# Example on Back Propagation (1)

Example: Given training data below



$C_1$: $(1,2)$, $(1,3)$, $(2,3)$,

$C_2$: $(2,1)$, $(3,1)$, $(3,2)$

as shown in figure 1.

Region $C_1$, $d=1$

Region $C_2$, $d=0$

Suppose a single neuron network is given below



$y = f\left(\sum w_i x_i + \theta\right)$

Now based on the notation in the class, see Figure 1. we have



1. input and output neurons

Input

$$\vec{z} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}_{M \times 1}$$

Output

$$\vec{s} = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_{no} \end{pmatrix}_{N_p \times 1}$$

So we have input and output vector as follows

$$\vec{\sigma} = \begin{pmatrix} \sigma_1 \\ \sigma_2 \end{pmatrix} \quad \dots (1)$$

and the output vector

$$\vec{s} = (s_1) \quad \dots (2)$$

$$w_{i,k} \quad \dots (3)$$

where i = 1, 2 for inputs and k =1 for the output

Harry Li, Ph.D.

# Example on Back Propagation (2)

Example: Given training data below

$C_1: (1,2), (1,3), (2,3),$

$C_2: (2,1), (3,1), (3,2)$
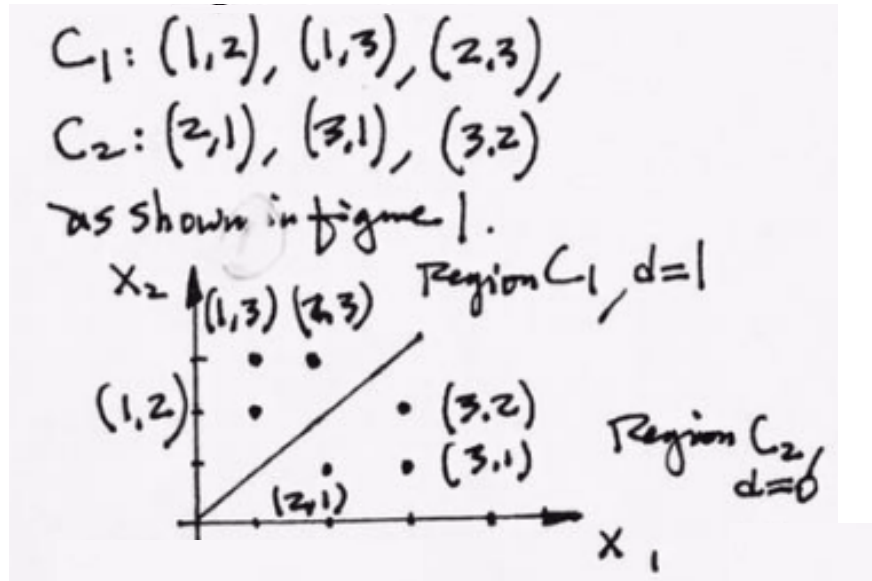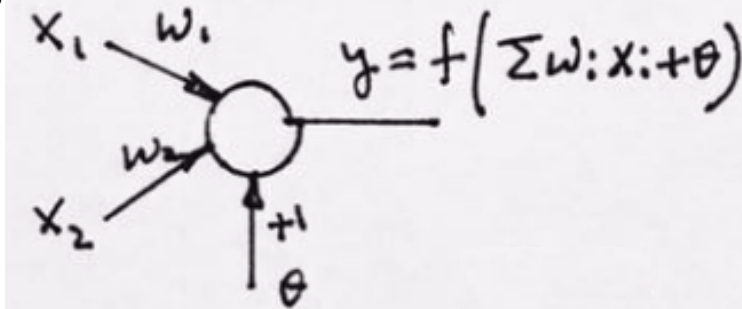
as shown in figure 1.



Suppose a single neuron network is given below



Now based on the notation in the class, see Figure 1. we have

Note: $(\sigma_1, \sigma_2) = (1, 2)$ for class $C_1$ for example.

Define activation function below

$$s_i = f(\sum_{i=1}^{2} w_{i,k} s_k + \phi) = f(\sum_{i=1}^{3} w_{i,k} s_k) \quad \dots (4)$$

Note the summation upper bound is set to 2 for explicit offset phi, and the upper bound is changed to 3 when changed to generalized form by making s_3 = phi, and its weigh $w_{i,k}$ for i = 3 and k=1. So the notation is in the unified summation form.

