

Backward Propagation

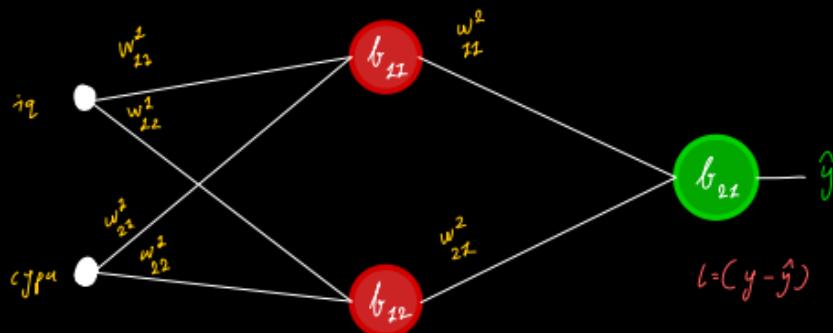
→ Backpropagation "Backward Propagation of errors", is an algorithm for supervised learning of artificial neural networks using gradient descent. Given an artificial neural network & an error function, the method calculates the gradient of error function with respect to the neural network's weights.

in simple

↳ An algorithm to train neural-networks.

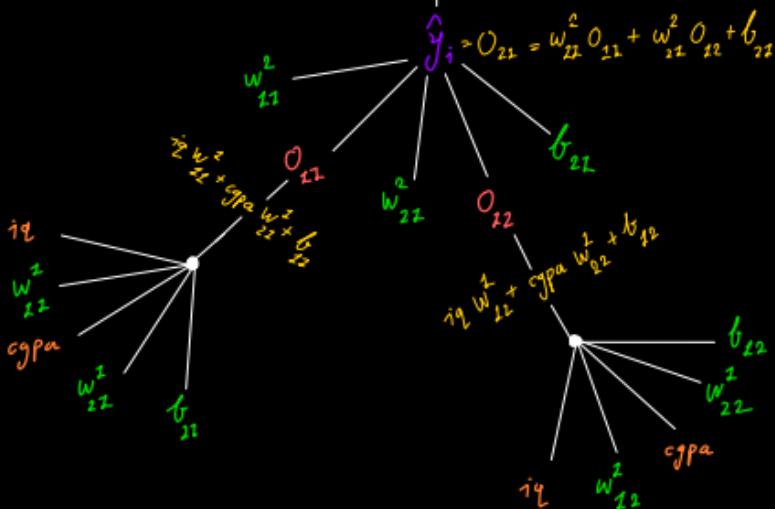
Steps

1. Initialize weights & biases.
2. Select a data row & feed it to network.
3. Make prediction.
4. Select loss function & calculate loss.
5. Adjust weights & biases based on calculated error.



Let's see how all these parameters depends for each other

$$\text{Loss} = (y_i - \hat{y}_i)$$



- As above dependency tree shows, in order to change loss we need to change \hat{y}_i
- To change \hat{y}_i change $w_{21}^2, w_{22}^2, O_{21}, O_{22}$ and b_{21} .
- Since w_{21}^2, w_{22}^2 and b_{21} are independent, we can change the value of those weight
- In order to change O_{21} , we need to update $iq, w_{21}^2, cgp, w_{22}^2$ and b_{21}
- In order to change O_{22} , we need to update $iq, w_{21}^2, cgp, w_{22}^2$ and b_{22} .
- Since iq and cgp are actual input delta, we won't update it hence remained parameters $w_{21}^2, w_{22}^2, w_{21}^2, w_{22}^2, b_{21}$ and b_{22} will get update.



formulas for weights / bias updation.

for weights

for bias

$$W_{new} = W_{old} - \alpha \frac{\partial L}{\partial W_{old}}$$

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

Total number of derivatives, we need to find....

$$\frac{\partial L}{\partial W_{22}} \quad \frac{\partial L}{\partial W_{22}^2} \quad \frac{\partial L}{\partial b_{22}} \quad \left| \begin{array}{ccc} \frac{\partial L}{\partial W_{22}} & \frac{\partial L}{\partial W_{22}^2} & \frac{\partial L}{\partial b_{22}} \\ \end{array} \right| \quad \left| \begin{array}{ccc} \frac{\partial L}{\partial W_{22}} & \frac{\partial L}{\partial W_{22}^2} & \frac{\partial L}{\partial b_{22}} \\ \end{array} \right|$$

derivative $\frac{dy}{dx}$ says that how 'y' is changing with respect to 'x'.

