$$y = f(x)$$

Neural Networks.



$$\underline{W^T X}$$

$W_1 x_1$

$+ W_2 x_2$

$+ W_3 x_3$

$+ W_4 x_4$

$+ W_5 x_5$

Thresholded perception



$$W^T x \geq 0 \quad \longrightarrow 1$$
$$\longrightarrow -1$$



$W^T x$

<u>Unit</u>

<u>Layer</u>

## Sigmoid unit

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

## March 12 Friday

- PA2 is out
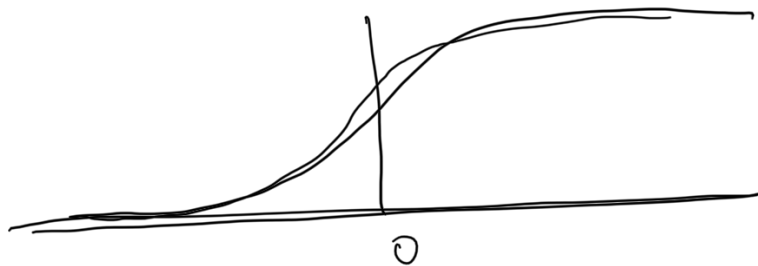
- Mid-term next friday

D = 2

Input

$x_1$

$x_n$

Input

Output

## # of units in a layer ≡ width of a layer

Sigmoid

tanh

Relu

$f()$

$\omega^T x$

$\max(0, w^T x)$

$\sigma\left(W_1^{(1)^T} \cancel{x}\right)$

$W_1^{(1)}$

$W_2^{(1)}$

$W_3^{(1)}$

$x_1$

$x_2$

$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$W_1^{(2)}$

$\begin{bmatrix} \sigma\left(W_1^{(1)^T} x\right) \\ \sigma\left(\omega_2^{(1)^T} x\right) \\ \left(\phantom{x}^{(1)^T} x\right) \end{bmatrix}$

$\sigma\left(w_2^{(1)T}x\right)$

$\sigma\left(w_3^{(1)T}x\right)$

$\sigma(w_3 \;\;)$

net

$\sigma(net)$

$x_1$

$w_1$
$w_2$
$w_3$
$w_4$

$w^T u$

$O$

$u_2$

$u_3$

$u_4$

unit, node, neuron

**Sigmoid**

$$sigmoid(z) = \frac{1}{1+e^{-z}}$$

**Tanh**

$$tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

# Simple example

$D = 2$ (2-d data)

$M = 2$

$k = 3$ (3- outputs/ classes)

## Assume sigmoid activation

Data: $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$



$\begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$  $\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$

$W_1^{(1)}$  $W_1^{(2)}$

$\begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$  $\begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$

$W_2^{(1)}$  $W_2^{(2)}$

$\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$

$W_3^{(r)}$

Input
Bias

hidden
Bias

## At hidden unit 1

$$net_1^{(1)} = W_1^{(1)T} x$$

$$= \begin{bmatrix} 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 4$$

$$Z_1 = \sigma(\text{net}_1^{(1)}) = \sigma(4)$$

$$= 0.98$$

## At hidden unit 2

$$\text{net}_2^{(1)} = w_2^{(1)} X$$

$$= \begin{bmatrix} 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 4$$

$$Z_2 = \sigma(\text{net}_2^{(1)}) = \sigma(4)$$

$$= 0.98 \qquad Z = \begin{bmatrix} z_1 \\ z_2 \\ 1 \end{bmatrix}$$

## At output unit 1

$$\text{net}_1^{(2)} = w_1^{(2)\,T} Z$$

$$= \begin{bmatrix} 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0.98 \\ 0.98 \\ 1 \end{bmatrix} = 1.98$$

$$O_1 = \sigma(1.98) = \frac{1}{1+\exp(-1.98)}$$

$$= 0.88$$

## At output unit 2

$$\text{net}_2^{(2)} \qquad w^{(2)\,T} \qquad Z$$

$$net_2 = \tilde{w}_2 \, z$$

$$= \begin{bmatrix} 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.98 \\ 0.98 \\ 1 \end{bmatrix} = 1.98$$