Define transfer function as

$$h_i = \sum_{i=1}^{2} w_{i,k} s_k + \phi = \sum_{i=1}^{3} w_{i,k} s_k \quad \dots (5)$$

Note subscript for h is index of the output i, in this example i = 1

# Example on Back Propagation (3)

Define transfer function, in case of single layer NN, we use $\sigma_k$ as input, and in case of multilayer NN, we use $s_k$ as notation for the input from the previous output layer.

$$h_i = \sum_{i=1}^{2} w_{i,k}\sigma_k + \phi = \sum_{i=1}^{3} w_{i,k}\sigma_k \quad ...(5\text{-}1)$$
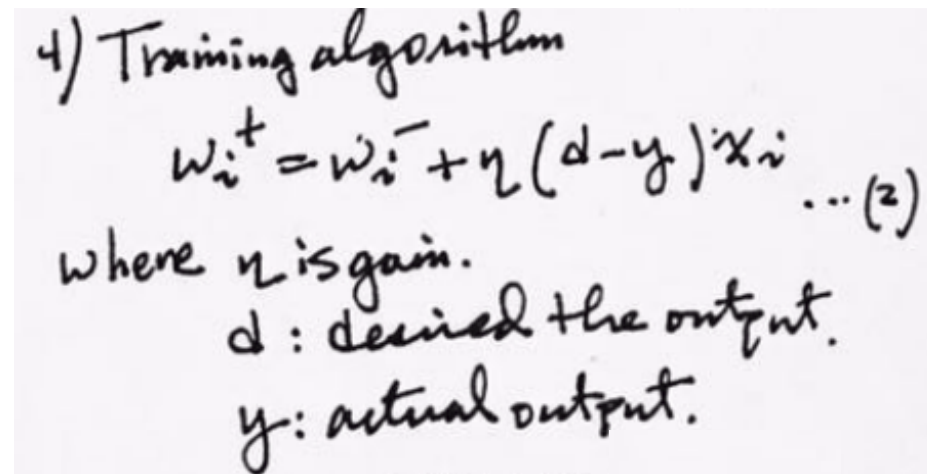
So the neuron output i is

$$s_i = f(h_i) = f(h_i(w_{i,k})) \quad ... (6)$$

4) Training algorithm
$$w_i^{+} = w_i^{-} + \eta (d-y)\, x_i \quad ...(2)$$
where $\eta$ is gain.
d : desired the output.
y : actual output.

Define an error as follows

$$\zeta_i{}^{\mu} - s_i{}^{\mu} \qquad ... (7)$$

The experiment $\mu$ in equation (7) refers to as the input data (1,2) from C1.

note where $\zeta_i{}^{\mu}$ is desired (correct) output i at experiment $\mu$, and $s_i{}^{\mu}$ is the actual output i at experiment $\mu$.

Let's take a look at the hand calculation example, from the reference copied on the right, where desired output is d and actual output is y, and since in this hand calculation example it is only have one output, so i = 1 in eqn(7), and

Harry Li, Ph.D.

# Example on Back Propagation (4)

5) Now, using data from $C_1$ and $C_2$ to perform training.

Choose $(1, 2)$ from $C_1$.

$$W_i^+ = W_i^- + \eta(1-y)x_1 \Big|_{x_1 = 1}$$

$$= 0.5 + \eta(1-y) \cdot 1$$

where

$$y = f\left(\sum_{i=1}^{2} W_i x_i + \theta\right)$$

$$= f\left(W_1 x_1 + W_2 x_2 + \theta\right)$$

$$= f\left(0.5 \times 1 + 0.5 \times 2 + 0.5\right) = f(2)$$

$$= sgn(2) = 1$$

hence, $W_i^+ = W_i^- + \eta(1-1)x_1 = W_i^- = 0.5$

See from hand calculation example,

The experiment $\mu_i$ in equation (7) refers to as the input data (1,2) from C1 as shown on the left.

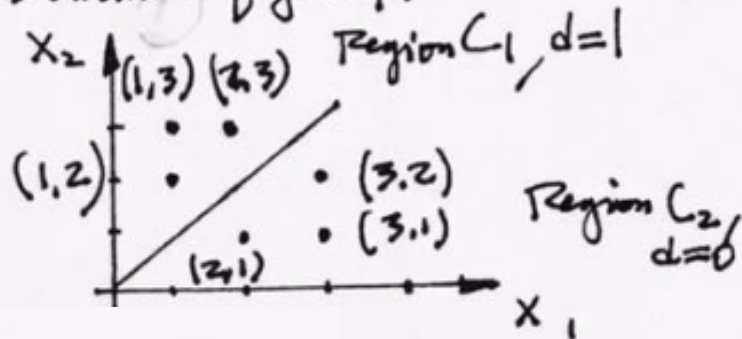Harry Li, Ph.D.