$\frac{\partial L}{\partial W_{22}^2}$: How Loss is changing with respect to W_{22}^2 .

$$\begin{array}{c} w_{22}^2 \\ \downarrow \\ \hat{y} \\ \downarrow \\ \text{Loss} \end{array} \quad \frac{\partial L}{\partial W_{22}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_{22}^2}$$

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

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

$$\frac{\partial \hat{y}}{\partial W_{22}^2} = \frac{\partial}{\partial W_{22}^2} (O_{22} w_{22}^2 + O_{22} w_{22}^2 + b_{22})$$

$$\frac{\partial \hat{y}}{\partial W_{22}^2} = \boxed{O_{22}}$$

$$\frac{\partial L}{\partial W_{22}^2} = -2(y - \hat{y}) O_{22}$$



$$\frac{\partial L}{\partial w_{22}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{22}^2}$$

$$\frac{\partial L}{\partial \hat{y}} = \boxed{-2(y - \hat{y})} \quad \frac{\partial \hat{y}}{\partial w_{22}^2} = \frac{\partial}{\partial w_{22}^2} (O_{22} w_{22}^2 + O_{22} w_{22}^2 + b_{22})$$

$$\frac{\partial \hat{y}}{\partial w_{22}^2} = \boxed{O_{22}}$$

$$\boxed{\frac{\partial L}{\partial w_{22}^2} = -2(y - \hat{y}) O_{22}}$$

$$\frac{\partial L}{\partial b_{22}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b_{22}}$$

$$\frac{\partial L}{\partial \hat{y}} = \boxed{-2(y - \hat{y})} \quad \frac{\partial \hat{y}}{\partial b_{22}} = \frac{\partial}{\partial b_{22}} (O_{22} w_{22}^2 + O_{22} w_{22}^2 + b_{22})$$

$$\frac{\partial \hat{y}}{\partial b_{22}} = \boxed{1}$$

$$\boxed{\frac{\partial L}{\partial b_{22}} = -2(y - \hat{y})}$$

Now we will calculate next derivatives.

$$\frac{\partial L}{\partial w_{11}^2} \quad \frac{\partial L}{\partial w_{21}^2} \quad \frac{\partial L}{\partial b_{11}}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial O_{11}} \quad \frac{\partial O_{11}}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{22}^2} = \frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial O_{22}} \quad \frac{\partial O_{22}}{\partial w_{22}^2}$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial O_{21}} \quad \frac{\partial O_{21}}{\partial w_{21}^2}$$

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

$$\frac{\partial L}{\partial b_{22}} = \frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial O_{22}} \quad \frac{\partial O_{22}}{\partial b_{22}}$$

$$\frac{\partial L}{\partial b_{21}} = \frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial O_{21}} \quad \frac{\partial O_{21}}{\partial b_{21}}$$

already known

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

Not known yet...

$$\frac{\partial \hat{y}}{\partial O_{11}} \quad \frac{\partial \hat{y}}{\partial O_{22}}$$

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

$$\frac{\partial \hat{y}}{\partial O_{11}} = \boxed{w_{11}^2}$$

$$\frac{\partial \hat{y}}{\partial O_{22}} = \frac{\partial}{\partial O_{22}} [O_{11} w_{11}^2 + O_{12} w_{21}^2 + b_{11}]$$

$$\frac{\partial \hat{y}}{\partial O_{12}} = \boxed{w_{21}^2}$$



$$\frac{\partial O_{11}}{\partial w'_{11}} = \frac{\partial}{\partial w'_{11}} [iqw'_{11} + cgpa w'_{21} + b_{11}]$$

