

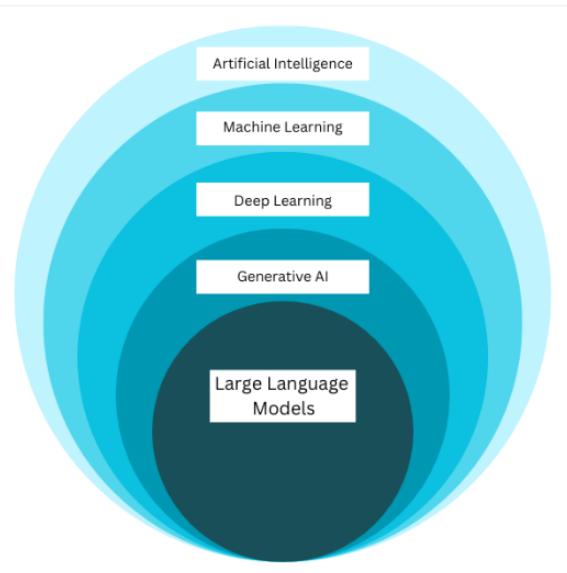
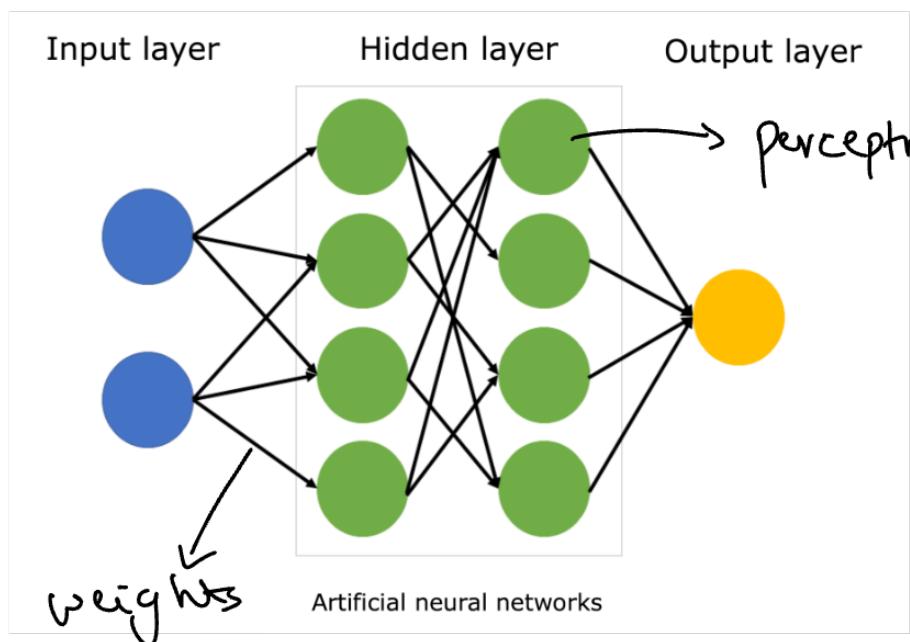
✓ A M N

Photos 

# Deep Learning

deep learning is a subfield of AI and ML that is inspired by structure of human brain.

Similar deep learning algorithm attempt to draw continually conclusions as humans would by logically analyzing data with a given structure called Neural Network.



DL is part of a broader family of ml methods based on ANN with representation learning

↓  
auto feature extraction

Deep learning algorithms uses multiple layers to progressively extract higher-level features from the raw input. for ex, in image processing, lower layers may identify edges (primitive features), while higher layers may identify concepts relevant to a human such as digits or letters or faces (complex features)

In DL, we do not need feature extraction because the architecture is designed in a way we dont need to extract feature.