$$O_2 = \sigma(1.98) = 0.88$$

### At output unit 3

$$net_3^{(2)} = W_3^{(2)T} z = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.98 \\ 0.98 \\ 1 \end{bmatrix}$$

$$= 2.96$$

$$O_3 = \sigma(2.96) = 0.95$$

Data: $\mathbf{x} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$



$\begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$

$\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$

$W_1^{(1)}$  0.98  0.98  0.98

$\begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$

$W_1^{(2)}$  $O_1 = 0.88$

$\begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$

$W_2^{(1)}$  0.98  0.98  0.98

$W_2^{(2)}$  $O_2 = 0.88$

$\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$

$W_3^{(r)}$  $O_3 = 0.95$

Input
Bias

hidden

Softmax output :

$$\left[ \frac{O_1}{O_1 + O_2 + O_3} , \frac{O_2}{O_1 + O_2 + O_3} , \frac{O_3}{O_1 + O_3 + O_3} \right]$$

---

$W_1^{(1)}$   $3 \times 1$

$W_2^{(1)}$   $3 \times 1$

$W^{(1)}$

$W_1^{(2)}$

$W_2^{(2)}$

$W_3^{(2)}$

$W^{(2)}$

$$\sigma \left( \underset{2 \times 3}{W^{(1)}} \underset{3 \times 1}{X} \right) = \left[ \frac{Z_1}{Z_2} \right] \qquad Z = \begin{bmatrix} Z_1 \\ Z_2 \\ 1 \end{bmatrix}$$

$$\sigma \left( W^{(2)} Z \right) = \begin{bmatrix} O_1 \\ O_2 \\ O_3 \end{bmatrix}$$

---

$O_1 = Y_1$
$O_2 = Y_2$
$O_3 = Y_3$

$$J \left( W_1^{(1)}, \cdots, W^{(2)}, \cdots \right)$$

$$= \sum_{i}^{N} J_i$$

$$J_i = \frac{1}{2} \sum_{\ell=1}^{k} (y_{i\ell} - O_{i\ell})^2$$

$O_{i\ell} \rightarrow$ output for the $i^{th}$ training example at output unit $\ell$

$$J = \frac{1}{2} \sum_{i=1}^{N} \sum_{\ell=1}^{k} (y_{i\ell} - O_{i\ell})^2$$

$y_{i\ell} \rightarrow$ true output for the $i^{th}$ example at $\ell^{th}$ output unit.

E.g.        3 class classifier.

| x  D= 2 | y |
|---------|---|
| 3·7, 4·8 | 2 |
| 3·4, 1·2 | 1 |

One-of-k encoding

Dummy encoding

Let K = 3

2 ⟶ | 0 | 1 | 0 |

1 ⟶ | 1 | 0 | 0 |

3 ⟶ | 0 | 0 | 1 |

# Gradient Descent

$J \longrightarrow$ is a function of all the weights.

$\partial J$

## Wednesday    March 17

## Notation

### Subscripts

| | | |
|---|---|---|
| $i$ | Training example | $X_i$ |
| $p$ | Feature | $X_{ip}$ |
| $j$ | Hidden layer unit | $W_j^{(1)}$ |
| $l$ | Output layer unit | $W_l^{(2)}$ |

