

Back Propagation (1)

1. input and output neurons

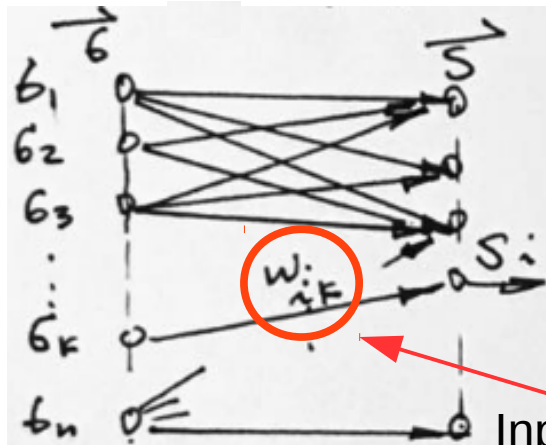
Input

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} n \times 1$$

Output

$$\vec{s} = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_{n_o} \end{bmatrix} n_o \times 1$$

The architecture



Input: k
output: i

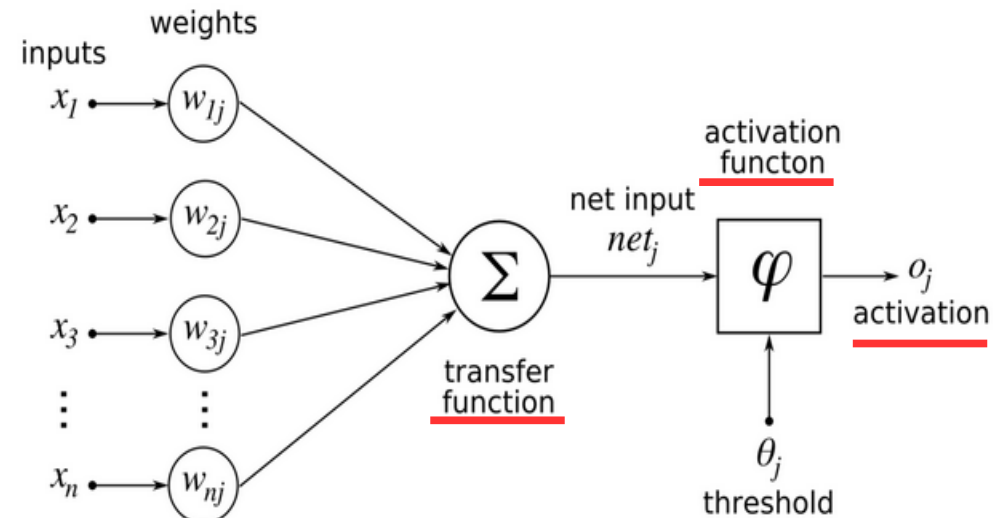
2. weights w_{ik}

i for output

k for input

Other popular notation

$$Y = \sum (\text{weight} * \text{input}) + \text{bias}$$



3. Activation function

$$s_i = f\left(\sum_k w_{ik} x_k\right)$$

... (1)

Note the activation function $f(\cdot)$ and the bias (offset can be written in the unified summation term)

Back Propagation (2)

4. Transfer function h_i

let

$$h_i = \sum_k w_{ik} \phi_k \quad \dots (2)$$

Index i for the output neuron

Neuron output function S_i

$$S_i = f(h_i) = f[h_i(w_{ik})] \quad \dots (2-1)$$

Neuron output is tied to activation function $f(\cdot)$, transfer function h_i and weights w_{ik}

5. Error at each neuron output (the difference between the true output Zeta (desired true output) and the current output S_i) at the experiment μ

$$z_i^\mu - S_i^\mu \quad \dots (2-2)$$

6. total error for all output neuron I and all experiments μ

$$D = \frac{1}{2} \sum_\mu \sum_i (z_i^\mu - S_i^\mu)^2 \quad \dots (3)$$

7. minimize the error wrt to w_{ik}

$$\begin{aligned} \frac{\partial D}{\partial w_{ik}} &= \frac{\partial}{\partial w_{ik}} \cdot \frac{1}{2} \sum_\mu \sum_i [z_i^\mu - f(h_i^\mu)]^2 \\ &= \sum_\mu [z_i^\mu - f(h_i^\mu)] f'(h_i^\mu) \frac{\partial h_i}{\partial w_{ik}} \end{aligned}$$

