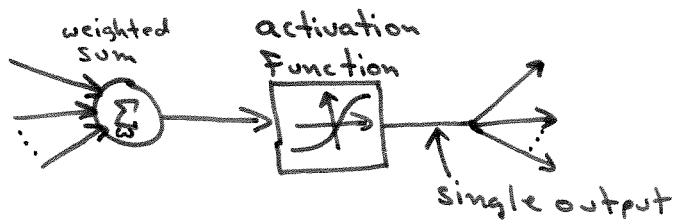
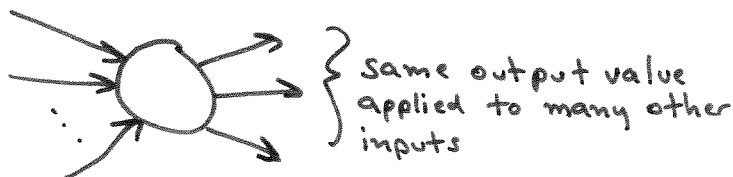


# Chapter 13      FeedForward Neural Networks

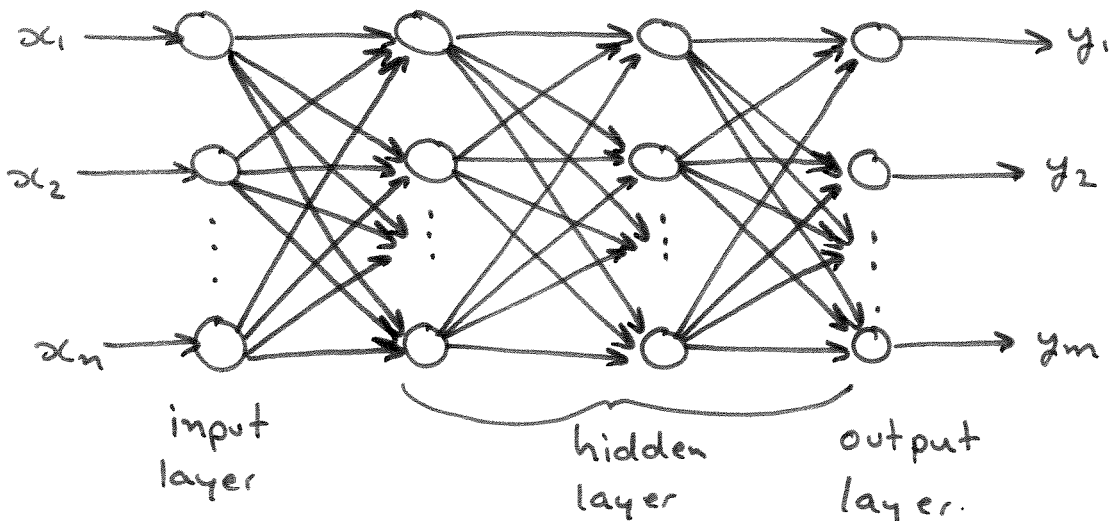
Single neuron model:



multi-input single-output system:



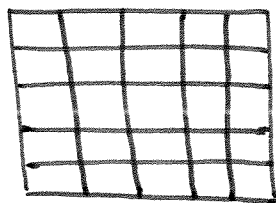
## FeedForward neural network



mapping:  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  (possibly nonlinear)

want to approximate  $F$  by adjusting the weights in the system.  $\Rightarrow$  "distributed representation" of  $F$

## Example: Pattern Recognition

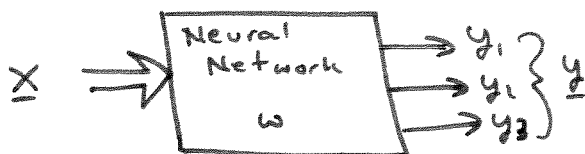


5x5 pixel image

let  $\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} \in \mathbb{R}^{25}$

$x_1, x_2, x_3, x_4, x_5$

$x_i \in \mathbb{R}^5 \quad i=1(1)5$



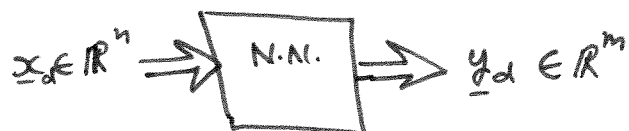
training set: "picture" dog  $\underline{x}_d$   $\underline{y}_d = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$   
 cat  $\underline{x}_c$   $\underline{y}_c = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$  } set weights to achieve this  
 bird  $\underline{x}_b$   $\underline{y}_b = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$

this is called supervised learning.

now input picture of fish  $\underline{x}_f \Rightarrow \boxed{\phantom{N.N.}} \rightarrow \underline{y}_f = \begin{pmatrix} .5 \\ 0 \\ .5 \end{pmatrix}$

interpretation??