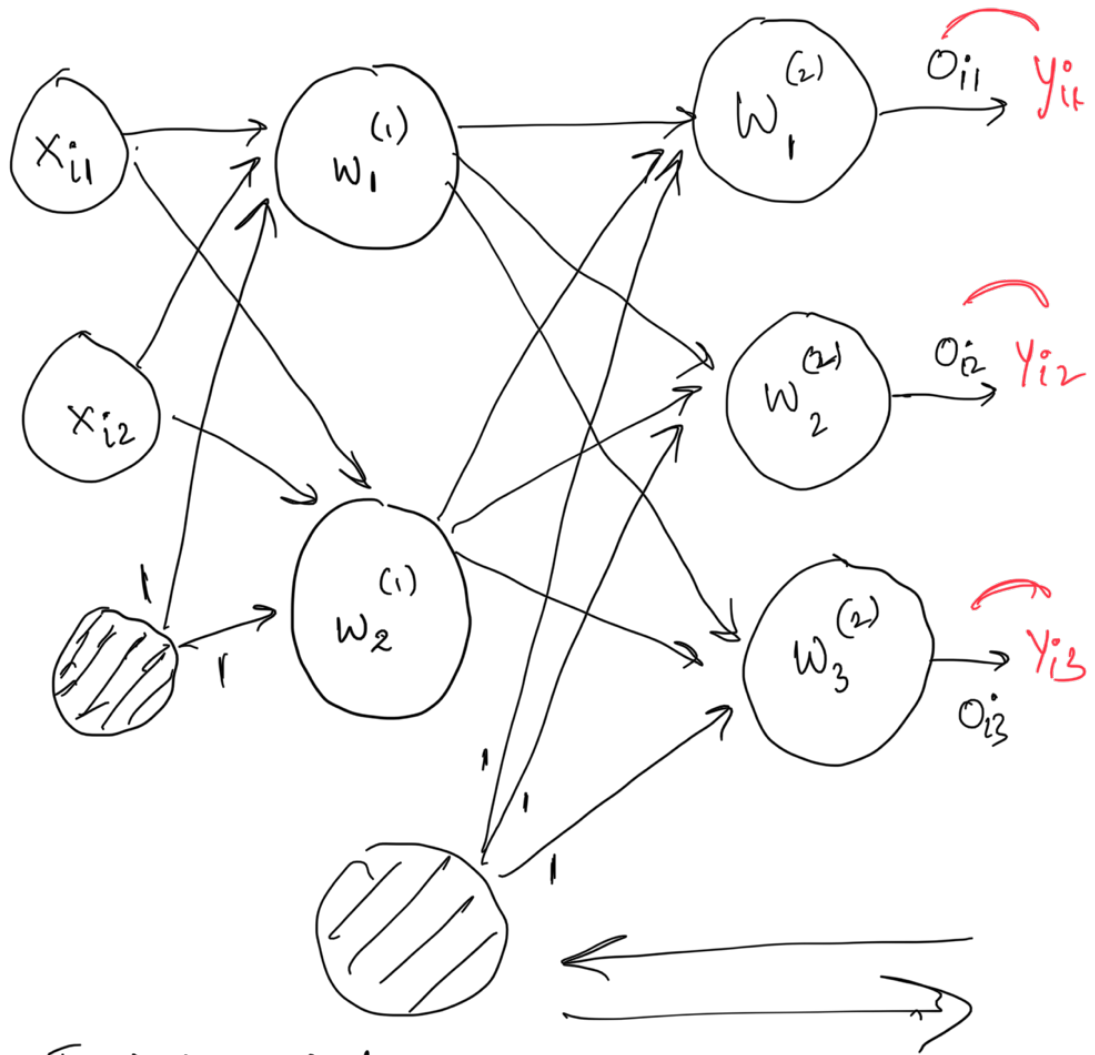### Variables

| | | |
|---|---|---|
| $D+1$ | $X_i$ | Input vector for $i^{th}$ training example |
| $D+1$ | $W_j^{(1)}$ | weight vector at $j^{th}$ hidden unit |
| $M+1$ | $W_l^{(2)}$ | weight vector at $l^{th}$ output unit |
| | $Z_j$ | output of the $j^{th}$ hidden unit |
| | $O_l$ | output    the $l^{th}$ output unit |
| $K$ | $Y_i$ | 1-of-K  true output for $i^{th}$ training example |