8. Training the NN by updating the weights

$$w_{ik}(t+1) = w_{ik}(t) + \delta w_{ik}(t) \quad \dots (4)$$

where

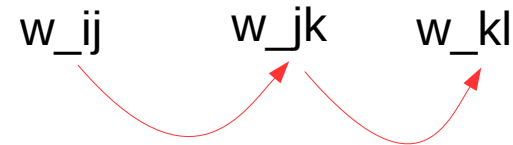
$$\delta w_{ik}(t) = -\epsilon \frac{\partial D}{\partial w_{ik}} \quad \dots (5)$$

Back Propagation (3)

9. computation of the derivatives to update the weight

$$\begin{pmatrix} \frac{\partial D}{\partial w_{i1}} \\ \frac{\partial D}{\partial w_{i2}} \\ \vdots \\ \frac{\partial D}{\partial w_{in}} \end{pmatrix} = \begin{pmatrix} [z_i^u - f(h_i^u)] f' \frac{\partial h_i^u}{\partial w_{i1}} \\ [z_i^u - f(h_i^u)] f' \frac{\partial h_i^u}{\partial w_{i2}} \\ \vdots \\ [z_i^u - f(h_i^u)] f' \frac{\partial h_i^u}{\partial w_{in}} \end{pmatrix}$$

Update sequence back to the front:

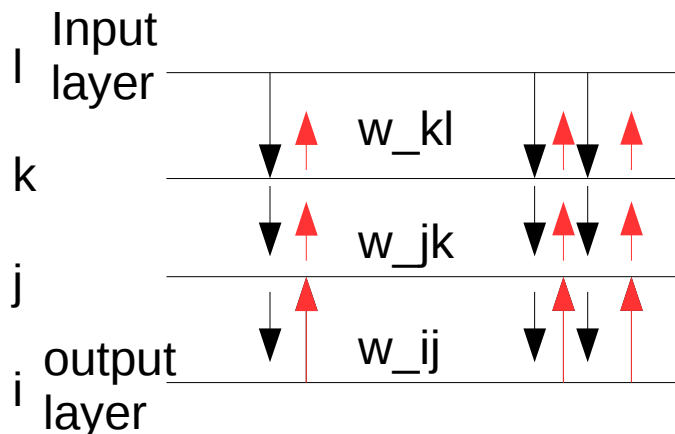


The state of the output functions for the layers i (output) and j input:

$$S_i = f(h_i)$$

$$h_i = \sum_j w_{ij} S_j - \theta_i$$

10. feed forward NN (4 layers)



Black for the input; red for the back prop training direction

Back Propagation (4)

11. chain rule, updating the weights (training)

For layer i

$$\frac{\partial D}{\partial w_{ij}} = \sum_n [y_i^n - f(h_i^n)] f'(h_i^n) \frac{\partial h_i^n}{\partial w_{ij}}$$

Note: the training described here all related to the derivative to the activation function, $f'(\cdot)$. So selection of the activation function is important and will be discussed in details next.

For layer j

$$\frac{\partial D}{\partial w_{jk}} = \sum_n \sum_i [y_i^n - f(h_i^n)] f'(h_i^n) \frac{\partial h_i^n}{\partial s_j} \cdot \frac{\partial s_j}{\partial w_{jk}}$$