$$\frac{\partial O_{11}}{\partial w'_{11}} = \boxed{iq} (x_{i1})$$

$$\frac{\partial O_{11}}{\partial w'_{21}} = \frac{\partial}{\partial w'_{21}} [iqw'_{11} + cgpa w'_{21} + b_{11}] \quad \frac{\partial O_{11}}{\partial b_{11}} = \frac{\partial}{\partial b_{11}} [iqw'_{11} + cgpa w'_{21} + b_{11}]$$

$$\frac{\partial O_{11}}{\partial w'_{21}} = \boxed{cgpa} (x_{i1}) \quad \frac{\partial O_{11}}{\partial b_{11}} = \boxed{1}$$

$$\frac{\partial O_{12}}{\partial w'_{12}} = \frac{\partial}{\partial w'_{12}} [iqw'_{12} + cgpa w'_{22} + b_{12}]$$

$$\frac{\partial O_{12}}{\partial w'_{12}} = \boxed{iq} (x_{i1})$$

$$\frac{\partial O_{12}}{\partial w'_{22}} = \frac{\partial}{\partial w'_{22}} [iqw'_{12} + cgpa w'_{22} + b_{12}] \quad \frac{\partial O_{12}}{\partial w'_{22}} = \boxed{cgpa} (x_{i1})$$

$$\frac{\partial O_{12}}{\partial b_{12}} = \frac{\partial}{\partial b_{12}} [iqw'_{12} + cgpa w'_{22} + b_{12}]$$

$$\frac{\partial O_{12}}{\partial b_{12}} = \boxed{1}$$

Put them all together

$$L. \quad \frac{\partial L}{\partial w''_{12}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{12}} \frac{\partial O_{12}}{\partial w'_{22}}$$

$$= -2(y - \hat{y}) w''_{12} x_{i2}$$

$$2. \frac{\partial L}{\partial w_{22}^z} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{22}} \frac{\partial O_{22}}{\partial w_{22}^z}$$

$$-2(y - \hat{y}) w_{11}^2 x_{i2}$$

$$3. \frac{\partial L}{\partial b_{22}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{22}} \frac{\partial O_{22}}{\partial b_{22}}$$

$$-2(y - \hat{y}) w_{11}^2$$

$$4. \frac{\partial L}{\partial w_{22}^z} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{22}} \frac{\partial O_{22}}{\partial w_{12}^z}$$

$$-2(y - \hat{y}) w_{21}^2 x_{i2}$$

$$5. \frac{\partial L}{\partial w_{22}^z} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{22}} \frac{\partial O_{22}}{\partial w_{22}^z}$$

$$-2(y - \hat{y}) w_{21}^2 x_{i2}$$

$$\frac{\partial L}{\partial b_{22}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{22}} \frac{\partial O_{22}}{\partial b_{22}}$$

$$-2(y - \hat{y}) w_{21}^2$$

Now everything works together
 $\text{epoch} = 10$

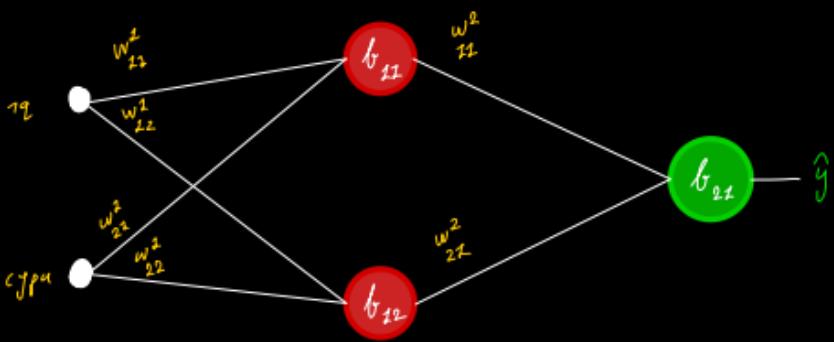
for i in range(epoch):