$X_{i1}$

$W_1^{(1)}$

$W_1^{(2)}$ → $O_{i1}$ → $y_{i1}$

$X_{i2}$

$W_2^{(2)}$ → $O_{i2}$ → $y_{i2}$

$W_2^{(1)}$

$W_3^{(2)}$ → $y_{i3}$

$O_{i3}$

Training data

$X_1$ $\quad$ $Y_1$ $\quad$ 1 $\quad$ | 1 0 0 |

$X_2$ $\quad$ $Y_2$ $\quad$ 1 $\quad$ | 1 0 0 |

$X_3$ $\quad$ $Y_3$ $\quad$ 2 $\quad$ | 0 1 0 |

$x_i$

$x_N$

$y_i$ ⟩ ₃

$y_N$ ₁

| 0 | 0 | 1 |

$J$

$W^{(1)}$
$W_2^{(1)}$
$W_M^{(1)}$
$W_1^{(2)}$
$W_2^{(2)}$
$W_K^{(2)}$

$\nabla$

$x_{i1}$

$x_{i2}$

$W_1^{(1)}$

$W_2^{(1)}$

$W_1^{(2)}$ → $O_{i1}$ → $y_{i1}$

$W_2^{(2)}$ → $O_{i2}$ → $y_{i2}$

$W_3^{(2)}$ → $y_{i3}$ , $O_{i3}$

$J =$

$J$

$$\frac{\partial J}{\partial w_{jp}^{(1)}} \qquad \frac{\partial J}{\partial w_{\ell j}^{(2)}}$$

scalar

$$w_{jp}^{(1)} \leftarrow w_{jp}^{(1)} - \eta \frac{\partial J}{\partial w_{jp}^{(1)}}$$

$$w_{\ell j}^{(2)} \leftarrow w_{\ell j}^{(2)} - \eta \frac{\partial J}{\partial w_{\ell j}^{(2)}}$$
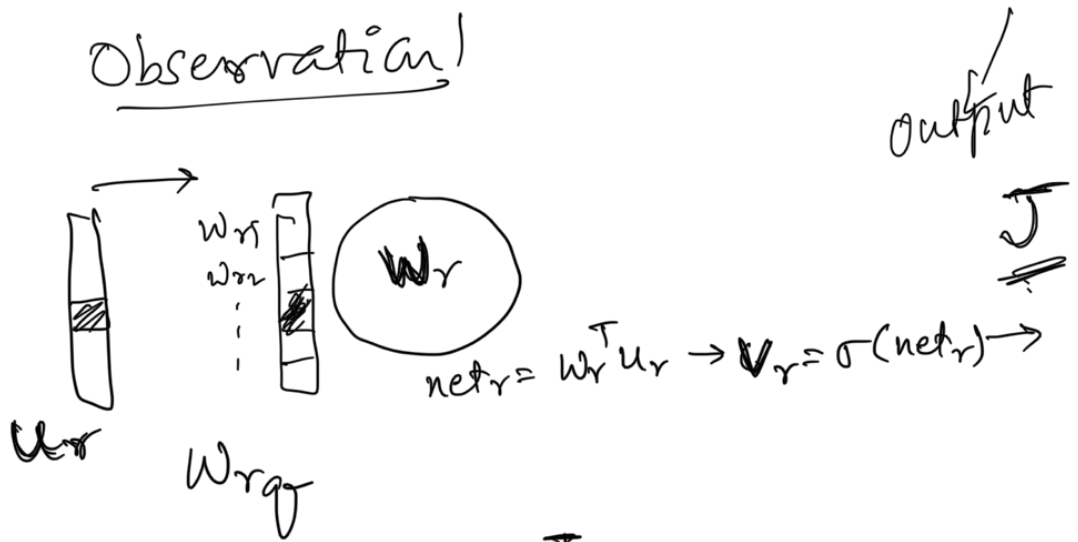
Loop 1 $\Rightarrow$ Assume we have
only 1 training example.

$$J = \frac{1}{2} \sum_{i=1}^{N} \sum_{\ell=1}^{K} (y_{i\ell} - O_{i\ell})^2$$

$$\boxed{J = \frac{1}{2} \sum_{\ell=1}^{K} (y_\ell - O_\ell)^2}$$

---

## Observation!

output

$J$



$net_r = W_r^T U_r \rightarrow V_r = \sigma(net_r) \rightarrow$

$W_{rq}$

$$net_r = W_r^T U_r$$

$$= \sum_{q} W_{rq} \, U_{rq}$$

$$\frac{\partial J}{\partial W_{rq}} = \frac{\partial J}{\partial net_r} \boxed{\frac{\partial net_r}{\partial W_{rq}}}$$