For layer k

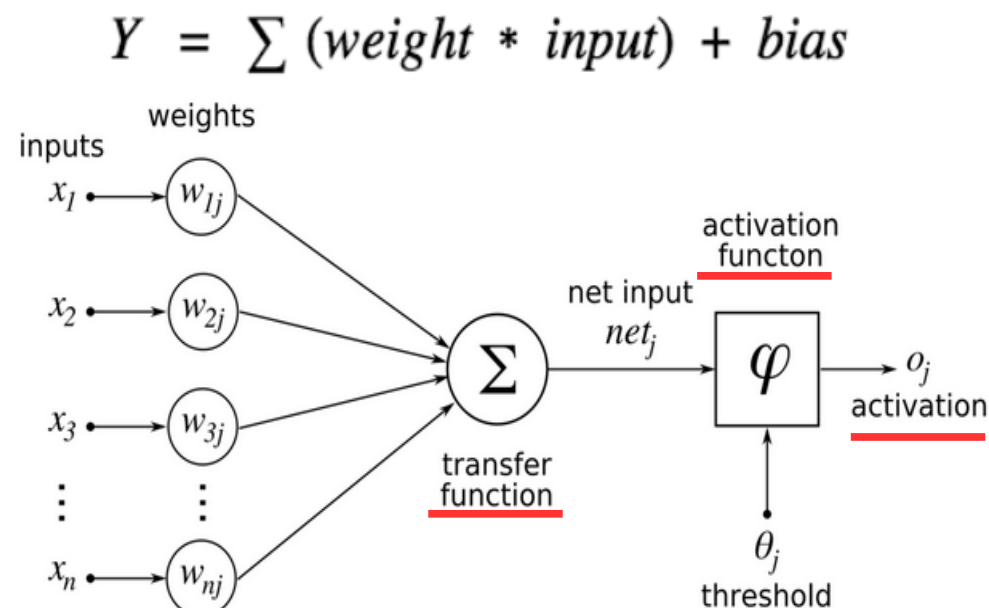
$$\frac{\partial D}{\partial w_{kl}} = \sum_n \sum_i \sum_j [y_i^n - f(h_i^n)] f'(h_i^n) \frac{\partial h_i^n}{\partial s_j} \cdot \frac{\partial s_j}{\partial s_k} \cdot \frac{\partial s_k}{\partial w_{kl}}$$

Activation Functions

Definition: for single neuron for the purpose of generating the output of the neuron.

Type: over hundreds proposed, we will focus On the following 4 types, Sigmoid, Softmax, Tanh, ReLU, (SSTR)

Characteristics and Comparison:



