# (MLP) → Backpropagation

nodes i
$l-1$ ↗   nodes $-j$
$l ↗ . j$

o/p



update     update

$$\frac{\partial(Loss)}{\partial(weight)}$$

$[W_{jk}^{(l)}$ = Weights of the connection that connects node K in layer $l-1$ to $j$ in $l.]$

$[Loss\ calculated$ and updates weight in backward direction based on the amount of Loss produced $]$

— Avg of weight if updated in BP.

## Steps :

1) passing data to model via forward propagation.

2) calculate loss on output.

3) SGD minimizes the loss.
   — Gradient is calculated via Backpropagation

$W_j^l$ — the vector that contains all weights connect

$Z_j^{(l)}$ — the input for node $j$ in layer $l$.

$g^{(l)}$ — the activation function used for layer $l$.

$a_j^{(l)}$ — activation o/p of node $j$ in layer $l$.

<u>Loss $C_0$</u>: given by.

$$(a_j^{(L)} - y_j)^2$$

↑ activation o/p @ node $j$

↗ desired o/p @ node '$j$'. } layer 'L'

To calculate total loss, we should sum. this squared difference for all the nodes $j$ in layer 'L'.

$$C_0 = \sum_{j=0}^{\eta-1} (a_j^{(L)} - y_j)^2$$

<u>Input $z_j^{(l)}$:</u>

I/P for node '$j$' in layer '$l$' is weighted sum of activation o/p from previous layer $(l-1)$

ex:

$$w_{jk}^{(l)} \, a_k^{(l-1)}$$

↳ o/p.

↳ weights.

Input to node '$j$' in layer '$l$' is expressed as:

$$z_j^{(l)} = \sum_{k=0}^{\eta-1} w_{jk}^{(l)} \, a_k^{(l-1)}$$

## Activation o/p : $a_j^{(\ell)}$

$a_j^{(\ell)}$ → It is the result of passing $Z_j^{(\ell)}$ to whatever activation function we choose to use.

Say $g^{(\ell)}$

$(\ell)$ → maps @ layer $\ell$.

Activation output of node $j$ in layer '$\ell$' is expressed as :

$$\boxed{a_j^{(\ell)} = g^{(\ell)}\left(z_j^{(\ell)}\right)}$$

The i/ps for node $j$ is a function of all weights connected to node $j$.

So, $z_j^{(L)}$ → func $\left(w_j^{(L)}\right)$.

$$z_j^{(L)}\left(w_j^{(L)}\right).$$

$$\boxed{\therefore C_{0j} = C_{0j}\left(a_j^{(L)}\left(z_j^{(L)}\left(w_j^{(L)}\right)\right)\right)} \quad \vdots$$

Loss @ node $j$

$$\boxed{C_0 = \sum_{j=0}^{n-1} C_{0j}}$$

Calculations : Derivative of Loss w·r·t weights. ④

$$\frac{\partial C_0}{\partial W_{12}^{(L)}}$$



o/p → $C_0$.

L-1        L

$$W_{12}^{(L)} \longrightarrow \boxed{z | a} \longrightarrow op - C_0$$

$C_0$ depends on $a_1^{(\ell)}$

$a_1^{(\ell)}$ depends on $z_1^{(\ell)}$      } Composition of functions

$z_1^{(\ell)}$ depends on $W_{12}^{(L)}$

$$\frac{\partial C_0}{\partial W_{12}^{(L)}} = \left(\overset{①}{\frac{\partial C_0}{\partial a_1^{(L)}}}\right) \cdot \left(\overset{Ⓓ}{\frac{\partial a_1^{(L)}}{\partial z_1^{(L)}}}\right) \cdot \left(\overset{②}{\frac{\partial z_1^{(L)}}{\partial W_{12}^{(L)}}}\right)$$

Consider this first term

$$C_0 = \sum_{d=0}^{n-1} (a_j^{(L)} - y_j)^2 \implies \frac{\partial C_0}{\partial a_j^{(L)}} = \frac{\partial}{\partial a_{ij}} \left((a_0^{(L)} - y_0)^2 + (a_1^{(L)} - y_0)^2 + (a_2^{(L)} - y_2)^2 + (a_3^{(L)} - y_3)^2\right)$$

derivative of Sum = Sum of derivatives :

$$\left[ \frac{\partial (\_ + \_ + \_)}{\partial x} = \frac{\partial (\ )}{\partial x} + \frac{\partial (\ )}{\partial x} + \frac{\partial (\ )}{\partial x} \right]$$

$$\frac{\partial C_0}{\partial a_1^{(L)}} = \frac{\partial}{\partial a_1^{(L)}} (a_0^{(L)} - y_0)^2 + \frac{\partial}{\partial a_2^{(L)}} (a_1^{(L)} - y_1)^2 + \frac{\partial}{\partial a_1^{(L)}} (a_2^{(L)} - y_2)^2 + \frac{\partial}{\partial a_2^{(L)}} (a_3^{(L)} - y_3)^2$$

$$\boxed{\frac{\partial C_0}{\partial a_1^{(L)}} \Rightarrow 2(a_1^{(L)} - y_1)} \quad \}\rightarrow \text{1st}$$

