8th October, 2021



$j = 1 \cdots$

$i = 0 \cdots d$

$w_{ji}$

$k = 1$ to $c$

$w_{ji}$

$i$

$w_{ji}$

$x_1$

$x_2$

$x_3$
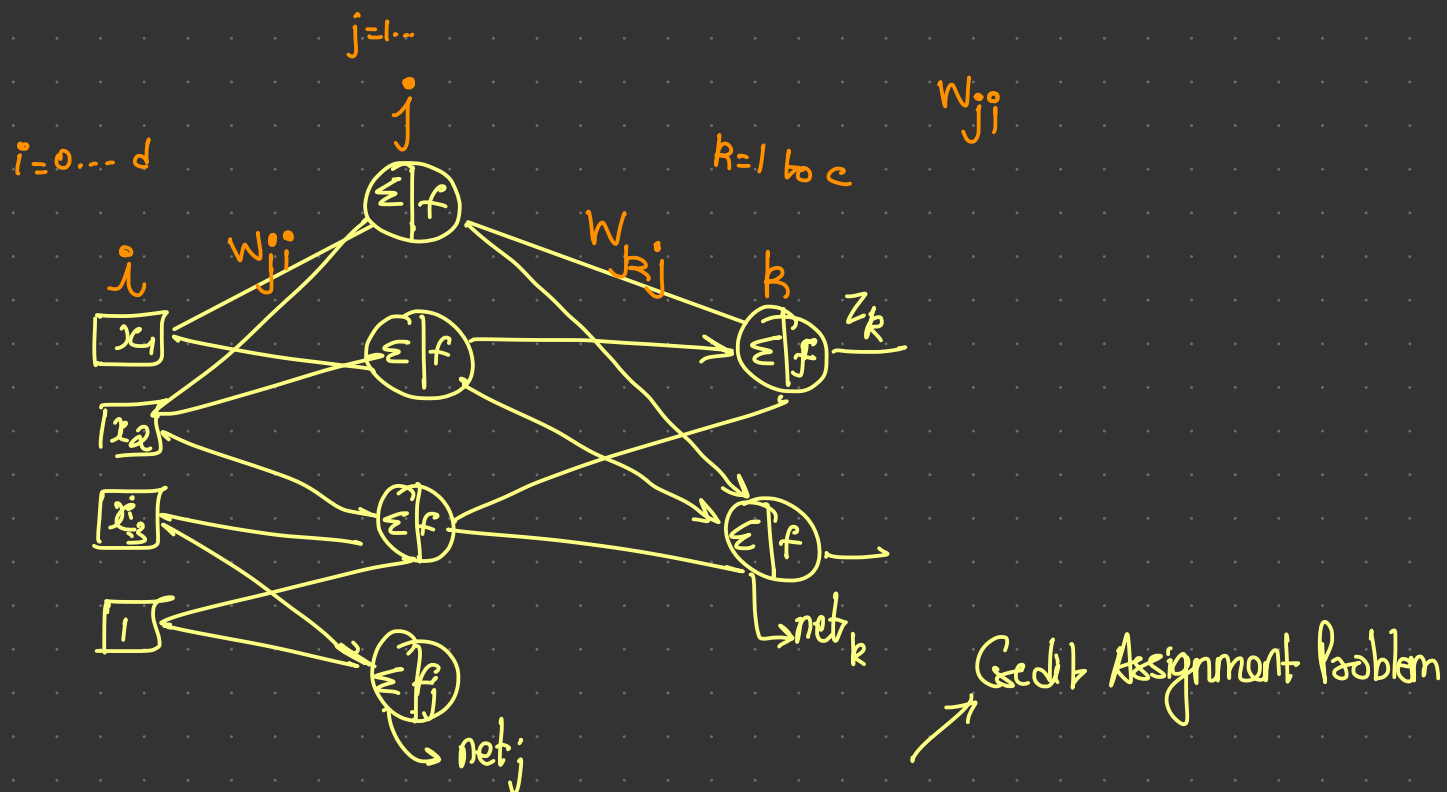
$1$

$W_{kj}$

$k$

$z_k$

$net_k$

$net_j$

→ Credit Assignment Problem

Weights affect neuron just like a single perceptron.