Sigmoid

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

Tanh

$$\tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}}$$

ReLU

$$f(x) = \max(0, x)$$

Softmax

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}, j = 1, 2, \dots, K$$

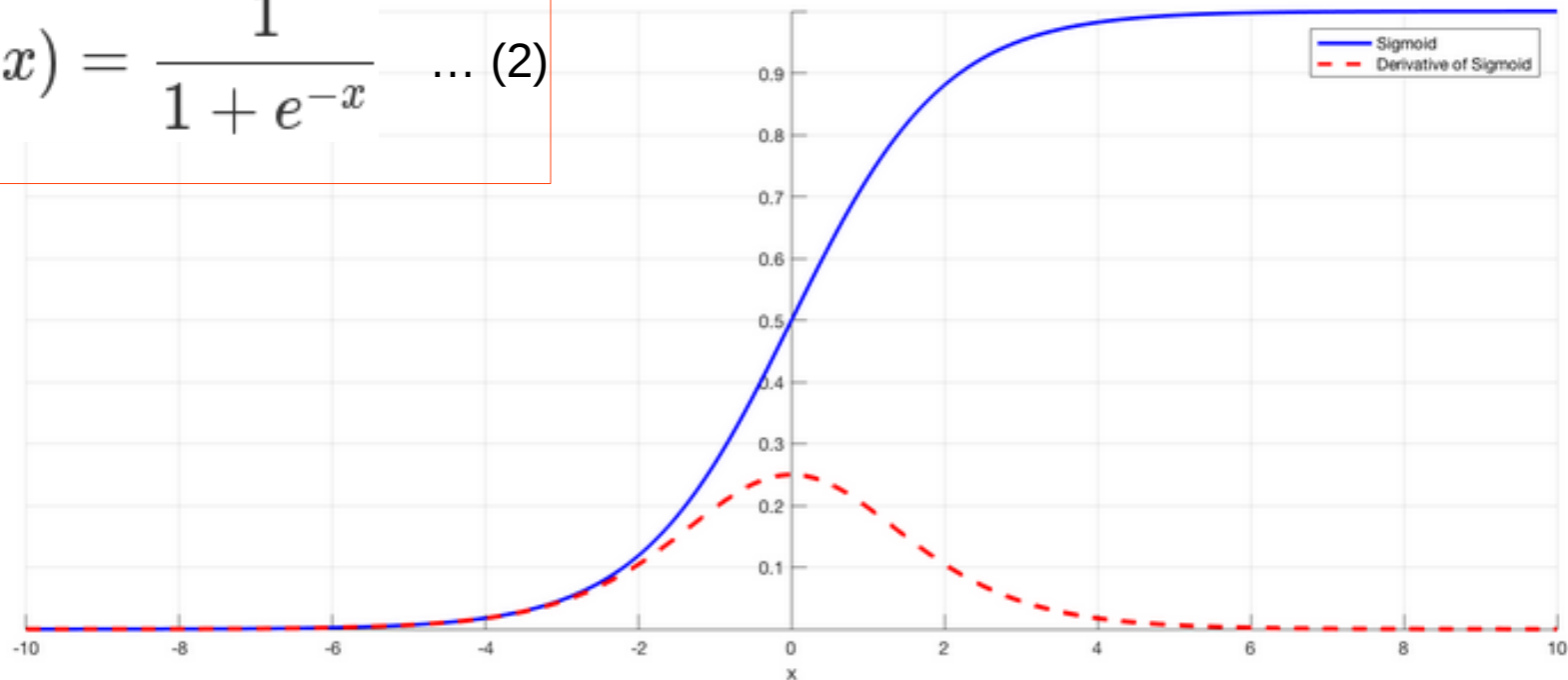
Sigmoid Functions

https://isaacchanghau.github.io/post/activation_functions/

$$\sigma(x) = \frac{L}{1 + e^{-k(x-x_0)}} \dots (1)$$

Logistic function in general as in equation (1)