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

1. Select i row randomly
2. forward propagation & prediction
3. Calculate loss
4. Update weights & biases using GD

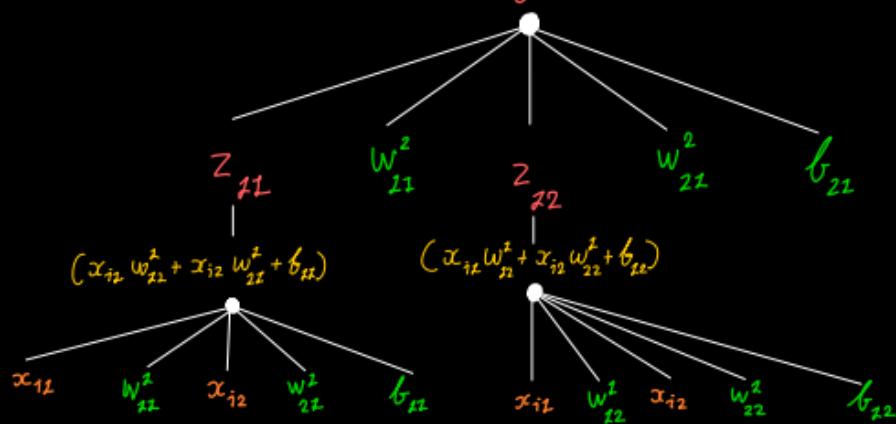
Calculate average loss for each epochs

Let's see the same thing in classification task.



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

$$\hat{y} = \sigma(O_{12} w_{11}^1 + O_{12} w_{12}^1 + b_{21})$$



$$\frac{\partial L}{\partial w_{zz}^2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{zz}^2} \quad \frac{\partial L}{\partial w_{zz}^2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{zz}^2}$$

$$\frac{\partial L}{\partial b_{zz}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial b_{zz}}$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} \left[-y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \right]$$

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

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

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

$$\frac{\partial L}{\partial \hat{y}} = \boxed{-\frac{(y - \hat{y})}{\hat{y}(1-\hat{y})}}$$

$$\frac{\partial \hat{y}}{\partial z} = \frac{\partial}{\partial z} (\sigma(z))$$

$$= \sigma(z) [z - \sigma(z)]$$

$\downarrow \hat{y} \leftarrow$

$$\frac{\partial \hat{y}}{\partial z} = \boxed{\hat{y}(z - \hat{y})}$$

$$\frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} = -\frac{(y - \hat{y})}{\hat{y}(1 - \hat{y})} \quad \hat{y}(1 - \hat{y})$$

$$= \boxed{-(y - \hat{y})}$$

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

$$= \boxed{O_{11}}$$

$$= \boxed{O_{12}}$$

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

$$= \boxed{1}$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{21}^2}$$

$$= \boxed{-(y - \hat{y}) O_{11}}$$

$$\frac{\partial L}{\partial w_{22}^2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{22}^2}$$

$$= \boxed{-(y - \hat{y}) O_{12}}$$

$$\frac{\partial L}{\partial b_{22}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial b_{22}}$$

$$= \boxed{-(y - \hat{y})}$$

$$\frac{\partial L}{\partial w_{11}^2}, \quad \frac{\partial L}{\partial w_{22}^2}, \quad \frac{\partial L}{\partial b_{22}}$$

$$z_{fin} = O_{fin}w_{12}^2 + O_{fin}w_{22}^2 + b_{22}$$

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_{fin}} \cdot \frac{\partial z_{fin}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}^1}$$

$$O_{11} = \sigma(z_{prev})$$

$$\frac{\partial Z_{fin}}{\partial O_{11}} = \boxed{w_{11}^2} \quad \frac{\partial O_{11}}{\partial Z_{phenv}} = \boxed{O_{11}(1 - O_{11})}$$

$$\frac{\partial Z_{phenv}}{\partial w_{11}^i} = \frac{\partial}{\partial w_{11}^i} [w_{11}^i x_{i1} + w_{21}^i x_{i2} + b_{11}]$$

$$= \boxed{x_{i2}}$$

$$\frac{\partial L}{\partial w_{11}^i} = \boxed{-(y - \hat{y}) w_{11}^2 O_{11} (1 - O_{11}) x_{i2}}$$