in general, for a neural network:

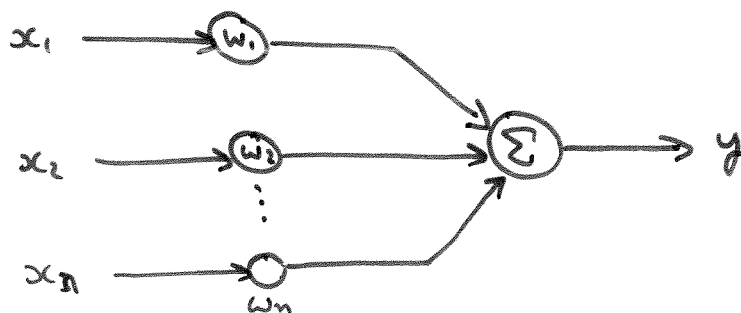


training set: p input/output pairs.

$(\underline{x}_{d,1}, \underline{y}_{d,1}), (\underline{x}_{d,2}, \underline{y}_{d,2}) \dots (\underline{x}_{d,p}, \underline{y}_{d,p}) \in \mathbb{R}^n \times \mathbb{R}^m$

where:  $\underline{y}_{d,i} = F(\underline{x}_{d,i})$  "desired" output

## Single - Neuron Training



single linear  
neuron  
 $\mathbb{R}^n \rightarrow \mathbb{R}$

$$y = \sum_{i=1}^n w_i x_i = \underline{x}^T \underline{w}$$

input  
 $\underline{x} \in \mathbb{R}^n$   $\underline{w} \in \mathbb{R}^n$  weights  
 $y \in \mathbb{R}$  output

given a map  $F: \mathbb{R}^n \rightarrow \mathbb{R}$  Find weights  $\underline{w} \ni$   
neuron approximates  $F$  as  
closely as possible.

use a training set of  $p$  pairs:  $\{(\underline{x}_{d,1}, y_{d,1}) \dots (\underline{x}_{d,p}, y_{d,p})\}$   
 $y_{d,i} = F(\underline{x}_{d,i}) \quad i=1(1)p$

we want to optimize  $\underline{w}$ :

$$\min_{\underline{w}} \frac{1}{2} \sum_{i=1}^p (y_{d,i} - \underline{x}_{d,i}^T \underline{w})^2$$

or in matrix form, let  $X_d = [\underline{x}_{d,1} \dots \underline{x}_{d,p}] \in \mathbb{R}^{n \times p}$

$$\underline{y}_d = \begin{bmatrix} y_{d,1} \\ \vdots \\ y_{d,p} \end{bmatrix}$$

$$\min_{\underline{w}} \frac{1}{2} \|\underline{y}_d - X_d^T \underline{w}\|^2$$

Case 1:  $p < n$

assuming that the  $p$  inputs of the training set are independent  $\Rightarrow \text{rank } X_d^T = p$

we're actually trying to solve  $\underline{y}_d = X_d^T \underline{w}$   $X_d^T \in \mathbb{R}^{p \times n}$

if  $p < n \Rightarrow$  fewer training pairs,  $p$ , than the number of weights,  $n$ .

$\Rightarrow$  infinitely many possible solutions for weights,  $\underline{w}$  i.e.  $\|\underline{y}_d - X_d^T \underline{w}\| = 0$

the solution with the minimum norm  $\|\underline{w}\|$  is given by  $(A^T(AA^T)^{-1})$

$$\underline{w}^* = X_d (X_d^T X_d)^{-1} \underline{y}_d$$

by Kaczmarz's algorithm

$$\left\{ \begin{array}{l} \underline{w}^{(k+1)} = \underline{w}^{(k)} + \mu \frac{e_k \underline{x}_{d,R(k)}}{\|\underline{x}_{d,R(k)}\|^2} \quad \underline{w}^{(0)} = \underline{0} \\ e_k = y_{d,R(k)} - \underline{x}_{d,R(k)}^T \underline{w}^{(k)} \end{array} \right.$$

$R(k)$  is remainder if we divide  $k$  by  $p$ , +1  
 $R(k) \in \{1, \dots, p\}$

case 2:  $p > n \Rightarrow$  more training points than number of weights.