$\leftarrow$ chain
ruled

$$\frac{\partial net_r}{\partial w_{rq}} = U_{rq}$$

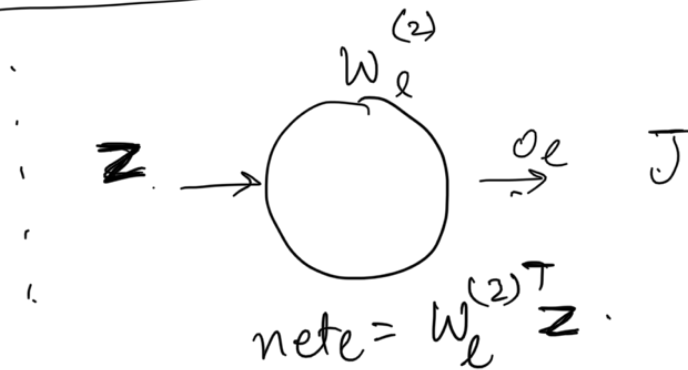$$\boxed{\frac{\partial J}{\partial W_{rq}} = \frac{\partial J}{\partial net_r} U_{rq} \underline{\hspace{1cm}}} \quad \swarrow$$

---

## Observation 2

For $l^{th}$ output unit :

$$\boxed{\frac{\partial J}{\partial net_l} = \frac{\partial J}{\partial O_l} \frac{\partial O_l}{\partial net_l}}$$

$W_l^{(2)}$

$\mathbf{Z} \longrightarrow \bigcirc \xrightarrow{O_l} J$

$$net_l = W_l^{(2)T} \mathbf{Z} .$$

$$O_l = \sigma(net_l) :$$

$$J = \frac{1}{2} \sum_{l=1}^{K} (y_l - O_l)^2$$

what is $\dfrac{\partial J}{\partial O_l}$ ?

$$J = \frac{1}{2}\left[ (y_1 - O_1) + (y_2 - O_2) + (y_3 - O_3) \right.$$
$$\left. + \cdots + (y_\ell - O_\ell)^2 + \cdots \right]$$

$$\boxed{\begin{aligned} \frac{\partial J}{\partial O_\ell} &= \frac{1}{2}\, 2\,(y_\ell - O_\ell)\,(-1) \\ &= -(y_\ell - O_\ell) \end{aligned}}$$ #1

what is $\dfrac{\partial O_\ell}{\partial net_\ell}$

$$O_\ell = \sigma(net_\ell)$$
$$= \frac{1}{1 + \exp(-net_\ell)}$$

$$\frac{\partial O_\ell}{\partial net_\ell} = \frac{\partial}{\partial net_\ell}\left[ \frac{1}{1 + e^{-net_\ell}} \right]$$

$$= -\frac{1}{(1 + e^{-net_\ell})^2}\,(-e^{-net_\ell})$$

$$\boxed{\frac{\partial O_\ell}{\partial net_\ell} = \frac{e^{-net_\ell}}{(1 + e^{-net_\ell})^2} = O_\ell(1 - O_\ell)}$$ #2

Combining #1 & #2 :

$$\frac{\partial J}{\partial n} = -O_\ell(1 - O_\ell)(y_\ell - O_\ell)$$

$$\frac{\partial net_\ell}{}$$

Let $\boxed{\delta_\ell = O_\ell(1-O_\ell)(y_\ell - O_\ell)}$

$$\frac{\partial J}{\partial W_{\ell j}^{(2)}} = \frac{\partial J}{\partial net_j} z_j$$

$$= - O_\ell(1-O_\ell)(y_\ell - O_\ell) z_j$$

$$= - \delta_\ell z_j$$

Update rule for output layer.