# Example on Back Propagation (5)

**Definition 6.** *Define total error for all neuron outputs and for all experiments*

$$D = \frac{1}{2} \sum_{\mu=1}^{6} \sum_{i=1}^{1} (\zeta_i^{\mu} - s_i^{\mu})^2 \tag{8}$$

note where $\zeta_i^{\mu}$ for i equal to 1, for only 1 output neuron, from the hand calculation example, and and experiment $\mu$ is equal to 6 due to total number of 6 feature points, e.g., 6 experiments.
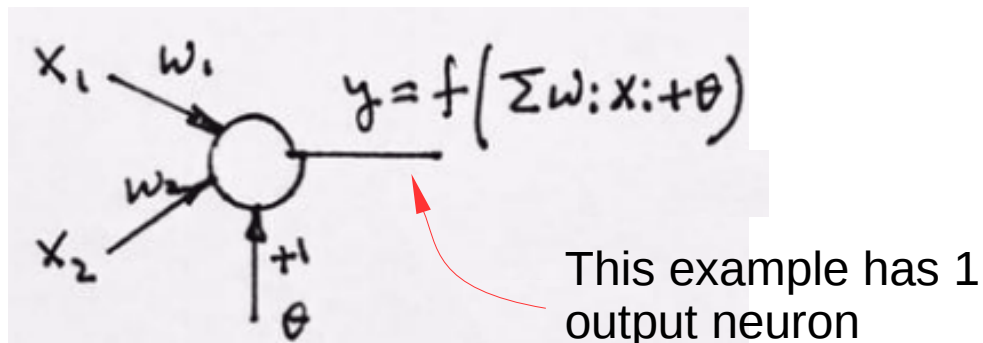
# Example on Back Propagation (6)

**Definition 7.** *Minimize error function*

$$\frac{\partial D}{\partial w_{i,k}} = \frac{\partial}{\partial w_{i,k}} \frac{1}{2} \sum_{\mu=1}^{6} \sum_{i=1}^{1} (\zeta_i^{\mu} - s_i^{\mu})^2 = \sum_{\mu=1}^{6} (\zeta_i^{\mu} - s_i^{\mu}) f'(h_i^{\mu}) \frac{\partial h_i}{\partial w_{i,k}} \quad \ldots (9)$$

So for the hand calculation example, we have 1 output, and input feature vector is 2 dimension,

$$\frac{\partial D}{\partial w_{1,k}} = \frac{\partial}{\partial w_{1,k}} \frac{1}{2} \sum_{\mu=1}^{6} \sum_{i=1}^{1} (\zeta_1^{\mu} - s_1^{\mu})^2 = \sum_{\mu=1}^{6} (\zeta_1^{\mu} - s_1^{\mu}) f'(h_1^{\mu}) \frac{\partial h_1}{\partial w_{1,k}} \quad \ldots (10)$$



This example has 1 output neuron

Harry Li, Ph.D.

# Example on Back Propagation (7)

Note the derivative of the activation function

$$f'(h_1{}^\mu) \qquad\qquad (11)$$
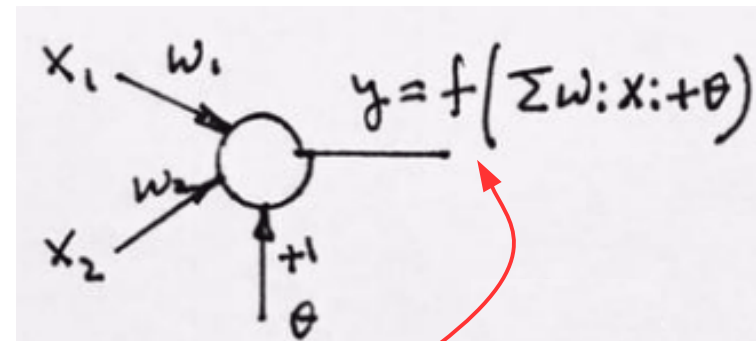
Suppose we choose RELU, then compute its derivative as

**One Relu example**

$$f(x) = \ln(1 + e^x) \qquad \dots (12)$$

Its derivative:

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

$$y = f\left(\Sigma w_i x_i + \theta\right)$$

Activation function
can be chosen as
RELU for example

Harry Li, Ph.D.

# Example on Back Propagation (8)

**Definition 8.** *Learning by updating the weights*

$$w_{i,k}(t+1) = w_{i,k}(t) + \delta w_{i,k}(t) \qquad (12)$$

where

$$\delta w_{i,k}(t) = -\epsilon \frac{\partial D}{\partial w_{i,k}} \qquad (13)$$

for the given example, we have

$$w_{1,k}(t+1) = w_{1,k}(t) + \delta w_{1,k}(t) \qquad (14)$$

where

$$\delta w_{1,k}(t) = -\epsilon \frac{\partial D}{\partial w_{1,k}} \qquad (15)$$

for k = 1, 2, see example feature vector dimension is 2

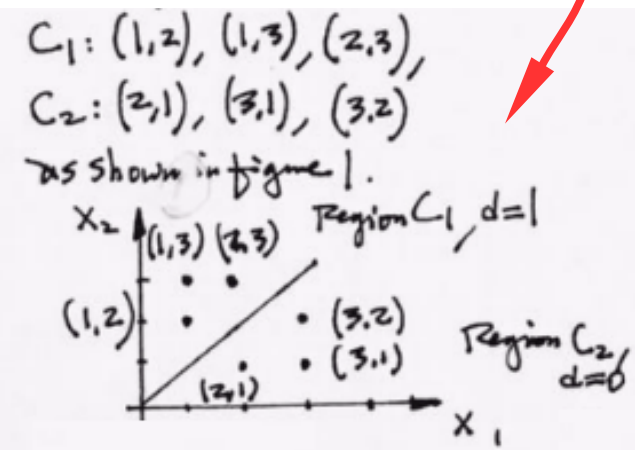Note equation (14) now is employed to replace hand calculation equation (2)