assume rank  $X_d^T = n$

$$\min_{\underline{w}} \frac{1}{2} \|\underline{y}_d - X_d^T \underline{w}\|^2$$

in general  $X_d^T \underline{w} \neq \underline{y}_d$   
unless  $\underline{y}_d \in \mathcal{R}(X_d^T)$

least squares solution:

$$\underline{w}^* = (X_d X_d^T)^{-1} X_d \underline{y}_d$$

we can also use any unconstrained optimization algorithm on

$$\begin{aligned} \|\underline{y}_d - X_d^T \underline{w}\|^2 &= (\underline{y}_d - X_d^T \underline{w})^T (\underline{y}_d - X_d^T \underline{w}) \\ &= \|\underline{y}_d\|^2 + \underline{w}^T X_d X_d^T \underline{w} - 2 \underline{w}^T X_d \underline{y}_d \end{aligned}$$

$$F(\underline{w}) = \frac{1}{2} \|\underline{y}_d - X_d^T \underline{w}\|^2 = \frac{1}{2} \underline{w}^T X_d X_d^T \underline{w} - \underline{w}^T X_d \underline{y}_d + \|\underline{y}_d\|^2$$

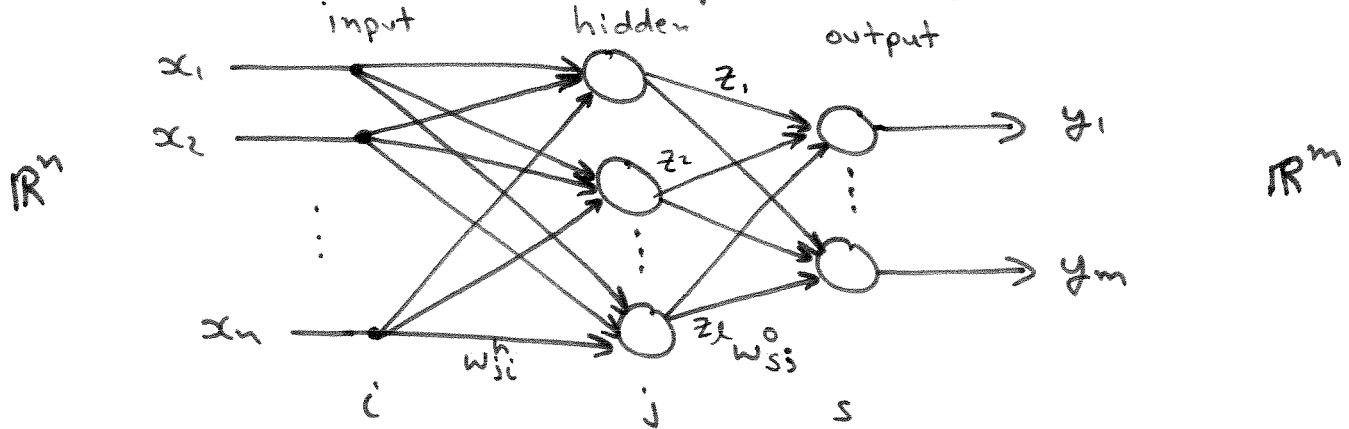
$$\nabla F(\underline{w}) = X_d X_d^T \underline{w} - X_d \underline{y}_d$$

descent direction  $\underline{d} = -\nabla F = X_d \underline{y}_d - X_d X_d^T \underline{w}$

descent algorithm:  $\underline{w}^{(k+1)} = \underline{w}^{(k)} + \alpha_k \underline{d}$

# Backpropagation Algorithm

consider the following three layer network:



$n$  inputs  $x_i$   $i=1(1)n$   $l$  neurons in hidden layer.  
 $m$  outputs  $y_s$   $s=1(1)m$   
 $l$  outputs from hidden layer  $z_j$   $j=1(1)l$

input layer  $\Rightarrow$  identity map for activation function

hidden layer  $\Rightarrow$  activation function  $F_j^h$   $j=1(1)l$

output layer  $\Rightarrow$  activation function  $F_s^o$   $s=1(1)m$

$$F: \mathbb{R} \rightarrow \mathbb{R}$$

weights for inputs into hidden layer:  $w_{ji}^h$   $i=1(1)n$   
 $j=1(1)l$

weights for inputs into output layer:  $w_{ss}^o$   $s=1(1)m$

input to a hidden layer neuron:  $\underline{x} \Rightarrow \text{node } w_{ji}^h \rightarrow \nu_j^h$

$$\nu_j^h = \underline{w}_{ji}^{hT} \underline{x} = \sum_{i=1}^n w_{ji}^h x_i$$

output from hidden layer:

$$z_j = F_j^h \left( \sum_{i=1}^n w_{ji}^h x_i \right)$$

output from  $s^{\text{th}}$  neuron of the output layer:

$$y_s = F_s^o \left( \sum_{j=1}^l w_{sj}^o z_j \right)$$

$$\begin{aligned} \therefore y_s &= F_s^o \left( \sum_{j=1}^l w_{sj}^o F_j^h(v_j) \right) \\ &= F_s^o \left( \sum_{j=1}^l w_{sj}^o F_j^h \left( \sum_{i=1}^n w_{ji}^h x_i \right) \right) \\ &= F_s(x_1, \dots, x_n) \end{aligned}$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} F_1(\underline{x}) \\ \vdots \\ F_m(\underline{x}) \end{bmatrix}$$

