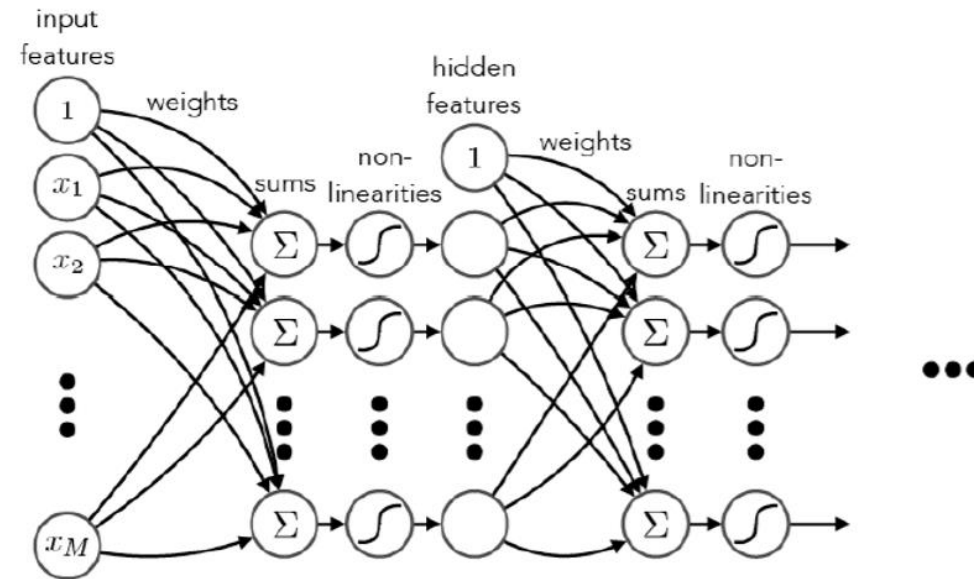
- It can be used in multiclass classification
- \mathbf{h} is a $N \times 1$ vector which represents class scores
- Can also apply softmax function to get normalized probabilities

Multi-layer network

- A **3-layer** neural network



Multi-layer network

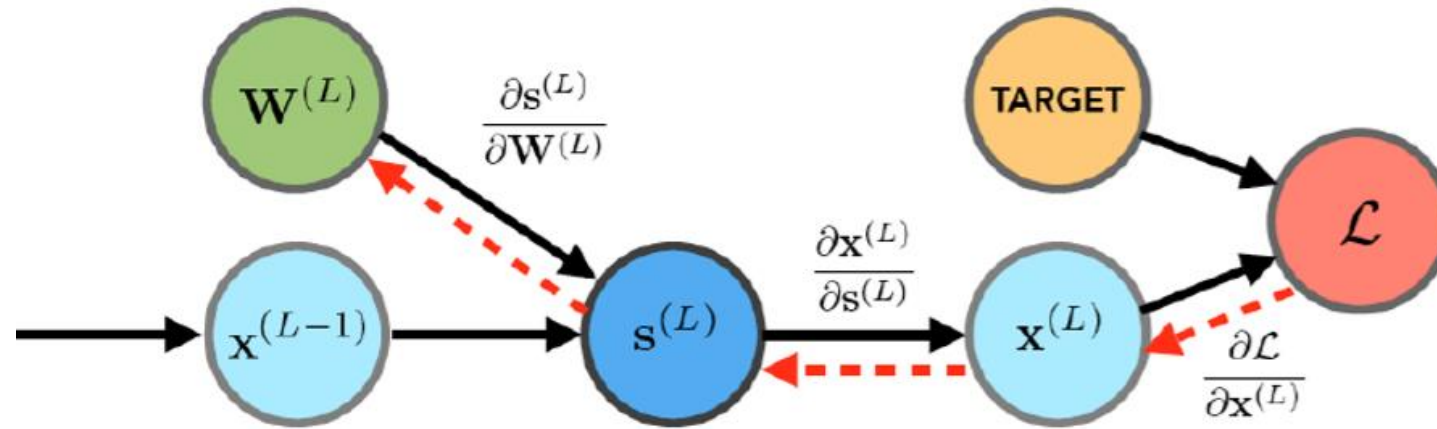


network: sequence of parallelized weighted sums and non-linearities

DEFINE $\mathbf{x}^{(0)} \equiv \mathbf{x}$, $\mathbf{x}^{(1)} \equiv \mathbf{h}$, ETC.

1st layer	2nd layer	
$\mathbf{s}^{(1)} = \mathbf{W}^{(1)} \mathbf{x}^{(0)}$	$\mathbf{s}^{(2)} = \mathbf{W}^{(2)} \mathbf{x}^{(1)}$	
$\mathbf{x}^{(1)} = \sigma(\mathbf{s}^{(1)})$	$\mathbf{x}^{(2)} = \sigma(\mathbf{s}^{(2)})$...

Multi-layer network(Back-propagation)



$$\mathcal{L} = \frac{1}{2} (x^{(L)} - t)^2$$

$$x^{(L)} = \sigma(s^{(L)})$$

$$s^{(L)} = (W^{(L)})^T x^{(L-1)}$$

$$\frac{\partial \mathcal{L}}{\partial W^{(L)}} = \frac{\partial \mathcal{L}}{\partial x^{(L)}} \frac{\partial x^{(L)}}{\partial s^{(L)}} \frac{\partial s^{(L)}}{\partial W^{(L)}}$$

$$x^{(L)} - t$$

$$(x^{(L-1)})^T$$

$$\sigma(s^{(L)}) (1 - \sigma(s^{(L)})) = x^{(L)} (1 - x^{(L)})$$

Multi-layer network(Gradient-based learning)

- Model initialization
- Forward propagate
- Loss function
- Differentiation
- Back-propagation
- Weight update: $\mathbf{W}^{(L)} = \mathbf{W}^{(L)} - \eta \frac{\partial \text{Loss}}{\partial \mathbf{W}^{(L)}}$

Batch GD; Mini-batch GD; SGD