$$\frac{\partial L}{\partial w_{21}^i} = \boxed{-(y - \hat{y}) w_{11}^2 O_{11} (1 - O_{11}) x_{i2}}$$

$$\frac{\partial L}{\partial w_{22}^2} \quad \frac{\partial L}{\partial w_{22}^i} \quad \frac{\partial L}{\partial b_{12}}$$

$$\frac{\partial L}{\partial w_{22}^i} = \frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial Z_{fin}} \quad \frac{\partial Z_{fin}}{\partial O_{12}} \quad \frac{\partial O_{12}}{\partial Z_{phenv}} \quad \frac{\partial Z_{phenv}}{\partial w_{12}^i}$$

$$\frac{\partial L}{\partial w_{12}^i} = \boxed{-(y - \hat{y}) w_{22}^2 O_{12} (1 - O_{12}) x_{i2}}$$

$$\frac{\partial L}{\partial w_{22}^2} = \boxed{-(y - \hat{y}) w_{22}^2 O_{12} (1 - O_{12}) x_{i2}}$$

$$\frac{\partial L}{\partial b_{12}} = \boxed{-(y - \hat{y}) w_{22}^2 O_{12} (1 - O_{12})}$$

Q. Why this works....?

$$y = f(x)$$

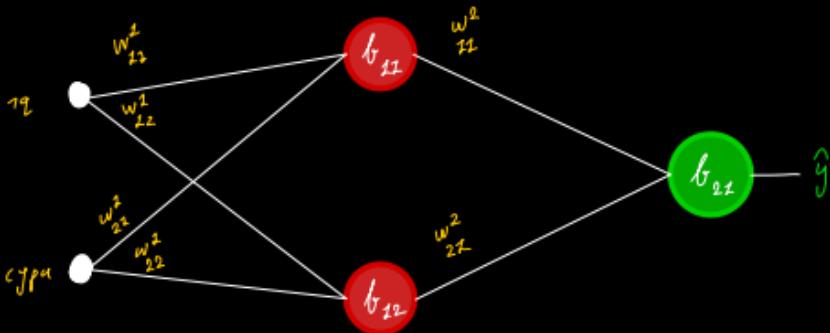
→ y is the function of x .

$$L = (y - \hat{y})$$

→ So, L is the function of whom...?

→ y is coming from data, hence it is constant

→ What is this \hat{y}?



$$\hat{y} = w_{11}^2 O_{11} + w_{12}^2 O_{12} + b_{11}$$

→ Hence w_{11}^2 and w_{12}^2 are constant, what about O_{11} and O_{12}

$$O_{11} = x_{i1} w_{11}^2 + x_{i2} w_{12}^2 + b_{11}$$

$$O_{12} = x_{i1} w_{12}^2 + x_{i2} w_{21}^2 + b_{12}$$

→ Hence x_{i1} and x_{i2} are input data-points i_1 and i_2 .

Hence...

$$\hat{y} = w_{11}^2 [x_{i1} w_{11}^2 + x_{i2} w_{12}^2 + b_{11}] + w_{12}^2 [x_{i1} w_{12}^2 + x_{i2} w_{21}^2 + b_{12}] + b_{21}$$

→ That's why L is function of \hat{y} .

$$L(\hat{y})$$



→ We can say that loss function is the function of all trainable parameters. ($w_{11}^2, w_{21}^2, w_{31}^2, w_{41}^2, b_{11}, b_{21}, b_{31}, b_{41}$)

Concept of gradient.

→ A nickname of derivative.