now how do we find the sets of weights  $w_{ji}^h$  and  $w_{sj}^o$  given a single pair of input/output vectors  $(\underline{x}_d, \underline{y}_d)$   $\underline{x}_d \in \mathbb{R}^n$   $\underline{y}_d \in \mathbb{R}^m$  (desired output)

$$\min_{\underline{w}^o, \underline{w}^h} \frac{1}{2} \sum_{s=1}^m (y_{ds} - y_s)^2 \quad y_s = F_s(\underline{x})$$

$$\begin{aligned} E(\underline{w}^o, \underline{w}^h) &= \frac{1}{2} \sum_{s=1}^m (y_{ds} - y_s)^2 \\ &= \frac{1}{2} \sum_{s=1}^m \left( y_{ds} - F_s^o \left( \sum_{j=1}^l w_{sj}^o F_j^h \left( \sum_{i=1}^n w_{ji}^h x_{di} \right) \right) \right)^2 \end{aligned}$$

We minimize  $E(\underline{w}^0, \underline{w}^h)$  using a gradient algorithm with fixed step size, i.e. no line-minimization.

$\therefore$  need gradient of  $E$  with respect to  $\underline{w}^0, \underline{w}^h$ ,

write: 
$$E(\underline{w}) = \frac{1}{2} \sum_{p=1}^m \left( y_{dp} - F_p^0 \left( \sum_{q=1}^l w_{pq}^0 z_q \right) \right)^2$$

where 
$$z_q = F_q^h \left( \sum_{i=1}^n w_{qi}^h x_{di} \right)$$

$$\frac{\partial E(\underline{w})}{\partial w_{sj}^0} = - (y_{ds} - y_s) F_s^{0'} \left( \sum_{q=1}^l w_{sq}^0 z_q \right) z_j$$

$F_s^{0'}: \mathbb{R} \rightarrow \mathbb{R}$  is derivative of  $F_s^0$ .

now define: 
$$\delta_s = (y_{ds} - y_s) F_s^{0'} \left( \sum_{q=1}^l w_{sq}^0 z_q \right)$$

error:  $(y_{ds} - y_s)$  scaled by  $F_s^{0'} \left( \sum_{q=1}^l w_{sq}^0 z_q \right)$

$$\frac{\partial E(\underline{w})}{\partial w_{sj}^0} = -\delta_s z_j$$

in order to take  $\partial E / \partial w_{ji}^h$  we re-write  $E(\underline{w})$

$$E(\underline{w}) = \frac{1}{2} \sum_{p=1}^m \left( y_{dp} - F_p^0 \left( \sum_{q=1}^l w_{pq}^0 F_q^h \left( \sum_{r=1}^n w_{qr}^h x_{dr} \right) \right) \right)^2$$



$$\frac{\partial E}{\partial w_{ji}^h}(\underline{w}) = - \sum_{p=1}^m (y_{d,p} - y_p) f_p^{o'} \left( \sum_{q=1}^l w_{pq}^o z_q \right) w_{pj}^o f_j^{h'} \left( \sum_{r=1}^n w_{jr}^h x_{dr} \right)$$

•  $x_{di}$

again  $f_j^{h'} : \mathbb{R} \rightarrow \mathbb{R}$  is derivative of  $f_j^h$

we can write the above as:

$$\frac{\partial E}{\partial w_{ji}^h}(\underline{w}) = - \left( \sum_{p=1}^m \delta_p w_{pj}^o \right) f_j^{h'}(v_j) x_{di}$$

the fixed-step size gradient algorithm for each weight will have the form

$$w^{(n+1)} = w^{(n)} - \eta \frac{\partial E}{\partial w} \quad \eta \text{ is the step-size.}$$

$\therefore$  we have what is called the "backpropagation algorithm"

$$w_{sj}^{o(k+1)} = w_{sj}^{o(k)} + \eta \delta_s^{(k)} z_j^{(k)}$$

$$w_{ji}^{h(k+1)} = w_{ji}^{h(k)} + \eta \left( \sum_{p=1}^m \delta_p^{(k)} w_{pj}^o \right) f_j^{h'}(v_j^{(k)}) x_{di}$$

$$v_j^{(k)} = \sum_{i=1}^n w_{ji}^h x_{di}^{(k)}$$

$$z_j^{(k)} = f_j^h(v_j^{(k)})$$

$$y_s^{(k)} = f_s^o \left( \sum_{q=1}^l w_{sq}^o z_q^{(k)} \right)$$

$$\delta_s^{(k)} = (y_{ds} - y_s^{(k)}) f_s^{o'} \left( \sum_{q=1}^l w_{sq}^o z_q^{(k)} \right)$$