$$\boxed{W_{\ell j}^{(2)} \leftarrow \underset{old}{W_{\ell j}^{(2)}} + \eta \, \delta_\ell z_j}$$

Preview:

$$\boxed{W_{jp}^{(1)} \leftarrow \underset{old}{W_{jp}^{(1)}} + \eta \, \boxed{\delta_j} x_u}$$

$$\delta_j = \text{function of} \quad \delta_\ell \text{ s.}$$

Monday Mar 22

What do we need?

$$\frac{\partial J}{\partial w_{rq}}$$

a general unit $(r)$



$u_1$
$u_2$
$u_3$
$\cdots$
$u_{qr}$
$\cdots$

$w_{r1}$
$w_{rr}$
$w_{r3}$
$\cdots$
$w_{r\sigma}$
$\cdots$

$\to net_r$
$= w_r^T u$

$\to \sigma(net_r)$

$W_r$

input to unit $(r)$

Observation 1

$$\frac{\partial J}{\partial w_{rq}} = \frac{\partial J}{\partial net_r} \frac{\partial net_r}{\partial w_{rq}} = \frac{\partial J}{\partial net_r} u_q$$

Observation 2

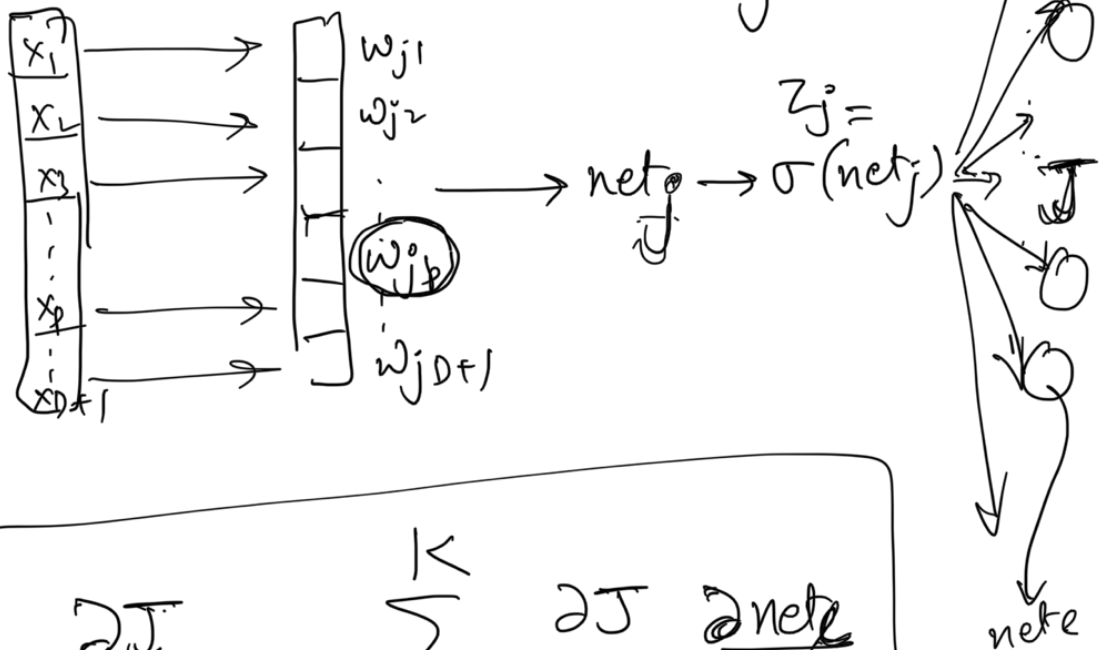For weight on an output unit:

$$\frac{\partial J}{\partial} = - \delta_\ell z_j$$

$$\frac{\partial}{\partial w_{\ell j}}$$

$$\boxed{\delta_\ell = (y_\ell - O_\ell)\, O_\ell (1 - O_\ell)}$$

for a hidden unit $\mathbf{w_j}$



$w_{j1}$
$w_{j2}$
$\vdots$
$w_{jb}$
$\vdots$
$w_{jD+1}$

$$z_j = net_j \to \sigma(net_j)$$

$net_\ell$

$$\boxed{\frac{\partial J}{\partial net_j} = \sum_{\ell=1}^{K} \frac{\partial J}{\partial net_\ell} \frac{\partial net_\ell}{\partial net_j}}$$

