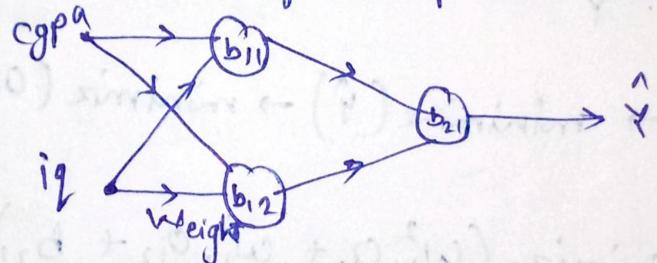


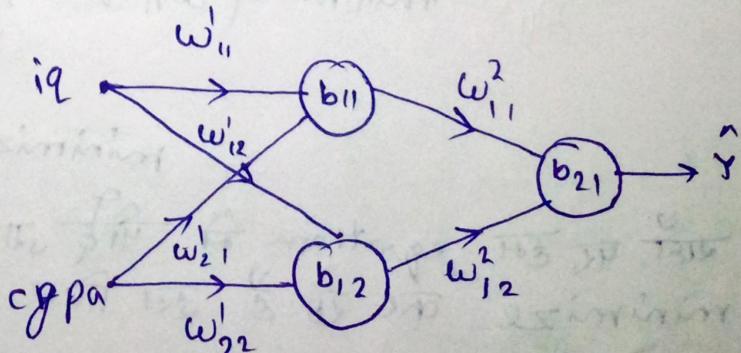
BACKPROPAGATION

- * Backpropagation, short for "backpropagation of errors", is an algorithm for supervised learning of ANN using GD.
- * Given an ANN and an error function, the method calculates the gradient of the error function with respect to the neural network's weights.
- * In simple it is used to algorithm used to train the Neural network.
- * here the term training mean. find the optimal value of the weights and bias where we got the best result of the prediction(?)



prerequisite: Gradient Descent, forward propagation.

iq	cgpa	Lpa
80	8	3
60	9	5
70	5	8
120	7	11



Here the activation function used is linear.

Total trainable parameters = 9

Backpropagation step:

1. initialize the weights (w) and bias (b), e.g. $w=1, b=0$
 2. You select the point (x_{in})
 3. predict (L_{pa}) \rightarrow forward propagation [Dot product]
 4. for adjust the weight and bias
choose a loss function example (MSE)
using GD, $L = (Y - \hat{Y})^2$
till convergence $L = (3 - 18)^2 = 225$ (error)
- go in case
also take
random

In the above we see that the error is high so that we have to minimize the error. It is only possible in MSE if we set or set the \hat{Y} .

$$\text{minimize}(L) \rightarrow \text{minimize}(\hat{Y}) \rightarrow \text{minimize}(O_{21})$$

$$O_{21} = \text{minimize}(w_{11}^2 O_1 + w_{21}^2 O_{12} + b_{21})$$

$$O_{11} = \text{minimize}(w_{11}^1 i_1 + w_{21}^1 e g p a + b_{11})$$

$$O_{12} = \text{minimize}(w_{12}^1 i_2 + w_{22}^1 e g p a + b_{12})$$

जहाँ पर हम equation में ~~पीछे~~ जो ~~कहा~~ values की minimize करते हैं उस की back propagation error.

तो क्या है?

η = learning rate

for weight update

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

for bias update

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$

for calculating $\frac{\partial L}{\partial w_{\text{old}}}$

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_{\text{old}}} \rightarrow \text{chain rule}$$

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

$$\frac{\partial \hat{y}}{\partial w_{\text{old}}} = \frac{\partial}{\partial w_{\text{old}}} [o_{11}w_{11}^2 + o_{12}w_{21}^2 + b_{21}]$$

$$= o_{11} \quad \text{Here } \hat{y} \approx o_{21}$$

$$\frac{\partial L}{\partial w_{11}^2} = -2(y - \hat{y})o_{11}$$

Hence

similarly we can calculate the other one. in the above example.