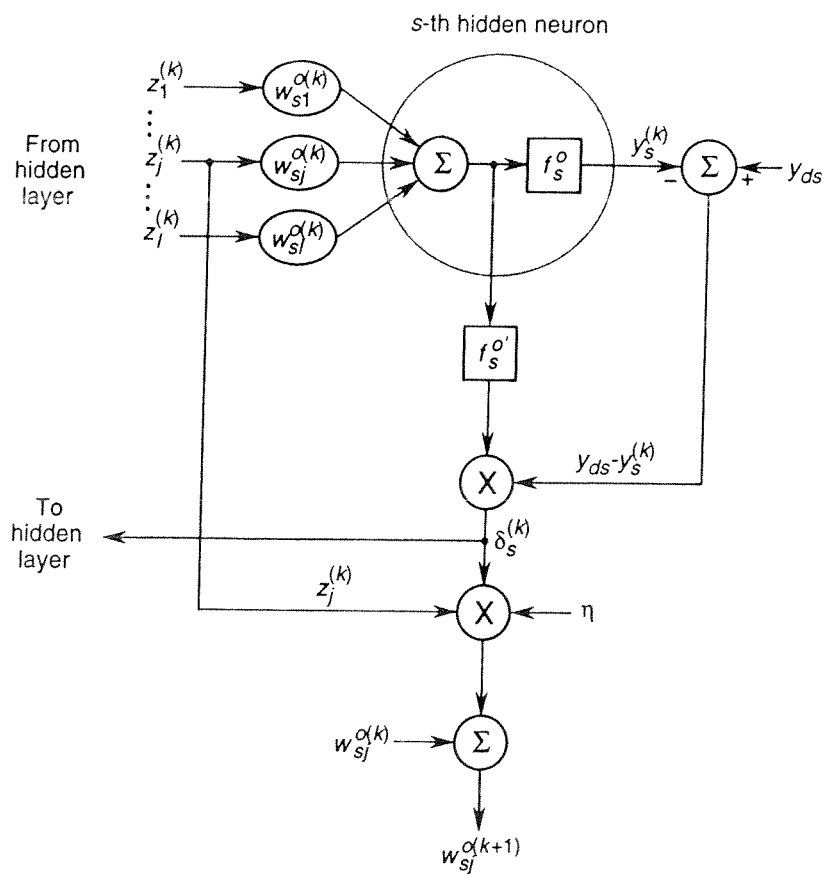


Figure 13.7. Illustration of the update equation for the output layer weights

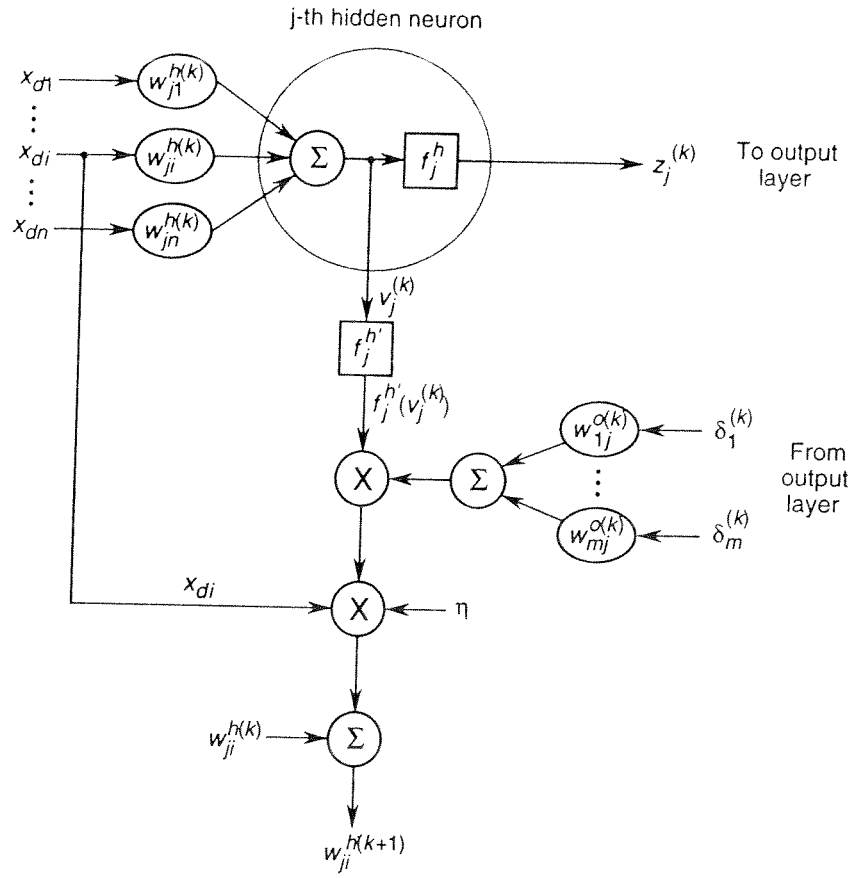
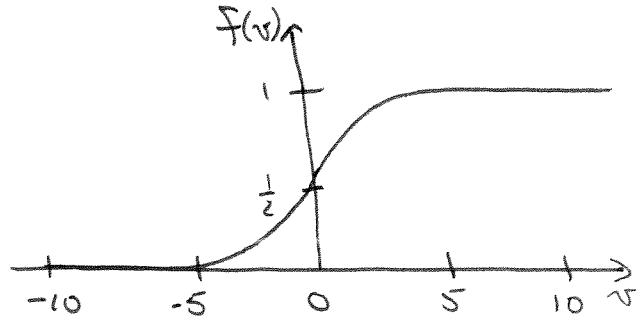