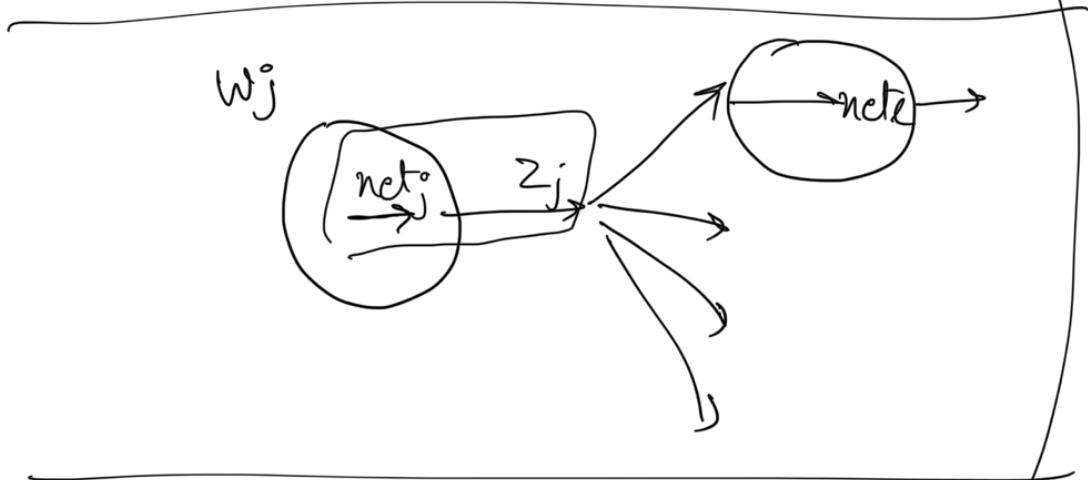We already know: $\dfrac{\partial J}{\partial net_\ell} = -\delta_\ell$

(Check Wed. derivation)

$$\frac{\partial J}{\partial ??} = -\sum^{K} \delta_\ell \frac{\partial net_\ell}{\partial ??}$$

$$\frac{\partial}{\partial net_j} = -\sum_{\ell=}^{K} \left[ \delta_\ell \left( \frac{\partial net_\ell}{\partial z_j^o} \right) \left( \frac{\partial z_j}{\partial net_j} \right) \right]$$

$l=1$ over the sum, $\frac{\partial net_j}{}$ at top



$$z_j^o = \sigma(net_j)$$

$$\Rightarrow \quad \boxed{\frac{\partial z_j^o}{\partial net_j} = z_j^o (1 - z_j^o)}$$

$$\frac{\partial J}{\partial net_j^o} = -\sum_{\ell=1}^{K} \left[ \delta_\ell \left( \frac{\partial net_\ell}{\partial z_j^o} \right) \right] z_j^o (1 - z_j^o)$$

$$\frac{\partial net_\ell}{\partial z_j^o}$$

Recall that $\quad net_\ell = \sum^{M+1} W_{\ell j} z_j^o$

$$\frac{\partial net_\ell}{\partial z_j} = w_{\ell j} \qquad j=1$$

$$\Rightarrow \quad \frac{\partial J}{\partial net_j} = -z_j(1-z_j)\left(\sum_{\ell=1}^{K} \delta_\ell \, w_{\ell j}\right)$$