$y = f(x) = x^2 + x$ task: Differentiate 'y' with respect to 'x'.

$$\frac{dy}{dx} = \frac{d}{dx} f(x) = \frac{d}{dx} (x^2 + x) = \boxed{2x + 1}$$

Derivative vs Gradient

Derivative: The opf of a function depends on only 1 parameter

ex
 $y = f(x)$

$$y \rightarrow x \quad \frac{d}{dx}$$

Gradient: The opf of a function depends on >1 parameters.

ex

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

$$y \rightarrow (x, z)$$

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

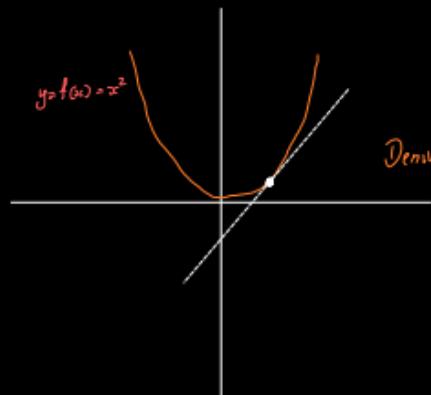
$$\frac{\partial y}{\partial z} = 2z$$

Notation = ∂

Our loss function depends on 9 parameters of network.

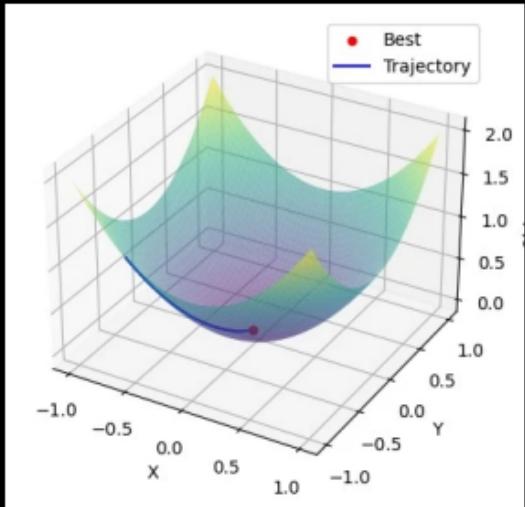
behind the scene, we will calculate gradient of each params.

in 2D



Derivative = Calculating slope at specific location.





3D :

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

Calculating slope with respect to x, y .

- Our loss function is a complex 9 dimensional function
- By calculating derivatives we're calculating 9 different slopes WRT each dimensions.

Concept of derivative ($\frac{dy}{dx}$)

- How y is changing by making changes in x .

$$y = x^2 + 2x$$

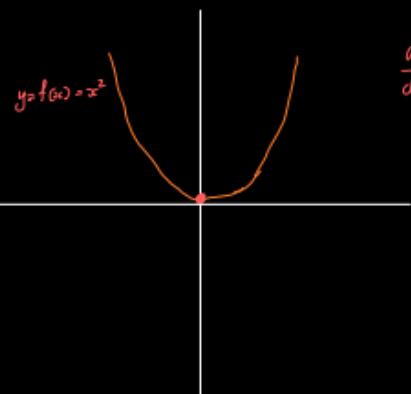
$$\text{for } x = 5$$

$$\frac{dy}{dx} = 2x + 2$$

$$\frac{dy}{dx} = 2(5) + 2$$

= 12 slope

Concept of minima



2D

$$\frac{dy}{dx} = 2x = 0$$

At $x = 0$ we are at global minima.

3D

$$\text{for } z = x^2 + y^2$$

$$\frac{\partial z}{\partial x} = 2x = 0 \quad \frac{\partial z}{\partial y} = 2y = 0$$

At $x = 0$ and $y = 0$, we're at global



Created with
Notewise

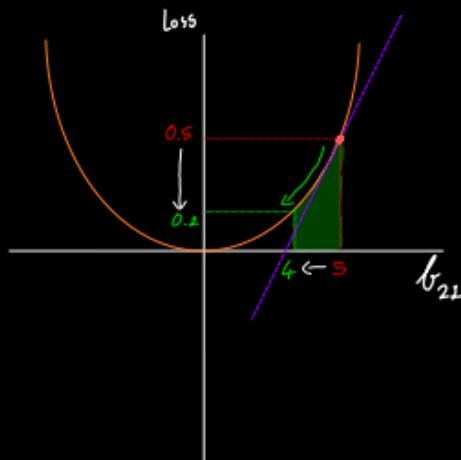
$$b_{21} = b_{21} - \frac{\partial L}{\partial b_{21}}$$