$$\frac{\partial L}{\partial w_{11}^2}, \frac{\partial L}{\partial w_{21}^2}, \frac{\partial L}{\partial b_{21}} \left| \frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{21}}, \frac{\partial L}{\partial b_{11}} \right| \frac{\partial L}{\partial w_{12}^2}, \frac{\partial L}{\partial w_{22}^2}, \frac{\partial L}{\partial b_{12}}$$

$$\frac{\partial L}{\partial w_{11}^2} = o_{11}; \frac{\partial L}{\partial w_{21}^2} = o_{12}; \frac{\partial \hat{y}}{\partial b_{21}} = -2(y - \hat{y})$$

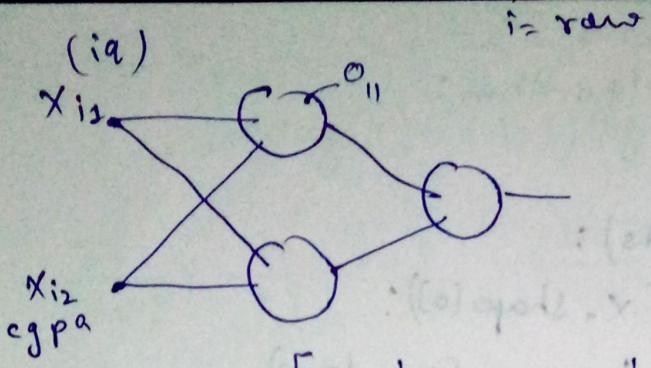
$$\left[\begin{array}{l} \frac{\partial L}{\partial w_{11}^2} = -2(Y - \hat{Y}) O_{11} \\ \frac{\partial L}{\partial w_{21}^2} = -2(Y - \hat{Y}) O_{21} \\ \frac{\partial L}{\partial b_{21}} = -2(Y - \hat{Y}) \end{array} \right] \quad (i)$$

$$\left[\begin{array}{l} \frac{\partial L}{\partial w_{11}'} = \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_{11}'} \\ \frac{\partial L}{\partial w_{21}'} = \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial w_{21}'} \\ \frac{\partial L}{\partial b_{11}} = \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial b_{11}} \end{array} \right] \quad (ii)$$

$$\left[\begin{array}{l} \frac{\partial L}{\partial w_{12}'} = \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial O_{12}} \times \frac{\partial O_{12}}{\partial w_{12}'} \\ \frac{\partial L}{\partial w_{22}'} = \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial w_{22}'} \\ \frac{\partial L}{\partial b_{12}} = \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial O_{12}} \times \frac{\partial O_{12}}{\partial b_{12}} \end{array} \right] \quad (iii)$$

where $\frac{\partial Y}{\partial O_{11}} = \frac{\partial}{\partial O_{11}} [w_{11}^2 O_{11} + w_{21}^2 O_{21} + b_{21}] = w_{11}^2$

$$\frac{\partial Y}{\partial O_{12}} = \frac{\partial}{\partial O_{12}} [w_{11}^2 O_{11} + w_{21}^2 O_{21} + b_{21}] = w_{21}^2$$



$$\left\{ \begin{array}{l} \frac{\partial O_{11}}{\partial W_{11}'} = \frac{\partial [iq w_{11}' + cgpa w_{21}' + b_{11}]}{\partial W_{11}'} = iq = x_{i1} \\ \frac{\partial O_{11}}{\partial W_{21}'} = x_{i2} \quad \frac{\partial O_{11}}{\partial b_{11}} = 1 \end{array} \right.$$

$$\left\{ \begin{array}{l} \frac{\partial O_{12}}{\partial W_{12}'} = \frac{\partial}{\partial W_{12}'} [iq w_{12}' + cgpa w_{22}' + b_{12}] = iq = x_{i2} \\ \frac{\partial O_{12}}{\partial W_{22}'} = x_{i2} \quad \frac{\partial O_{12}}{\partial b_{12}} = 1 \end{array} \right.$$