Figure 13.8. Illustration of the update equation for the hidden layer weights

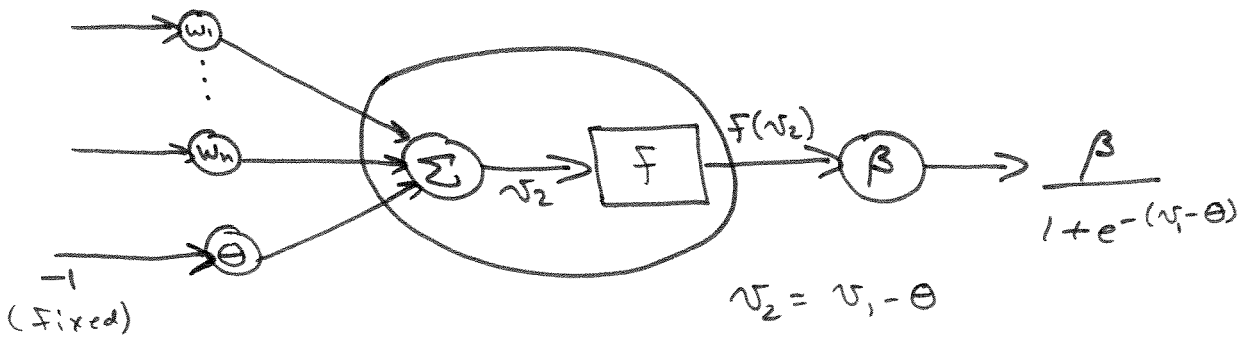
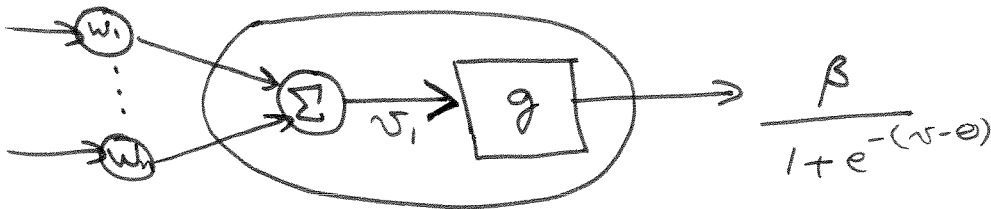
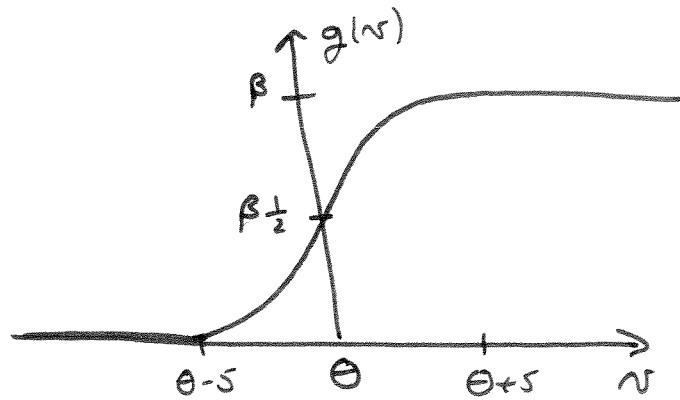
# Sigmoid Function

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



$$g(v) = \frac{\beta}{1 + e^{-(v-\theta)}}$$

$\beta, \theta$  are parameters  
 $\beta$ , scale parameter  
 $\theta$ , shift parameter