From hand calculation example

4) Training algorithm

$$w_i^+ = w_i^- + \eta(d-y)x_i \quad \cdots (2)$$

where $\eta$ is gain.

d : desired the output.

y : actual output.

$C_1$: $(1,2)$, $(1,3)$, $(2,3)$,
$C_2$: $(2,1)$, $(3,1)$, $(3,2)$
as shown in figure 1.

Region $C_1$, d=1
Region $C_2$, d=0

$(1,3)$ $(2,3)$
$(1,2)$
$(3,2)$
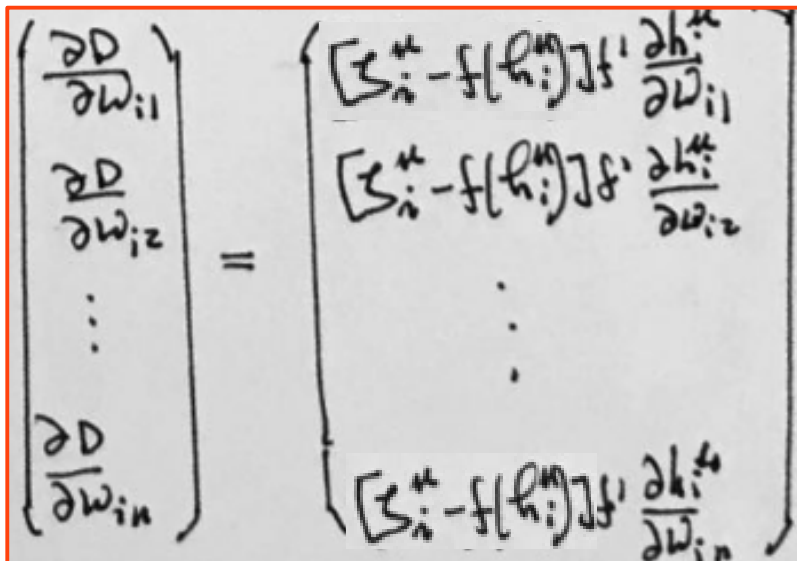$(2,1)$
$(3,1)$

# Example on Back Propagation (9)

So the gradient is defined as follows for the given example

$$\begin{pmatrix} \frac{\partial D}{\partial w_{i,1}} \\ \frac{\partial D}{\partial w_{i,2}} \end{pmatrix} = \begin{pmatrix} (\zeta_i{}^\mu - s_i{}^\mu) f'(h_i{}^\mu) \frac{\partial h_i}{\partial w_{i,1}} \\ (\zeta_i{}^\mu - s_i{}^\mu) f'(h_i{}^\mu) \frac{\partial h_i}{\partial w_{i,2}} \end{pmatrix} \qquad (16)$$

Since in our example, we have only 1 output, so

$$\begin{pmatrix} \frac{\partial D}{\partial w_{1,1}} \\ \frac{\partial D}{\partial w_{1,2}} \end{pmatrix} = \begin{pmatrix} (\zeta_1{}^\mu - s_1{}^\mu) f'(h_1{}^\mu) \frac{\partial h_1}{\partial w_{1,1}} \\ (\zeta_1{}^\mu - s_1{}^\mu) f'(h_1{}^\mu) \frac{\partial h_1}{\partial w_{1,2}} \end{pmatrix} \qquad (17)$$

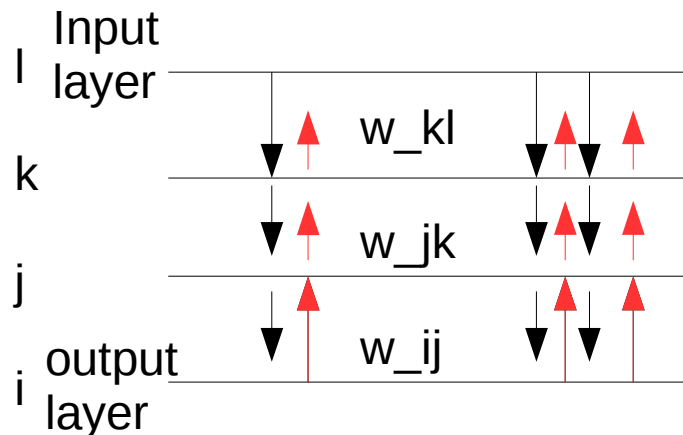Hence, the following gradient for n inputs is now replaced by equation (17)