Now put all these values in equations (i), (ii) and (iii)
Hence our nine derivatives are

$$\frac{\partial L}{\partial W_{11}'} = -2(\hat{Y} - Y) W_{11}^2 X_{i1}$$

$$\frac{\partial L}{\partial W_{21}'} = -2(\hat{Y} - Y) W_{11}^2 X_{i2}$$

$$\frac{\partial L}{\partial b_{11}} = -2(\hat{Y} - Y) W_{11}^2$$

$$\left| \begin{array}{l} \frac{\partial L}{\partial W_{12}'} = -2(\hat{Y} - Y) W_{21}^2 X_{i1} \\ \frac{\partial L}{\partial W_{22}'} = -2(\hat{Y} - Y) W_{21}^2 X_{i2} \\ \frac{\partial L}{\partial b_{12}} = -2(\hat{Y} - Y) W_{21}^2 \end{array} \right.$$

$$\frac{\partial L}{\partial W_{11}''} = -2(\hat{Y} - Y) O_{11}$$

$$\frac{\partial L}{\partial W_{21}''} = -2(\hat{Y} - Y) O_{12}$$

$$\frac{\partial L}{\partial b_{21}} = -2(\hat{Y} - Y)$$

Back propagation Algorithm:

epoch=5

for i in range(epochs):

 for j in range(x, shape[0]):

 → select 1 row (random)

 → predict (using forward loop)

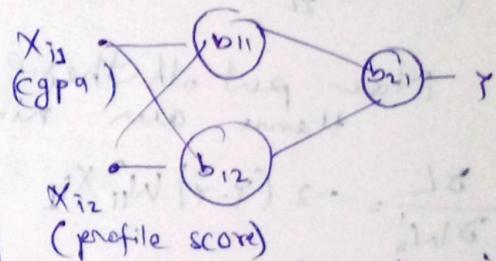
 → Calculate loss (using loss function → MSE)

 → update weight and bias using GD

$$\text{new } w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

calculate avg loss for the epoch

cgpa	profile score placement
8	8
7	9
6	10
5	5

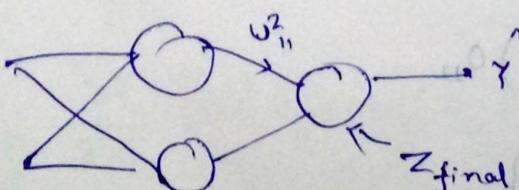


Activation function = sigmoid
trainable parameters = 9

$$\text{Loss} = -y \log(y) - (1-y) \log(1-y)$$

$$z = 0.1 \times 8 + 0.1 \times 8 + 0$$

$$\sigma(z) = 0.1$$



$$L \rightarrow \hat{Y} \rightarrow z_{\text{final}} \rightarrow z = \omega_{11}^2 \sigma_{11} + \omega_{21}^2 \sigma_{12} + b_{21}$$

$$L = -\gamma \log(\hat{Y}) - (1-\gamma) \log(1-\hat{Y})$$

$$\frac{\partial L}{\partial \omega_{11}^2} = \frac{\partial L}{\partial \hat{Y}} \times \frac{\partial \hat{Y}}{\partial z} \times \frac{\partial z}{\partial \omega_{11}^2} \left. \right\} = -(\hat{Y} - \gamma) \sigma_{11}$$

$$\frac{\partial L}{\partial \omega_{21}^2} = \frac{\partial L}{\partial \hat{Y}} \times \frac{\partial \hat{Y}}{\partial z} \times \frac{\partial z}{\partial \omega_{21}^2} \left. \right\} = -(1-\hat{Y}) \sigma_{12}$$

$$\frac{\partial L}{\partial b_{21}} = \frac{\partial L}{\partial \hat{Y}} \times \frac{\partial \hat{Y}}{\partial z} \times \frac{\partial z}{\partial b_{21}} \left. \right\} = -(\hat{Y} - \gamma)$$

Now $\frac{\partial L}{\partial \hat{Y}} = \frac{\partial [-\gamma \log(\hat{Y}) - (1-\gamma) \log(1-\hat{Y})]}{\partial \hat{Y}}$

$$= -\frac{\gamma}{\hat{Y}} + \frac{(1-\gamma)}{1-\hat{Y}} = \frac{-\gamma(1-\hat{Y}) + \hat{Y}(1-\gamma)}{\hat{Y}(1-\hat{Y})}$$

$$\frac{\partial L}{\partial \hat{Y}} = -\frac{(\hat{Y} - \gamma)}{\hat{Y}(1-\hat{Y})} \quad \text{--- (i)}$$

$$\frac{\partial \hat{Y}}{\partial z} = \frac{\partial (\sigma(z))}{\partial z} = \sigma(z)[1 - \sigma(z)] = \hat{Y}(1-\hat{Y}) \quad \text{--- (ii)}$$

$$\frac{\partial L}{\partial Y} \times \frac{\partial \hat{Y}}{\partial z} = -(\hat{Y} - \gamma)$$