why we're subtracting $\frac{\partial L}{\partial b_{21}} \dots ?$

Suppose $b_{21} = 5$ & d/p of $\frac{\partial L}{\partial b_{21}}$ is +ve

Hence d/p of $\frac{\partial L}{\partial b_{21}} = +ve$ means run loss

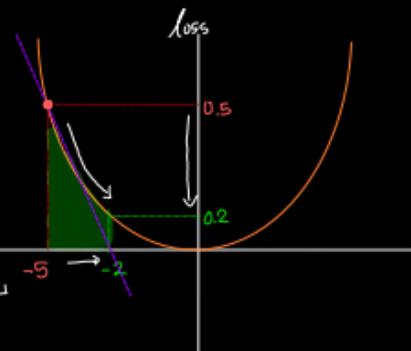
is increasing as we increase the value of b_{21} . hence we will subtract it.



\rightarrow Similar in opposite, if $\frac{\partial L}{\partial b_{21}}$ is -ve

then the value will be added in b_{21}

\rightarrow It shows that as b_{21} increases b_{21} , loss is getting negative.

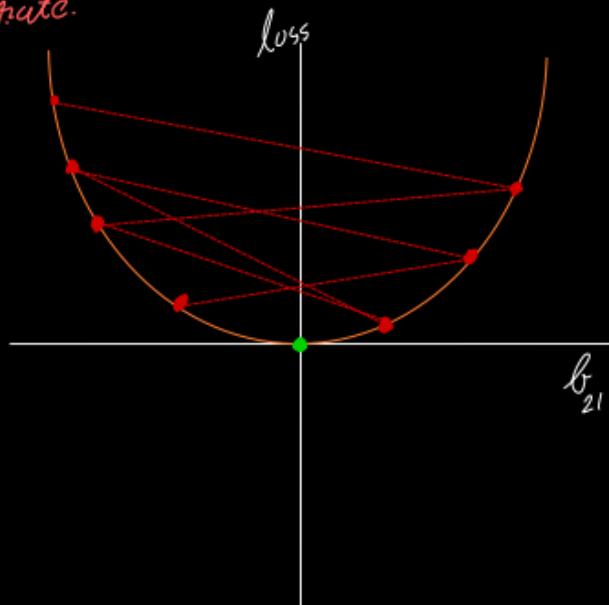


What about learning rate...?

↳ learning rate makes process smoother & makes process stable.

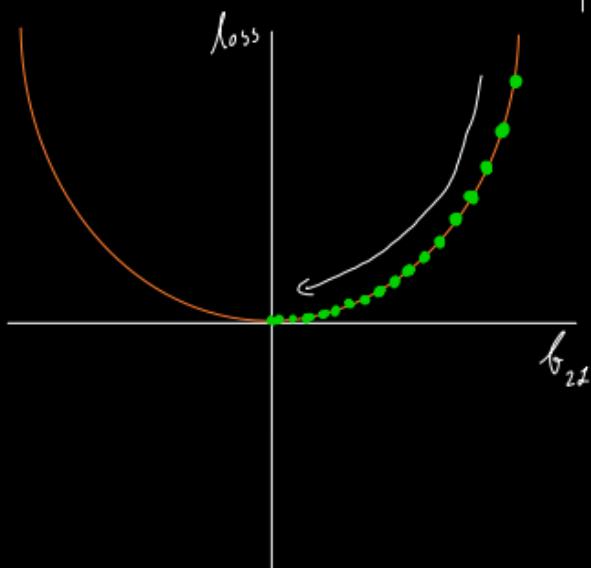
Without learning rate.

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



with learning rate

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



What is convergence...?

$$b_{\text{new}} = b_{\text{old}} - \alpha \frac{\partial L}{\partial b_{\text{old}}} = 0$$

→ $b_{\text{new}} = b_{\text{old}}$ says that now we've reached at global minima

→ We should keep updating parameters until

$$b_{\text{new}} \approx b_{\text{old}}$$

→ As you reach at minima, stop because further updation won't give any meaningful results.