Harry Li, Ph.D.
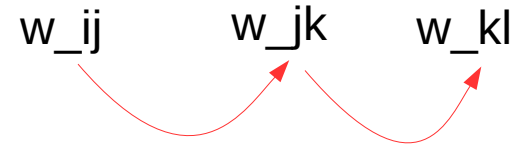
# Example on Back Propagation (10)

Now generalize the discussion, we can extend the mathematical formulation to multiple layers, for example in this lecture, we have 4 layers

10. feed forward NN (4 layers)

l  Input layer

k

j

i  output layer

w_kl

w_jk
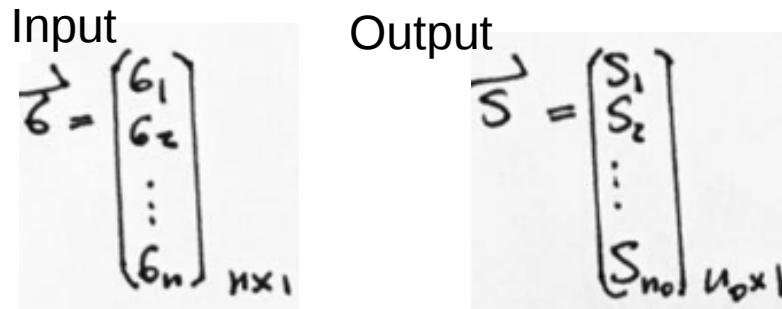
w_ij

w_ij        w_jk        w_kl

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

$$S_i = f(h_i)$$

$$h_i = \sum_{\partial} w_{ij} S_{\partial} - \theta_i$$
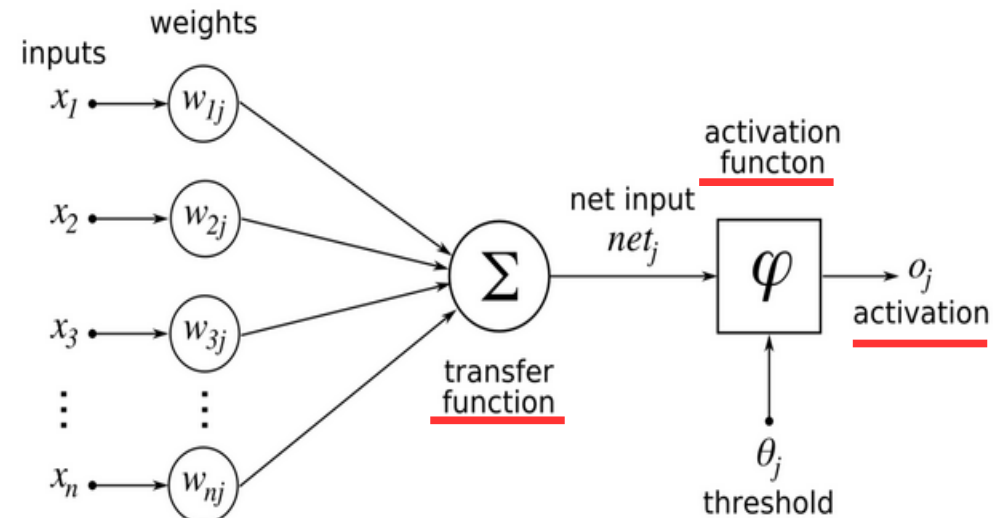
# Back Propagation (1)

1. input and output neurons

Input


Output


The architecture


Input: k
output: i

2. weights w_ik

i for output          k for input

Other popular notation

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



inputs
weights
$x_1 \quad w_{1j}$
$x_2 \quad w_{2j}$
$x_3 \quad w_{3j}$
$x_n \quad w_{nj}$

transfer function

net input $net_j$

$\Sigma$

activation functon

$\varphi$

$o_j$ activation

$\theta_j$ threshold

3. Activation function

$$S_i = f\left( \sum_k \omega_{ik} 6_k \right) \qquad \dots (1)$$

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

Harry Li, Ph.D.

