# Backpropagation in ANN (MLP)

———————X———————X———————

## * Chain Rule of Derivatives :-

The Chain Rule of Derivatives is used to find derivative of composite function (i.e function of function).

Suppose,

$$y = f(u) \; ; \; y \text{ is a function of } u$$

& $u = g(x) \; ; \; u$ is function of $x$ then,

$$y = f(g(x)) \; ; \; y \text{ is function of } x$$

i.e

If $y$ is function of $u$ and $u$ is function of $x$ then chain rule of derivatives is,

$$\boxed{\dfrac{dy}{dx} = \dfrac{dy}{du} \cdot \dfrac{du}{dx}}$$

## * Backpropagation :-

Backpropagation is the algorithm used to train ANN models.

In Backpropagation we calculate the gradients of loss function w.r.t. weights and biases of the model. (i.e $\dfrac{\partial L}{\partial \omega}, \dfrac{\partial L}{\partial b}$). and update

the values of Weights and Biases (W & b) until we get the values of weights and biases for which the loss function is minimum.

We use the weight updation formula as,

$$W_{new} = \cancel{\pm} W_{old} - \eta * \left(\frac{\partial L}{\partial W_{old}}\right)$$

$$b_{new} = b_{old} - \eta * \left(\frac{\partial L}{\partial b_{old}}\right)$$

And the method used to calculate gradients is the 'Gradient Descent'.

## * Intution :-

Let's suppose we have dataset with two input columns iq & cgpa and our aim is to predict the salary of student (placed).

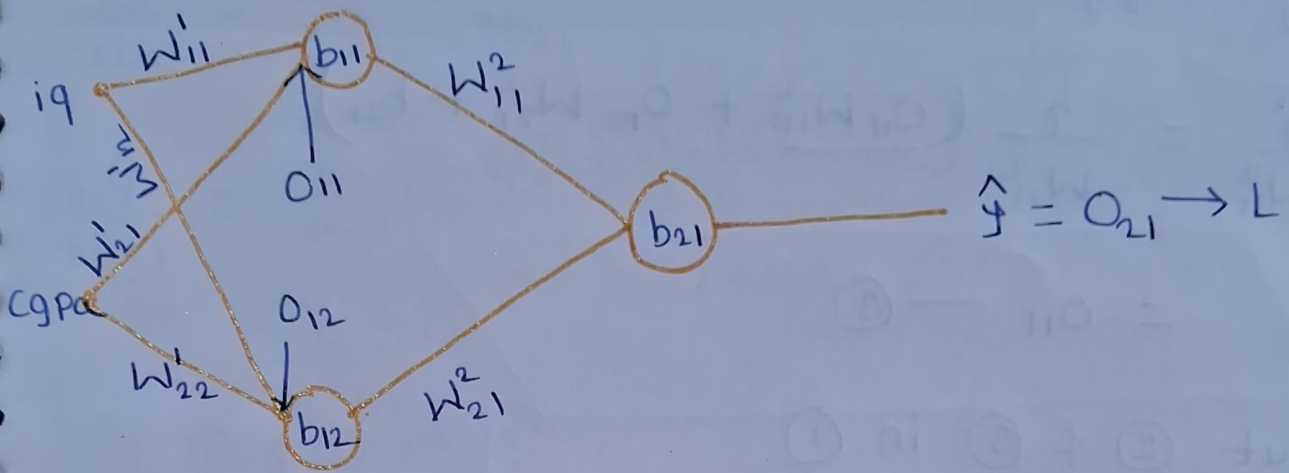So this is regression problem.

consider,

A.F (each layer) = linear

L.F (loss) = MSE

So the architecture is,

Two inputs, one hidden layer with two nodes and one output layer

Here, No. of trainable parameters are,

Layer 1 $= (2 \times 2) + 2 = 4 + 2 = 6$

Layer 2 $= (2 \times 1) + 1 = 2 + 1 = 3$

Total trainable parameters $= 6 + 3 = \underline{9}$

In this problem we have to calculate 9 derivatives.