~~similarly:~~

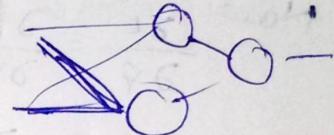
similarly:

* $L \rightarrow \hat{Y} \rightarrow z_f \rightarrow o_{11} \rightarrow z_{\text{pre}} \rightarrow w'_{11}$

$$\frac{\partial L}{\partial w'_{11}} = \underbrace{\frac{\partial L}{\partial \hat{Y}} \times \frac{\partial \hat{Y}}{\partial z_f}}_{-(\hat{Y} - \hat{Y})} \times \underbrace{\frac{\partial z_f}{\partial o_{11}}}_{w'^2_{11}} \times \underbrace{\frac{\partial o_{11}}{\partial z_{\text{pre}}}}_{o_{11}(1-o_{11})} \times \underbrace{\frac{\partial z_{\text{pre}}}{\partial w'_{11}}}_{x_{i1}}$$

$$\frac{\partial L}{\partial w'_{11}} = -(\hat{Y} - \hat{Y}) \cdot w'^2_{11} \cdot o_{11}(1-o_{11}) x_{i1}$$

$$\frac{\partial L}{\partial w'_{12}} = -(\hat{Y} - \hat{Y}) \cdot w'^2_{11} \cdot o_{11}(1-o_{11}) x_{i2}$$

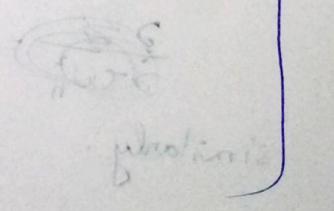
$$\frac{\partial L}{\partial b_{11}} = -(\hat{Y} - \hat{Y}) \cdot w'^2_{11} \cdot o_{11}(1-o_{11})$$


* $L \rightarrow \hat{Y} \rightarrow z_f \rightarrow o_{12} \rightarrow w'_{12} \rightarrow z_p$

$$\frac{\partial L}{\partial w'_{12}} = \underbrace{\frac{\partial L}{\partial \hat{Y}} \times \frac{\partial \hat{Y}}{\partial z_f}}_{-(\hat{Y} - \hat{Y})} \times \underbrace{\frac{\partial z_f}{\partial o_{12}}}_{z_f = w'^2_{11} o_{11} + w'^2_{12} o_{12} + b_{11}} \times \underbrace{\frac{\partial o_{12}}{\partial z_p}}_{o_{12}(1-o_{12})} \times \underbrace{\frac{\partial z_p}{\partial w'_{12}}}_{x_{i1}}$$

$$\frac{\partial L}{\partial w'_{12}} = -(\hat{Y} - \hat{Y}) \cdot w'^2_{11} o_{12} (1-o_{12}) x_{i1}$$

$$\frac{\partial L}{\partial w'_{12}} = -(\hat{Y} - \hat{Y}) \cdot w'^2_{12} o_{12} (1-o_{12}) x_{i2}$$

$$\frac{\partial L}{\partial b_{12}} = -(\hat{Y} - \hat{Y}) \cdot w'^2_{12} o_{12} (1-o_{12})$$


→ loss function is a function of all trainable parameters

→ for example Backpropagation Regression.

$$\text{Loss function} = (Y - \hat{Y})^2 = \text{MSE}$$

$$\text{OR } L(\hat{Y})$$

because $\gamma = \text{const.}$

$$\hat{Y} = W_{11}^2 O_{11} + W_{12}^2 O_{12} + b_{12}$$

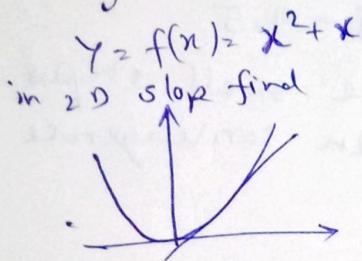
$$\hat{Y} = W_{11}^2 [W_{11}' X_{11} + W_{12}' X_{12} + b_{11}] + W_{12}^2 [W_{12}' X_{11} + W_{22}' X_{12} + b_{12}] + b_{12}$$