# Back Propagation (2)

**4. Transfer function h_i**

let

$$h_i = \sum_k w_{ik} 6_k \qquad \dots (2)$$

Index i for the output neuron

Neuron output function S_i

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

Neuron output is tied to activation function f( . ), 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 Mu**

$$\zeta_i^\mu - S_i^\mu \qquad \dots (2\text{-}2)$$

**6. total error for all output neuron i and all experiments Mu**

$$D = \frac{1}{2} \sum_\mu \sum_i \left( \zeta_i^\mu - S_i^\mu \right)^2 \qquad \dots (3)$$

**7. minimize the error wrt to w_ik**

$$\frac{\partial D}{\partial w_{ik}} = \frac{\partial}{\partial w_{ik}} \cdot \frac{1}{2} \sum_\mu \sum_i \left[ \zeta_i^\mu - f(h_i^\mu) \right]^2$$

$$= \sum_\mu \left[ \zeta_i^\mu - f(h_i^\mu) \right] f'(h_i^\mu) \frac{\partial h_i}{\partial w_{ik}}$$

**8. Training the NN by updating the weights**

$$w_{ik}(t+1) = w_{ik}(*) + \delta w_{ij}(*) \qquad \dots (4)$$

where

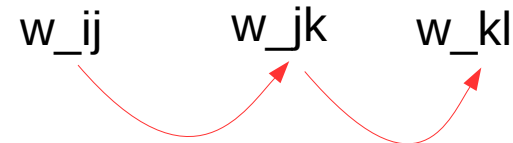$$\delta w_{ij}(t) = -\epsilon \frac{\partial D}{\partial w_{ij}} \qquad \dots (5)$$

Harry Li, Ph.D.

# Back Propagation (3)

9. computation of the derivatives to update the weight

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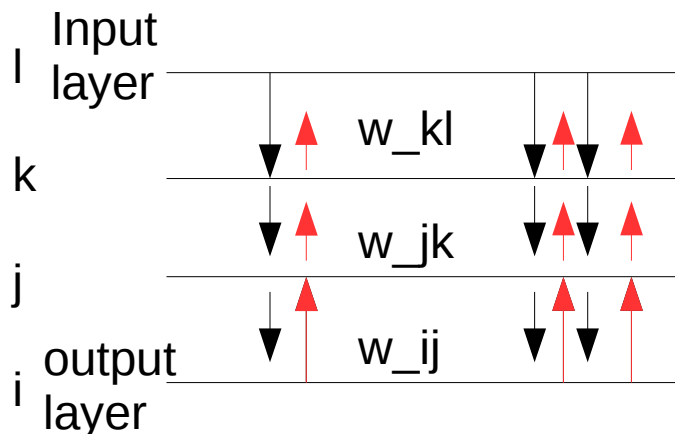$$

Update sequence back to the front:

w_ij     w_jk     w_kl

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

$$S_i = f(h_i)$$

$$h_i = \sum_{\partial} w_{ij} S_{\partial} - \theta_i$$

10. feed forward NN (4 layers)

Input
layer

l

w_kl

k

w_jk

j

output
w_ij
layer

i

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

Harry Li, Ph.D.

# Back Propagation (4)

11. chain rule, updating the weights (training)

For layer i

$$\frac{\partial D}{\partial w_{ij}} = \sum_{m} \left[ S_i^m - f(\ell_{n,i}^m) \right] f'(\ell_{n,i}^m) \frac{\partial \ell_{n,i}^m}{\partial w_{ij}}$$

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

For layer j

$$\frac{\partial D}{\partial w_{jk}} = \sum_{m} \sum_{i} \left[ S_i^m - f(\ell_{n,i}^m) \right] f'(\ell_{n,i}^m) \frac{\partial \ell_{n,i}^m}{\partial S_j} \cdot \frac{\partial S_j}{\partial w_{jk}}$$

For layer k

$$\frac{\partial D}{\partial w_{k\ell}} = \sum_{m} \sum_{i} \sum_{j} \left[ S_i^m - f(\ell_{n,i}^m) \right] f'(\ell_{n,i}^m) \frac{\partial \ell_{n,j}^m}{\partial S_j} \cdot \frac{\partial S_j}{\partial S_k} \cdot \frac{\partial S_k}{\partial w_{k\ell}}$$

Harry Li, Ph.D.