Note:  $f'(v) = \frac{-1(-e^{-v})}{(1+e^{-v})^2} = \frac{e^{-v}}{(1+e^{-v})^2} = \frac{1+e^{-v}-1}{(1+e^{-v})^2}$

$$= \frac{1}{(1+e^{-v})} \left( \frac{1+e^{-v}}{1+e^{-v}} - \frac{1}{(1+e^{-v})} \right) = f(v) (1 - f(v))$$

with the Sigmoid Function as the activation Function,

$$y_s = \beta_s^o f_s^o \left( \sum_{j=1}^l w_{sj}^o \beta_j^h f_j^h \left( \sum_{i=1}^n w_{ji}^h x_i - \theta_j^h \right) - \theta_s^o \right)$$

$$E(\underline{w}, \underline{\theta}, \underline{\beta}) = \frac{1}{2} \sum_{s=1}^m (y_{ds} - y_s)^2$$

$$= \frac{1}{2} \sum_{s=1}^m \left( y_{ds} - \beta_s^o f_s^o \left( \sum_{j=1}^l w_{sj}^o \beta_j^h f_j^h \left( \sum_{i=1}^n w_{ji}^h x_i - \theta_j^h \right) - \theta_s^o \right) \right)^2$$

input to  
j<sup>th</sup> hidden neuron

output from  
j<sup>th</sup> hidden neuron

output from s<sup>th</sup>  
output neuron

again letting  $z_q = f_q^h \left( \sum_{i=1}^n w_{qi}^h x_{di} - \theta_q^h \right)$

output from  
hidden neuron  
before multiplication  
with  $\beta_q^h$

$$E = \frac{1}{2} \sum_{p=1}^m \left( y_{dp} - \beta_p^o f_p^o \left( \sum_{q=1}^l w_{pq}^o \beta_q^h z_q - \theta_p^o \right) \right)^2$$

$$\frac{\partial E}{\partial w_{sj}^o} = - (y_{ds} - y_s) \beta_s^o f_s^o \left( \sum_{q=1}^l w_{sq}^o \beta_q^h z_q - \theta_s^o \right) \beta_j^h z_j$$

let  $\delta_s = (y_{ds} - y_s) f_s^o \left( \sum_{q=1}^l w_{sq}^o \beta_q^h z_q - \theta_s^o \right) \beta_s^o$

$$\therefore \frac{\partial E}{\partial w_{sj}^o} = - \delta_s \beta_j^h z_j$$

$$\frac{\partial E}{\partial \theta_s^0} = + (y_{ds} - y_s) \beta_s^0 F_s^0' \left( \sum_{q=1}^l w_{sq}^0 \beta_q^h z_q - \theta_s^0 \right) = \delta_s$$

$$\frac{\partial E}{\partial \beta_s^0} = - (y_{ds} - y_s) F_s^0 \left( \sum_{q=1}^l w_{sq}^0 \beta_q^h z_q - \theta_s^0 \right)$$

now we need  $\frac{\partial E}{\partial w_{ji}^h}$ ,  $\frac{\partial E}{\partial \theta_j^h}$ ,  $\frac{\partial E}{\partial \beta_j^h}$

we re-write  $E(\underline{w}, \underline{\theta}, \underline{\beta})$  as:

$$E = \frac{1}{2} \sum_{p=1}^m \left( y_{dp} - \beta_p^0 F_p^0 \left( \sum_{q=1}^l w_{pq}^0 \beta_q^h F_q^h \left( \sum_{r=1}^n w_{qr}^h x_{dr} - \theta_q^h \right) - \theta_p^0 \right) \right)^2$$

$$\frac{\partial E}{\partial w_{ji}^h} = - \sum_{p=1}^m (y_{dp} - y_p) \beta_p^0 F_p^0' \left( \sum_{q=1}^l w_{pq}^0 \beta_q^h z_q - \theta_p^0 \right) w_{pj}^0 \beta_j^h F_j^h' \left( \sum_{r=1}^n w_{jr}^h x_{dr} - \theta_j^h \right) \cdot x_{di}$$

$$= - \left( \sum_{p=1}^m \delta_p w_{pj}^0 \right) \beta_j^h F_j^h'(\nu_j) x_{di}$$

where  $\nu_j = \sum_{i=1}^n w_{ji}^h x_{di} - \theta_j^h$

$$\frac{\partial E}{\partial \theta_j^h} = + \sum_{p=1}^m (y_{dp} - y_p) \beta_p^0 F_p^0' \left( \sum_{q=1}^l w_{pq}^0 \beta_q^h z_q - \theta_p^0 \right) w_{pj}^0 \beta_j^h F_j^h' \left( \sum_{r=1}^n w_{jr}^h x_{dr} - \theta_j^h \right)$$