$\frac{\partial J}{\partial w_{ji}}$ → problem because in general there would be subsequent changes in $w_{kj} \cdot \Delta T$

Back propagation Algorithms

$$J(w) = \frac{1}{2}\sum_{k=1}^{c}(t_k - z_k)^2 = \frac{1}{2}\|t - z\|^2$$

↳ Vector form

$$\Delta w = -\eta \frac{\partial J}{\partial w} \quad \Big/ \quad w = w - \eta \frac{\partial J}{\partial w}$$

$$\boxed{\Delta w_{pq} = -\eta \frac{\partial J}{\partial w_{pq}}} \rightarrow \text{Compute for both I/P \& O/P layer}$$

# Updating $w_{kj}$

$$\frac{\partial J}{\partial w_{kj}} = \underbrace{\frac{\partial J}{\partial net_k} \cdot \frac{\partial net_k}{\partial w_{kj}}}_{\text{Chain rule}} = -\delta_k \frac{\partial net_k}{\partial w_{kj}}$$

$$\delta_k = \frac{-\partial J}{\partial net_k} = \frac{-\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial net_k} = (t_k - z_k) f'(net_k) \implies \delta_k = (t_k - z_k) f'(net_k) \quad ①$$

$$\frac{\partial J}{\partial w_{kj}} = -\delta_k \frac{\partial net_k}{\partial w_{kj}}$$

$$\frac{\partial J}{\partial w_{kj}} = -\delta_k \cdot y_j \qquad \boxed{\Delta w_{kj} = \eta \delta_k y_j} \quad ②$$

# B) Updating $w_{ji}$

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial y_j} \cdot \underbrace{\frac{\partial y_j}{\partial net_j}}_{} \cdot \underbrace{\frac{\partial net_j}{\partial w_{ji}}}_{} = \frac{\partial J}{\partial y_j} \cdot f'(net_k) \cdot x_i$$

$$f'(net_k) \qquad x_i$$

$$\frac{\partial J}{\partial y_j} = \frac{\partial}{\partial y_j}\left[\frac{1}{2}\sum_{k=1}^{C}(t_k - z_k)^2\right) = -\sum_{k=1}^{C}(t_k - z_k)\frac{\partial z_k}{\partial y_j}$$

$$= -\sum_{k=1}^{C}(t_k - z_k)\frac{\partial z_k}{\partial net_k} \frac{\partial net_k}{\partial y_j}$$

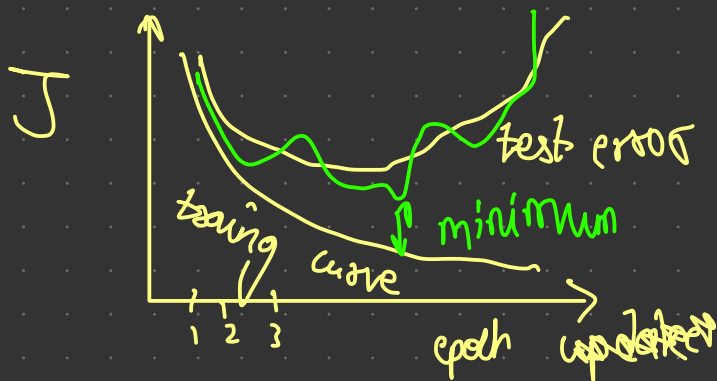$$= -\sum_{k=1}^{C}(t_k - z_k) \cdot f'(net_k) \cdot w_{kj}$$

$$\boxed{\delta_j = f'(net_j)\sum_{k=1}^{C} w_{kj}\, \delta_k} \quad ③$$

$$\frac{\partial J}{\partial w_{ji}} = -\delta_j x_i \quad \left| \quad \boxed{\Delta w_{ji} = \eta \, \delta_j x_i} \right. \;\; —\!2\ell$$

$$w^{t+1} = w^t + \Delta w$$

On an average increasing weights can increase/decrease

## Learning curves



Validation set : Split part of training data
Iterate only till the minima is reached

## Backprop Algorithms

a) Stochastic backprops
    ↳ Update after each sample
  / ↳ Initialize $n_H$, $w$, $\theta$, $\eta$, $m=0$
        while ( ) do
            $m = m+1$
            $x^m \rightarrow$ Randomly chosen
            $w_{kj} = w_{kj} + \eta \delta_k y_j$

$$W_{ji} \leftarrow W_{ji} + \eta \, \delta_j x_i$$

$$y$$

return $(W)$

Epoch's $\rightarrow$ Iterations

## b) Batch Backprop

for $(m = 0$ to $N$

$\{ \quad x^n \leftarrow$ selecting one training sample

$\Delta W_{ji} = \Delta W_{ji} + \eta \, \delta_j x_i$

$\Delta W_{kj} = \Delta W_{kj} + \eta \, \delta_k y_j$

$y$

$W_{ji} += \Delta W_{ji}$ | $W_{kj} += \Delta W_{kj}$

$\rightarrow$ Update only once per epoch

## c) Minibatch backprop



| | 1 | 2 | $\hat{}$ | - | - | - | - | $N$ |
|---|---|---|---|---|---|---|---|---|
| MB$_1$ | MB$_2$ | MB$_3$ | $\vdash$ | - - - | | MB$_{N/B}$ | | |