# 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:

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



inputs weights

$x_1 \bullet \longrightarrow w_{1j}$

$x_2 \bullet \longrightarrow w_{2j}$

$x_3 \bullet \longrightarrow w_{3j}$

$x_n \bullet \longrightarrow w_{nj}$

$\Sigma$

transfer function

net input $net_j$

activation functon

$\varphi$

$\longrightarrow o_j$

activation

$\theta_j$ threshold

Sigmoid

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

Tanh

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

Softmax

ReLU

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

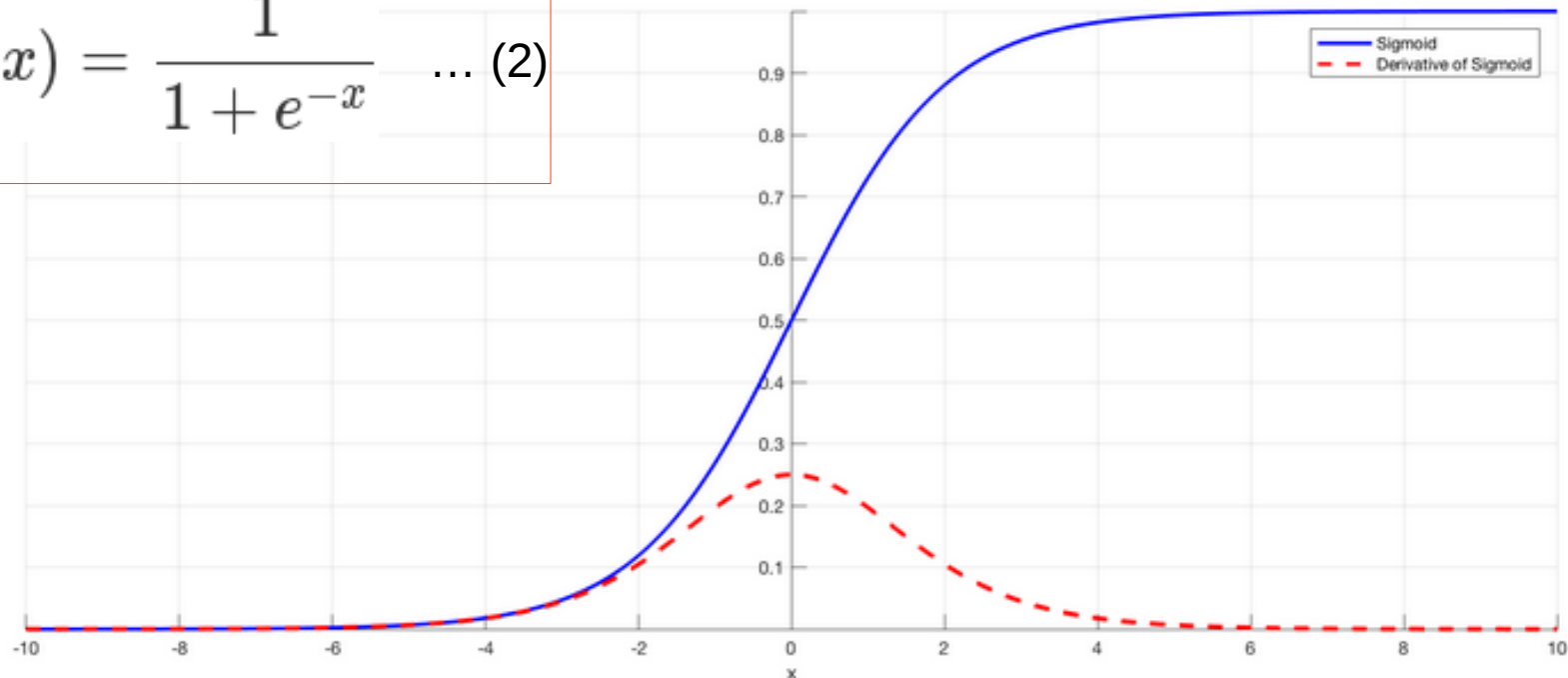$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}}, j = 1, 2, \ldots, K$$

# Sigmoid Functions

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

Logistic function in general as in equation (1)

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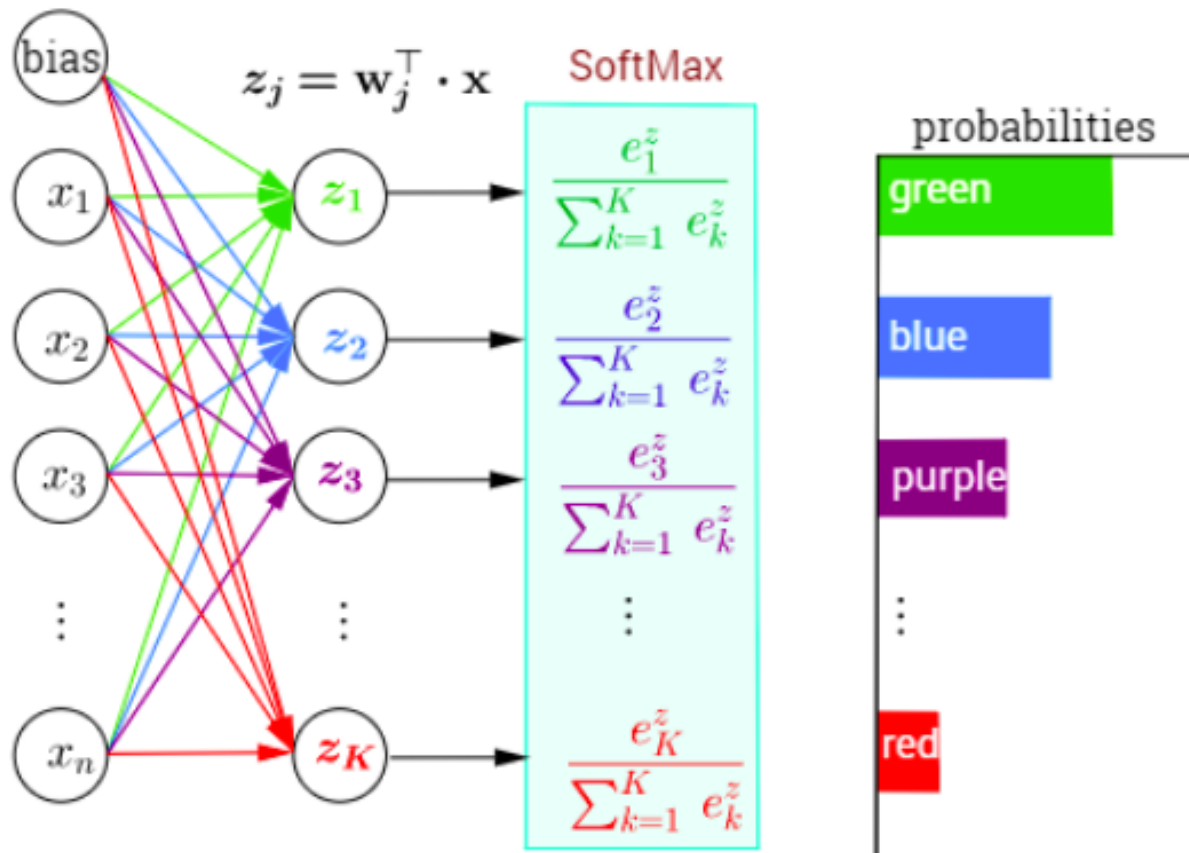
# Softmax Functions