0	0	0
0	0	0
0	0	0

tune the value for find minimum loss

Gradient

derivative function depends only on one variable



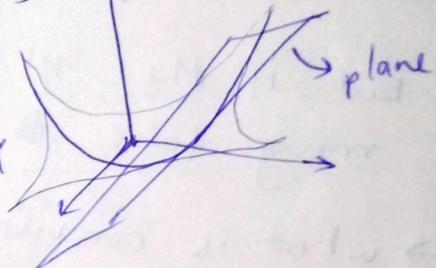
$$z = f(x, y) = x^2 + y^2$$

slope find with each dim

in 3-D

$$\frac{\partial z}{\partial x} = 2x$$

$$\frac{\partial z}{\partial y} = 2y$$



* concept of minima $\frac{\partial y}{\partial x} = 0$

at that point of x minima is lies

example. $z = x^2 + y^2$
 $(x, y) \rightarrow \min(z)$

$$\frac{\partial z}{\partial x} = 0 \Rightarrow 2x = 0 \Rightarrow x = 0$$

$$\frac{\partial z}{\partial y} = 0 \Rightarrow 2y = 0 \Rightarrow y = 0$$

If $(x, y) = (0, 0)$ then z is minimum.

→ effect of Learning Rate(η):

$$\ln \ln = k_0 - \frac{m^2 L}{M^2}$$

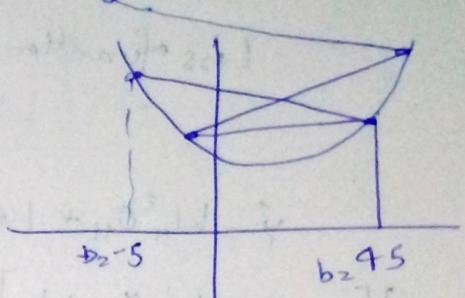
learning rate steps की smooth तरीका

$$\text{let } b = -5$$

$$\text{let } \frac{\partial L}{\partial b} = -50$$

update 'b'

$$b = -5 - (-5)$$
$$= 45$$



- for control this type of zigzag motion.
use learning rate (η)

$$\text{let } b = -5 \quad \frac{\partial L}{\partial b} = -5_0 \quad m = 0.01$$

update 'b'

$$b = -5 + (0.01 \times 50) = 5 - 4.5$$

Here we take small steps

— but if the ' η ' is very small then convergence may take a long time.

→ what is convergence?

(i) if $\text{isnew} = \text{World}$

- (i) if b knew = world
- (ii) if run for loop in 10^6 or 10^{10} times

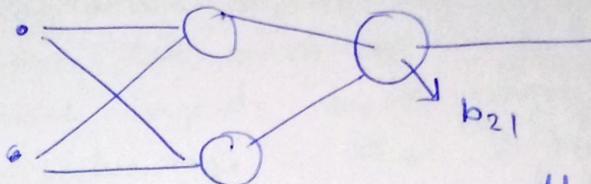
Backpropagation. Intuition:

- Now we understand

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

- Here learning rate we assume $\eta = 1$

$$\therefore w_{\text{new}} = w_{\text{old}} - \frac{\partial L}{\partial w}$$

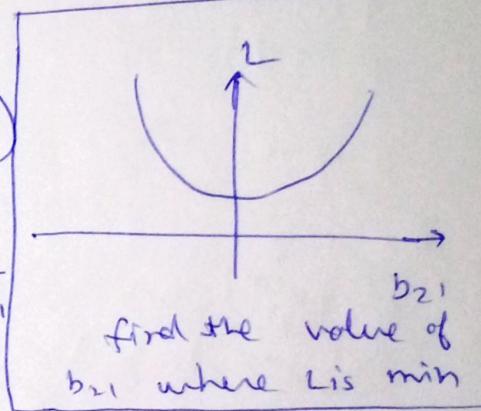


Now here less vs assume all are constant except b_{21}

$$L = (y - \hat{y})^2 \Rightarrow L(b_{21})$$

Now update b_{21}

$$b_{21} = b_{21} - \frac{\partial L}{\partial b_{21}}$$



, L derivative w.r.t b_{21}

if $\frac{\partial L}{\partial b_{21}} = +ve$. $b_{21} \downarrow L \downarrow$

if $\frac{\partial L}{\partial b_{21}} = -ve$ $b_{21} \uparrow L \downarrow$