$\Delta W$         $\Delta W$

$W + = \Delta W$    $W + = \Delta W$

## 11/10/2021

### ① Activation function

derivative

a) Continuity of $f()$ $f'()$ & smooth as we need to get ∧

b) $f()$ should saturate. [have some type of max & min]
$\qquad - f() \in [0, 1]$

c) To make error function simple,
$\qquad$ so $f()$ is monotonic & simple

d) $f()$ should be linear around $0$

→ Sigmoid satisfies all three properties

② Parameters of sigmoid :

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

To control unsaturated range, add scaling factor

$$f(net) = \frac{2a}{1 + e^{-b \cdot net}} - a$$

$a, b$ are parameters for appropriate range

Popular values : $\underbrace{a = 1.716, \; b = 2/3}$

⇓

$min/max \left( f''(net) \right)$ is $\pm 2$

$f( \; ) \rightarrow$ Antisymmetric

③ Input values :

If input values are too small, net is too small, sigmoid will become linear perceptrons
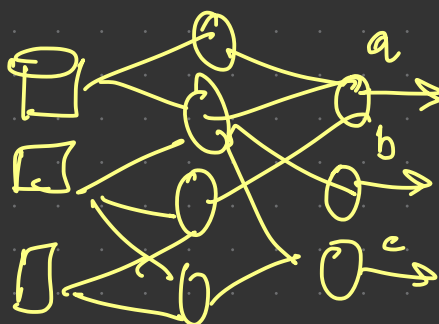↝ 0 mean, unit variance

④ Weight initialization
→ Should be in appropriate range
→ Should be random
→ Uniform weights from

$$\sum_{j=1}^{d} w_j y_j$$

$$w = U\left[ -\hat{w}, \hat{w} \right]$$
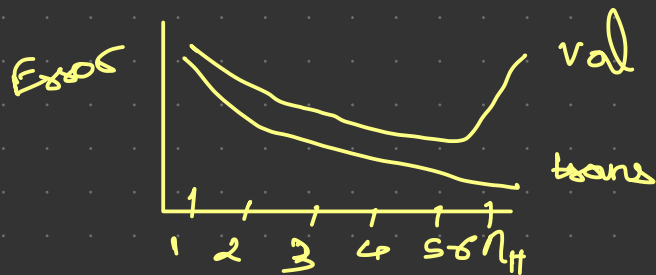
$$\hat{w} = 1/\sqrt{d}$$

⑤ Target Value :



⑥ Training with noise

⑦ Manufacture new data :

Adding new data ⇒ Data Aug

## ⑧ # of hidden units:

| I/P layer | hidden layer | O/P layer |
|-----------|--------------|-----------|
| $d+1$ | $n_H$ | # of class |

Error



val

trans

1 2 3 4 5 6 $n_H$

## ⑨ Learning rate

$$\eta_{opt} = \left( \frac{\delta^2 J}{\delta w^2} \right)^{-1} \qquad \eta_{max} = 2\eta_{opt}$$

In practice:

$\eta = 0.1$ as start

$J$ diverges → reduce $\eta$

$J$ is slow → increase $\eta$

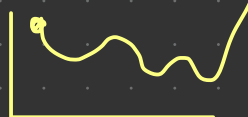$\eta$ schedule → $\eta = 0.1$

## ⑩ Momentum

Multiple local minimas are possible

$$\Delta W_M \Rightarrow (1-\alpha)\Delta W_m + \alpha \Delta W_{m-1}$$

How much momentum we have

# (1) Weight Decay [Revise]

$$w \leftarrow w(1 - \varepsilon)$$

Reduce everytime a little

Unwanted/ useless weights will all go to $0$

Useful weights are those

More expressive power

$$J_w = J(\ ) + \frac{2\varepsilon}{n} w^T w$$

Regulariser

# (2) Use of Hints (Revise)

 $\rightarrow$ Adding attributes using domain knowlege

# (1B) Loss Fn / Criterion fn:

a) Square Error $\quad J(w) = \sum_R (t_R - z_R)^2 \quad R \rightarrow$ even

b) Minkowski $\quad J(w) = \sum_R |(t_R - z_R|)$

c) Cross entropy

# (1c) # of hidden layers $J(w) = \sum_{R=1}^{c} t_R \ln (t_R / z_R)$

# hidden layer $> 1 \Rightarrow$ Deep network