Second term in $\dfrac{\partial C_0 \ ②}{\partial W_{12}^{(L)}}$

we know:
$$a_j^{(L)} = g^{(L)}(z_j^{(L)}) \rightarrow \text{general}$$

$$a_1^{(L)} = g^{(L)}(z_1^{(L)}) \rightarrow \text{for node 'i' in layer L}.$$

$$\frac{\partial a_1^{(L)}}{\partial z_1^{(L)}} = \frac{\partial}{\partial z_1^{(L)}} g^{(L)}(z_1^{(L)})$$

$$\boxed{\frac{\partial a_1^{(L)}}{\partial z_1^{(L)}} = g'^{(L)}(z_1^{(L)})} \rightarrow \text{grad}$$

Third term : $\dfrac{\partial z_1^{(L)}}{\partial W_{12}^{(L)}}$   we know,
$$z_j^{(L)} = \sum_{K=0}^{n-1} W_{jK}^{(L)} a_K^{(L-1)}$$

$$\underline{J=1},$$

$$z_1^{(L)} = \sum_{k=0}^{n-1} w_{1k}^{(L)} q_k^{(L-1)}$$

$$\frac{\partial \left( \sum_{k=0}^{n-1} w_{1k}^{(L)} q_k^{(L-1)} \right)}{\partial w_{12}^{(L)}} = \frac{\partial \left( w_{10}^{(L)} q_0^{(L-1)} \right)}{\partial w_{12}^{(L)}}^{0} + \frac{\partial \left( w_{11}^{(L)} q_1^{(L-1)} \right)}{\partial w_{12}^{(L)}}^{0}$$

$$\left( + \frac{\partial \left( w_{12}^{(L)} q_2^{(L-1)} \right)}{\partial w_{12}^{(L)}} + \right) \dots$$

expand summation

$$+ \dots \dots \text{rest will be zero}$$

only this will remain:

$$\boxed{\frac{\partial z_1^{(L)}}{\partial w_{12}^{(L)}} = q_2^{(L-1)}} \to 3rd$$

Substitute 1st, 2nd, 3rd in Loss derivative;

$$\boxed{\frac{\partial C_0}{\partial w_{12}^{(L)}} = \left[ 2\left( q_1^{(L)} - y_1 \right) \right] \cdot \left[ g'^{(L)} \cdot \left( z_1^{(L)} \right) \right] \cdot \left[ q_2^{(L-1)} \right]}$$

↳ This is Loss Calculation w.r.t weight $w_{12}$ for 1 training Sample.

To consider all the training samples, we have to Avg:

$$\boxed{\frac{\partial C_0}{\partial w_{12}^{(L)}} = \frac{1}{n} \sum_{i=0}^{n-1} \frac{\partial C_i}{\partial w_{12}^{(L)}}}$$

→ We can calculate similarly for all other weights in Network or Layer.

⑤

Scanned with CamScanner

...processing the handwritten content...

Loss

$(L-2) \quad (L-1) \quad (L)$

$$\frac{\partial C_0}{\partial W_{22}^{(L-1)}} = \left(\frac{\partial C_0}{\partial q_2^{(L-1)}}\right) \cdot \left(\frac{\partial q_2^{(L-1)}}{\partial z_2^{(L-1)}}\right) \cdot \left(\frac{\partial z_2^{(L-1)}}{\partial W_{22}^{(L-1)}}\right)$$

node 2.

$$\boxed{\frac{\partial C_0}{\partial q_2^{(L-1)}}}$$

$$\boxed{\begin{array}{c} \text{In some books} \\ a = h \end{array}}$$

↳ To find this we need to find the product of the derivatives of the composed functions.

$$\underbrace{Z_j^{(L)}}_{\downarrow} = \sum_{K=0}^{n-1} W_{jK}^{(L)} \, a_K^{(L-1)}$$

$$\boxed{\therefore \frac{\partial z_j^{(L)}}{\partial q_2^{(L-1)}} = \frac{\partial}{\partial q_2^{(L-1)}} \sum_{K=0}^{n-1} W_{jK}^{(L)} q_K^{(L-1)}} \longrightarrow \text{expand}$$