(1) $\dfrac{\partial L}{\partial W_{11}^2}$    (2) $\dfrac{\partial L}{\partial W_{21}^2}$    (3) $\dfrac{\partial L}{\partial b_{21}}$    (4) $\dfrac{\partial L}{\partial W_{11}'}$    (5) $\dfrac{\partial L}{\partial W_{12}'}$

(6) $\dfrac{\partial L}{\partial W_{21}'}$    (7) $\dfrac{\partial L}{\partial W_{22}'}$    (8) $\dfrac{\partial L}{\partial b_{11}}$    (9) $\dfrac{\partial L}{\partial b_{12}}$

Let's calculate one by one

(1) $\dfrac{\partial L}{\partial W_{11}^2} = \dfrac{\partial L}{\partial \hat{y}} \times \dfrac{\partial \hat{y}}{\partial W_{11}^2}$ —— ①

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}}(y-\hat{y})^2 = -2(y-\hat{y}) - \text{②}$$

$$\frac{\partial \hat{y}}{\partial W_{11}^2} = \frac{\partial}{\partial W_{11}^2}(O_{11}W_{11}^2 + O_{12}W_{21}^2 + b_{21})$$

$$= O_{11} - \text{③}$$

Put ② & ③ in ①

so the first derivative is,

$$\boxed{\frac{\partial L}{\partial W_{11}^2} = -2(y-\hat{y})O_{11}}$$

② $$\frac{\partial L}{\partial W_{21}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial W_{21}^2} - \text{④}$$

$$\frac{\partial \hat{y}}{\partial W_{21}^2} = \frac{\partial}{\partial W_{21}^2}(O_{11}W_{11}^2 + O_{12}W_{21}^2 + b_{21})$$

$$= O_{12} - \text{⑤}$$

use ② & ⑤ in ④

$$\boxed{\frac{\partial L}{\partial W_{21}^2} = -2(y-\hat{y})O_{12}}$$

$$\frac{\partial L}{\partial b_{21}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial b_{21}} \quad\text{———}\;\text{⑥}$$

$$\frac{\partial \hat{y}}{\partial b_{21}} = \frac{\partial}{\partial b_{21}} (O_{11} W_{11}^2 + O_{12} W_{21}^2 + b_{21})$$

$$= 1 \quad\text{———}\;\text{⑦}$$

use ⑦ & ② in ⑥

$$\boxed{\frac{\partial \hat{y}}{\partial b_{21}} = -2(y - \hat{y})}$$

$$\frac{\partial L}{\partial W_{11}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial W_{11}'} \quad\text{———}\;\text{⑧}$$

$$\frac{\partial \hat{y}}{\partial O_{11}} = \frac{\partial}{\partial O_{11}} (W_{11}^2 O_{11} + W_{21}^2 O_{21} + b_{21}) = W_{11}^2 \quad\text{———}\;\text{⑨}$$

$$\frac{\partial O_{11}}{\partial W_{11}'} = \frac{\partial}{\partial W_{11}'} (iq\, W_{11}' + cgpa\, W_{21}' + b_{11})$$

$$= iq\,(x_{i1}) \quad\text{———}\;\text{⑩}$$

use ⑨ ⑩ & ② in ⑧

$$\boxed{\frac{\partial L}{\partial W_{11}'} = -2(y - \hat{y})\, W_{11}^2\, x_{i1}}$$

similarly,

$$\frac{\partial L}{\partial W_{21}^1} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial W_{21}^1}$$

$$\boxed{\frac{\partial L}{\partial W_{21}^1} = -2\left(y - \hat{y}\right) W_{11}^2 X_{12}}$$

$$\frac{\partial L}{\partial b_{11}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial b_1}$$

$$\boxed{\frac{\partial L}{\partial b_{21}^1} \quad \frac{}{\partial b_{11}} = -2\left(y - \hat{y}\right) W_{11}^2}$$

$$\frac{\partial L}{\partial W_{12}^1} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{12}} \times \frac{\partial O_{12}}{\partial W_{12}^1}$$

$$\boxed{\frac{\partial L}{\partial W_{12}^1} = -2\left(y - \hat{y}\right) W_{21}^2 X_{11}}$$

$$\frac{\partial L}{\partial W_{22}^1} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{12}} \times \frac{\partial O_{12}}{\partial W_{22}^1}$$

$$\boxed{\frac{\partial L}{\partial W_{22}^1} = -2\left(y - \hat{y}\right) W_{21}^2 X_{12}}$$

$$\frac{\partial L}{\partial b_{12}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{12}} \times \frac{\partial O_{12}}{\partial b_{12}}$$

$$\boxed{\frac{\partial L}{\partial b_{12}} = -2(y-\hat{y})W_{21}^2}$$

After calculating all these weights, We update the weights using weight updation formula.