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

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

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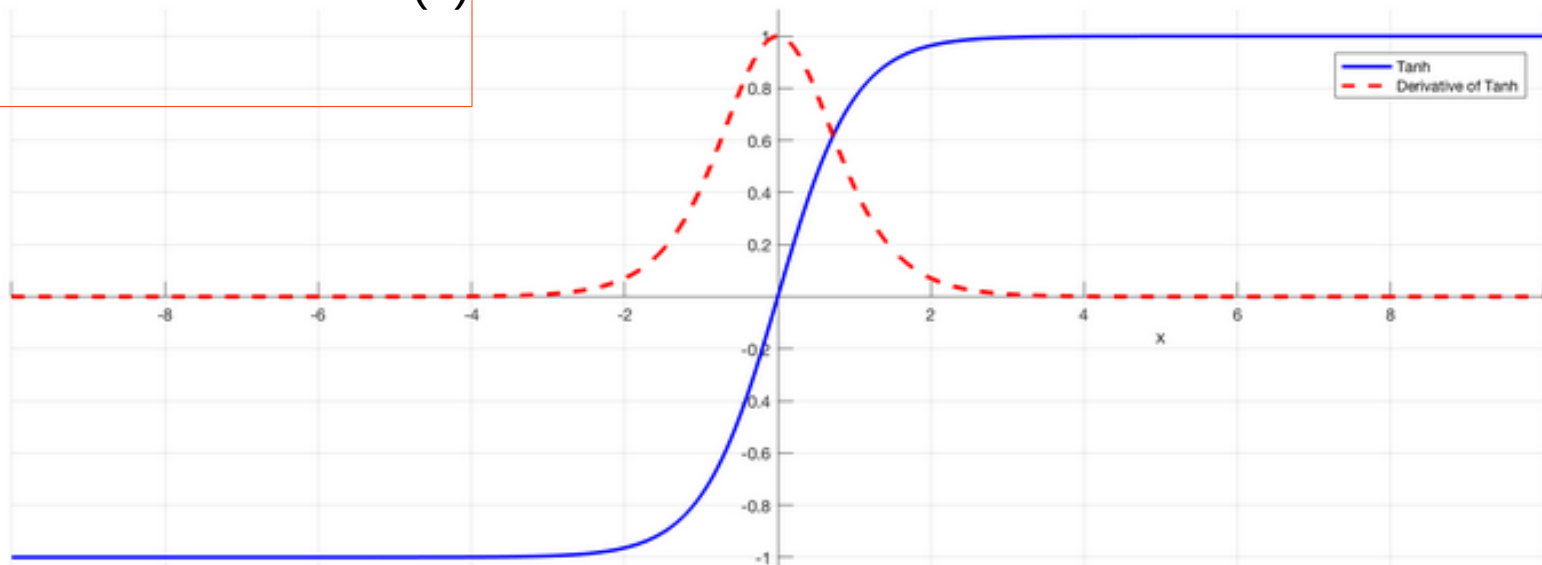
# tanh Functions

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

The derivative (red curve)

… (2)

# Relu Functions

$$f(x) = max(0, x)$$ ... (1)

One Relu example (green)

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

Its derivative:

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