$$= \left( \sum_{p=1}^m \delta_p w_{pj}^0 \right) \beta_j^h F_j^h'(\nu_j)$$

$$\frac{\partial E}{\partial \beta_j^h} = - \sum_{p=1}^m (y_{dp} - y_p) \beta_p^o F_p^o \left( \sum_{q=1}^l w_{pq}^o \beta_q^h z_q - \Theta_p^o \right) w_{pj}^o F_j^h \left( \sum_{r=1}^m w_{jr} x_{dr} - \Theta_j^h \right)$$

$$= - \left( \sum_{p=1}^m \delta_p w_{pj}^o \right) F_j^h(v_j)$$

therefore the backpropagation algorithm becomes

$$w_{sj}^o(k+1) = w_{sj}^o(k) + \eta \delta_s^{(k)} \beta_j^h(k) z_j^{(k)}$$

$$\Theta_s^o(k+1) = \Theta_s^o(k) - \eta \delta_s^{(k)}$$

$$\beta_s^o(k+1) = \beta_s^o(k) + \eta (y_{ds} - y_s^{(k)}) y_s^{(k)} / \beta_s^o(k)$$

$$w_{ji}^h(k+1) = w_{ji}^h(k) + \eta \left( \sum_{p=1}^m \delta_p^{(k)} w_{pj}^o(k) \right) \beta_j^h(k) F_j^{h'}(v_j^{(k)}) x_{di}$$

$$\Theta_j^h(k+1) = \Theta_j^h(k) - \eta \left( \sum_{p=1}^m \delta_p^{(k)} w_{pj}^o(k) \right) \beta_j^h(k) F_j^{h'}(v_j^{(k)})$$

$$\beta_j^h(k+1) = \beta_j^h(k) + \eta \left( \sum_{p=1}^m \delta_p^{(k)} w_{pj}^o(k) \right) z_j^{(k)}$$

$$v_j^{(k)} = \sum_{i=1}^n w_{ji}^h(k) x_{di} - \Theta_j^h(k)$$

$$z_j^{(k)} = F_j^h(v_j^{(k)})$$

$$y_s^{(k)} = \beta_s^o(k) F_s^o \left( \sum_{q=1}^l w_{sq}^o(k) \beta_q^h(k) z_q^{(k)} - \Theta_s^o(k) \right)$$

$$\delta_s^{(k)} = (y_{ds} - y_s^{(k)}) F_s^{o'} \left( \sum_{q=1}^l w_{sq}^o(k) \beta_q^h(k) z_q^{(k)} - \Theta_s^o(k) \right) \beta_s^o(k)$$

we can now use the fact that  $F'(v) = F(v)(1-F(v))$   
 For the sigmoid function and the  
 algorithm becomes:

$s = 1(1)m$  output neurons.  
 $j = 1(1)l$  hidden neurons.  
 $i = 1(1)n$  input neurons.

$$\left\{ \begin{aligned} \omega_{sj}^{o(k+1)} &= \omega_{sj}^{o(k)} + \eta \delta_s^{(k)} \beta_j^{h(k)} z_j^{(k)} \\ \Theta_s^{o(k+1)} &= \Theta_s^{o(k)} - \eta \delta_s^{(k)} \\ \beta_s^{o(k+1)} &= \beta_s^{o(k)} + \eta (y_{ds} - y_s^{(k)}) \frac{y_s^{(k)}}{\beta_s^{o(k)}} \end{aligned} \right.$$

$$\left\{ \begin{aligned} \omega_{ji}^{h(k+1)} &= \omega_{ji}^{h(k)} + \eta \Delta_j^{(k)} \beta_j^{h(k)} F_j^h(v_j^{(k)}) (1 - F_j^h(v_j^{(k)})) \alpha_{di} \\ \Theta_j^{h(k+1)} &= \Theta_j^{h(k)} - \eta \Delta_j^{(k)} \beta_j^{h(k)} F_j^h(v_j^{(k)}) (1 - F_j^h(v_j^{(k)})) \\ \beta_j^{h(k+1)} &= \beta_j^{h(k)} + \eta \Delta_j^{(k)} z_j^{(k)} \end{aligned} \right.$$

$$v_j^{(k)} = \sum_{i=1}^n \omega_{ji}^{h(k)} \alpha_{di} - \Theta_j^{h(k)}$$

$$z_j^{(k)} = F_j^h(v_j^{(k)})$$

$$y_s^{(k)} = \beta_s^{o(k)} F_s^o \left( \sum_{q=1}^l \omega_{sq}^{o(k)} \beta_q^{h(k)} z_q^{(k)} - \Theta_s^{o(k)} \right)$$

$$\delta_s^{(k)} = (y_{ds} - y_s^{(k)}) F_s^o(u_s^{(k)}) (1 - F_s^o(u_s^{(k)})) \beta_s^{o(k)}$$

$$\Delta_j^{(k)} = \sum_{p=1}^m \delta_p^{(k)} \omega_{pj}^{o(k)}$$