## * Steps in Backpropagation :-

① Initialise values of weights & Bias.

② for j in (epochs): (eg: 100, 1000)

.... for i in rang (n):

$\longrightarrow$ 2a) 1-stud. $\rightarrow$ Forward prop. $\rightarrow$ Predict
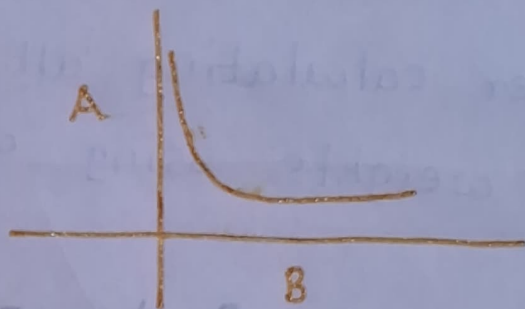
2b) Loss calculation

n
times

2c) Adjust all $\omega$. & b &

$$W_n = W_0 - n\left(\frac{\partial L}{\partial W_0}\right) \rightarrow \text{one for each traina-ble para.}$$

# * key concepts :-

## ① Derivative :-

Derivative is the method used to observe how 'A' changes when we make changes in 'B'.



suppose we write $\frac{dy}{dx}$

that means, if we twick $x$ then how $y$ changes according to $x$.

and we called it as derivative of $y$. w.r.t $x$.

and we can also say that,

y is a function of $x$. $\left(\frac{dy}{dx}\right)$

## ② Gradient :-

As like derivative gradient is just an extension of derivative,

i.e derivative term is refered when ~~the~~ ~~funct~~ something is function of single variable.

i.e $y$ is function of $x$

and derivative of y w.r.t x means,
How y changes according to x.

but,
if y is function of $x_1, x_2 \ldots x_n$
then, y is function of multiple variables, in this
case if we want to know how y changes w.r.
t. all these variables, we calculate partial
derivatives and these partial derivative are called
as gradients.

i.e $\underline{\underline{\dfrac{\partial y_i}{\partial x}}}$

So,
In Deep Learning our loss function is function
of all trainable parameters in the network
ie (all weights W & Biases b) therefore here
we calculate partial derivatives i.e gradients
using a method called Gradient Descent.

* Weight Updation Formula intution:-
We know the weight updation formula is,

$$W_{new} = W_{old} - \eta * \dfrac{\partial L}{\partial W_{old}}$$

Why (-) :-
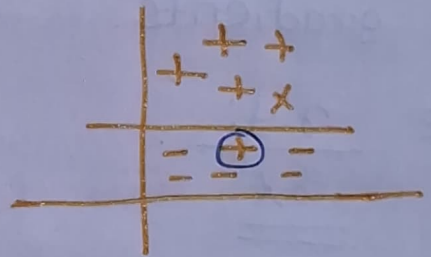We will take classification example,
Suppose the actual value is + but model
predicts -
i.e the positive point stuck in negative region

To bring back this point we have to add
something in weights.

our prediction is -tive & the value of gradi-
et also become - (because the value of gradient
is less than 0)

Therefore, it become positive
and that exactly we want.

and it is same for when model predict - but act-
ual value is +.

And for correct prediction gradient value is zero.

* What is convergence :-

We have converged when,

$$W_{new} \approx W_{old}$$

This happens because as we near the minimum, the slope
(gradient) approaches zero, making the update negligible.
We usually set a fixed number of epochs (eg.100)
rather than waiting for perfect mathematical converg-
ence.