$$\sigma(x) = \frac{1}{1 + e^{-x}} \dots (2)$$

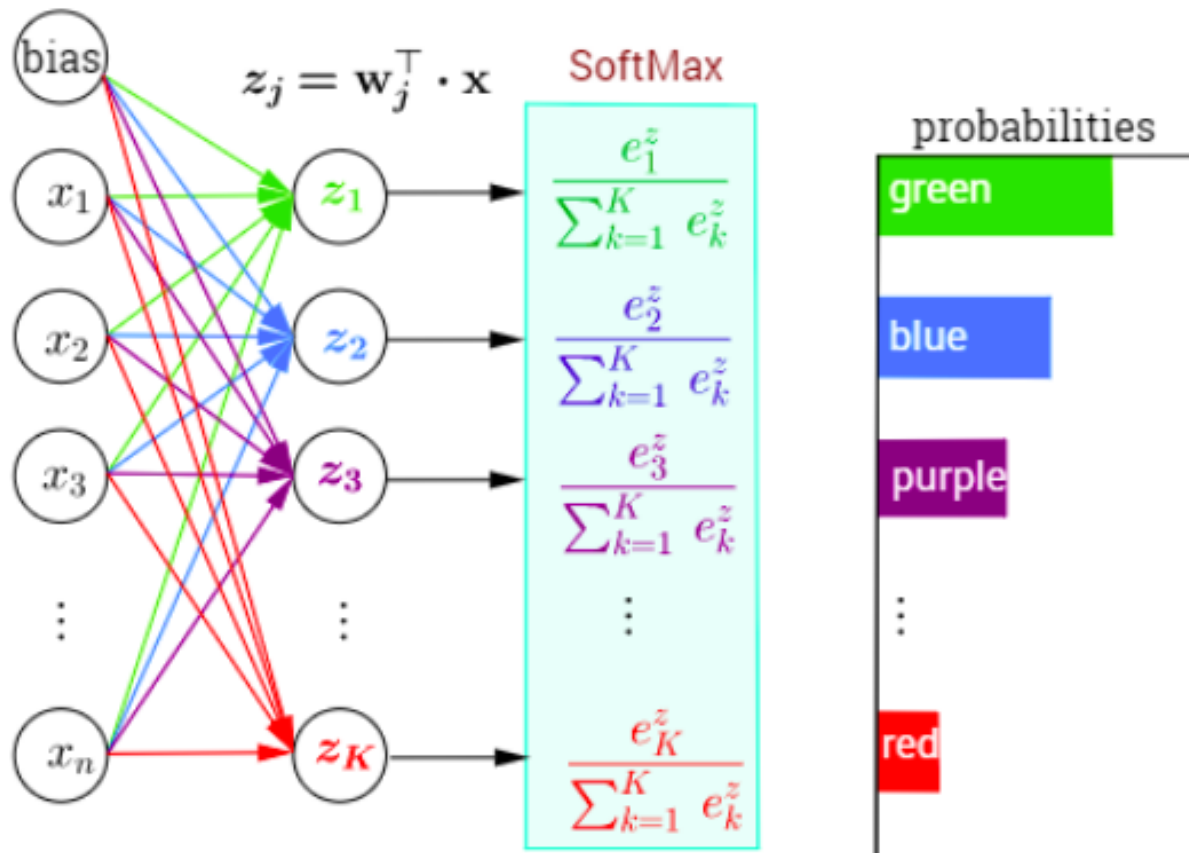


Softmax Functions

https://isaacchanghau.github.io/post/activation_functions/

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}, j = 1, 2, \dots, K \quad \dots (1)$$

used in various multiclass classification methods, such as multinomial logistic regression, multiclass linear discriminant analysis, naive Bayes classifiers, and artificial neural networks.



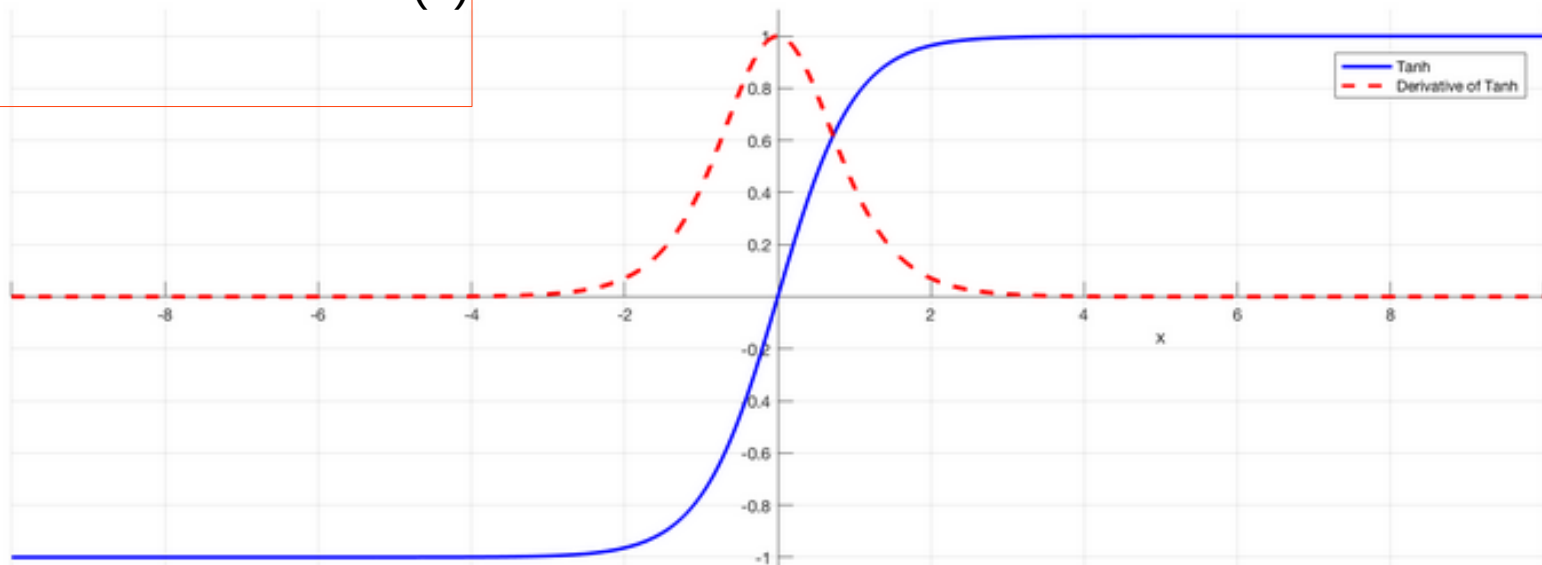
tanh Functions

https://isaacchanghau.github.io/post/activation_functions/

$$\tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}} \dots (1)$$

The derivative (red curve)

... (2)



Relu Functions

https://isaacchanghau.github.io/post/activation_functions/

$$f(x) = \max(0, x) \quad \dots (1)$$

One Relu example (green)

$$f(x) = \ln(1 + e^x) \quad \dots (2)$$

Its derivative:

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