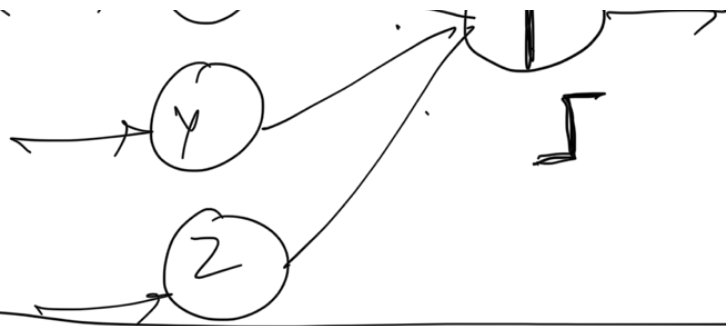$$\frac{\partial J}{\partial w_{jp}} = -z_j(1-z_j)\left(\sum_{\ell=1}^{K} \delta_\ell \, w_{\ell j}\right) x_p$$

$$\delta_j = z_j(1-z_j)\left(\sum_{\ell=1}^{K} \delta_\ell \, w_{\ell j}\right)$$

$$\boxed{\frac{\partial J}{\partial w_{jp}} = -\delta_j \, x_p}$$

$$\boxed{\begin{array}{l} w_{jp} = w_{jp} + \eta \, \delta_j \, x_p \\[2mm] w_{\ell j} = w_{\ell j} + \eta \, \delta_\ell \, z_j \end{array}}$$
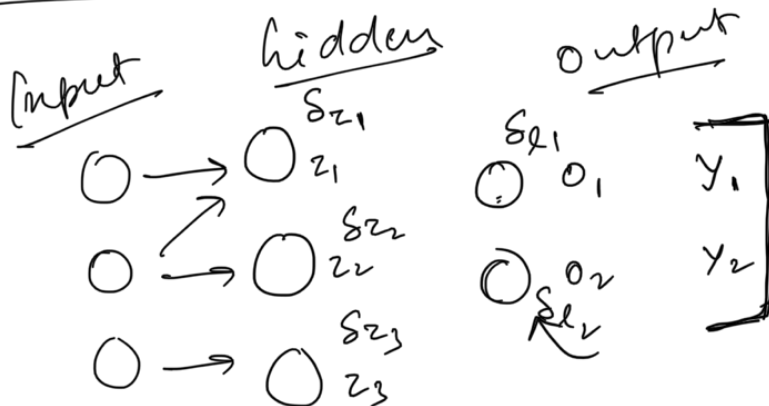
Wed Mar 24

## Announcements

- Mid-term grading questions
    - Chandola Office Hours
      3-5 PM on Friday

- PA1 Grading - TAs

- Gradiance 7 - delayed

- Friday lecture "Cancelled"
    ↳ Dean's office hours.
        1.50 PM - 2.40 PM

    $\boxed{\text{No code}}$    $\boxed{\text{PA2 only}}$

---

Input     hidden     output



$\delta_{z1}$
$z_1$

$\delta_{z2}$
$z_2$

$\delta_{z3}$
$z_3$

$\delta_{o1}$
$o_1$   $y_1$

$o_2$   $y_2$
$\delta_{o2}$

$$\xrightarrow{\hspace{3cm}}$$

Feed

First Calculate all $\delta_s$ $(\delta_{l_s} \text{ \& } \delta_{z_s})$

Then update all the weights.

objective fn./ loss function

error $(\delta) \rightarrow$ unit

Momentum

$$\underset{\text{new}}{\underbrace{W_{jp}^{(1)}}} \leftarrow \underset{\text{old}}{\underbrace{W_{jp}^{(1)}}} - \eta \frac{\partial J}{\partial w_{jp}}$$

$$\underset{\text{new\_2}}{\underbrace{W_{jp}^{(1)}}} = \alpha \underset{\text{old}}{\underbrace{W_{jp}^{(1)}}} + (1-\alpha) \underset{\text{new}}{\underbrace{W_{jp}^{(1)}}}$$

$$\boxed{0 \leq \alpha \leq 1}$$

Universal function approximator

$$x \xrightarrow{\quad f(\,) \quad} y$$

$$\underset{\bullet}{x} \quad \underset{}{y}$$

$$\frac{\partial \tilde{J}}{\partial w_{jp}^{(1)}} = \frac{\partial J}{\partial w_{jp}^{(1)}} + \frac{\lambda}{2N} 2 w_{jp}^{(1)}$$



$\rightarrow$ MLP

Convolutional neural networks.