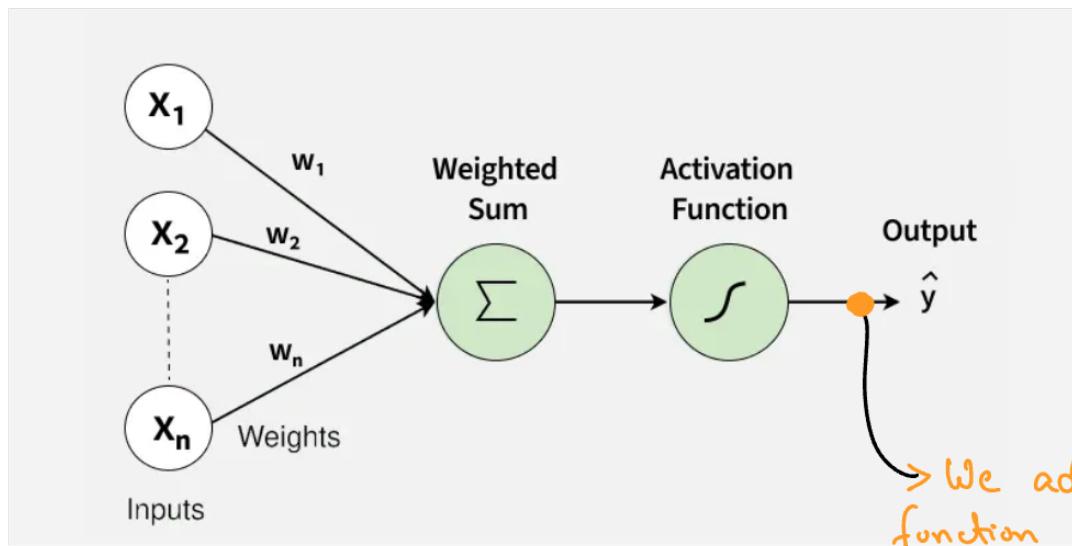
# Deep Learning vs Machine Learning

1. Data Dependency : The more data provided to DL, the more better it gets whereas in ML, after some threshold, the performance remains stegnant.
2. Hardware Dependency : GPU are required for DL model training because of a lot of complex matrix operations whereas, in ML, it's not required
3. Training Time : DL has very high training time while ML is very less
4. Feature Selection : No need of this in DL because of representative learning where DL layers automatically extract features from data whereas for ML, features are provided prior training.
5. Interpretability : DL models are generally black box models which is a downside.

# Perceptron

A perceptron is a very simple model of a neuron used for binary classification.

Simply, perceptron takes several inputs, multiplies each by a weight, adds them up together (including bias) and then passes the result through an activation function to output either 0 or 1.



We add loss function here during training.

The biggest drawback of a perceptron is that it only works on binary output which is valid for linear data.

# Perceptron trick

The simple mental model:

1. Randomly pick up a point  $x$ ,  
compute score  $z = w \cdot x + b$
  2. Check sign:  $z \geq 0 \rightarrow \text{predict } +1 \text{ else } 0$
  3. Compare it to true label  $y$ :
    - if correct then skip
    - Wrong then change the boundary (shift line)
  4. Repeat for all points until mostly correct.
- ∴ Start with a random line  
pick a point  
if misclassified point  
push line towards correct side

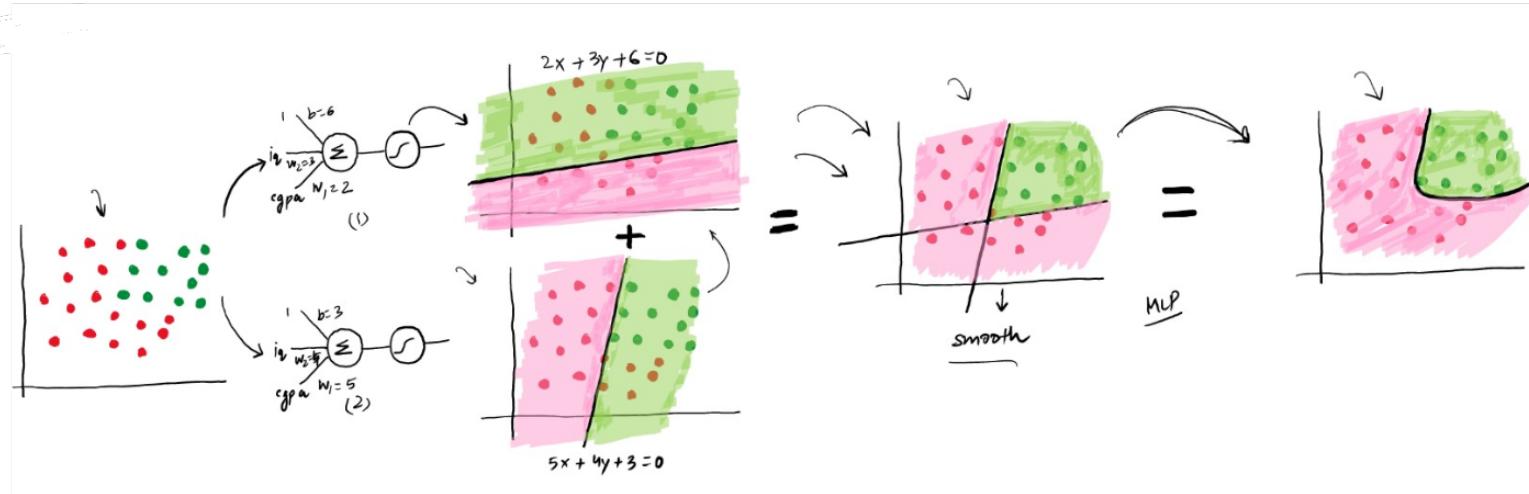
MODEL	LOSS FUNCTION	FORMULA	ACTIVATION	OUTPUT	TASK
Perceptron	Hinge Loss	$L = \max(0, 1 - y \cdot \hat{y})$	Step	{-1, +1} Hard labels	Binary Classification
Logistic Regression	Log-Loss (Binary Cross-Entropy)	$L = -[y \cdot \log(\hat{y}) + (1-y) \cdot \log(1-\hat{y})]$	Sigmoid	{0, 1} Probability → class	Binary Classification
Softmax Regression	Categorical Cross-Entropy	$L = -\sum y_k \cdot \log(\hat{y}_k)$	Softmax	[0, 1] <sup>k</sup> Probability distribution	Multiclass Classification
Linear Regression	MSE (Mean Squared Error)	$L = (1/n) \sum (y - \hat{y})^2$	Linear	$\mathbb{R}$ Continuous values	Regression

If we change the activation function or loss function or both, perceptron will behave differently.

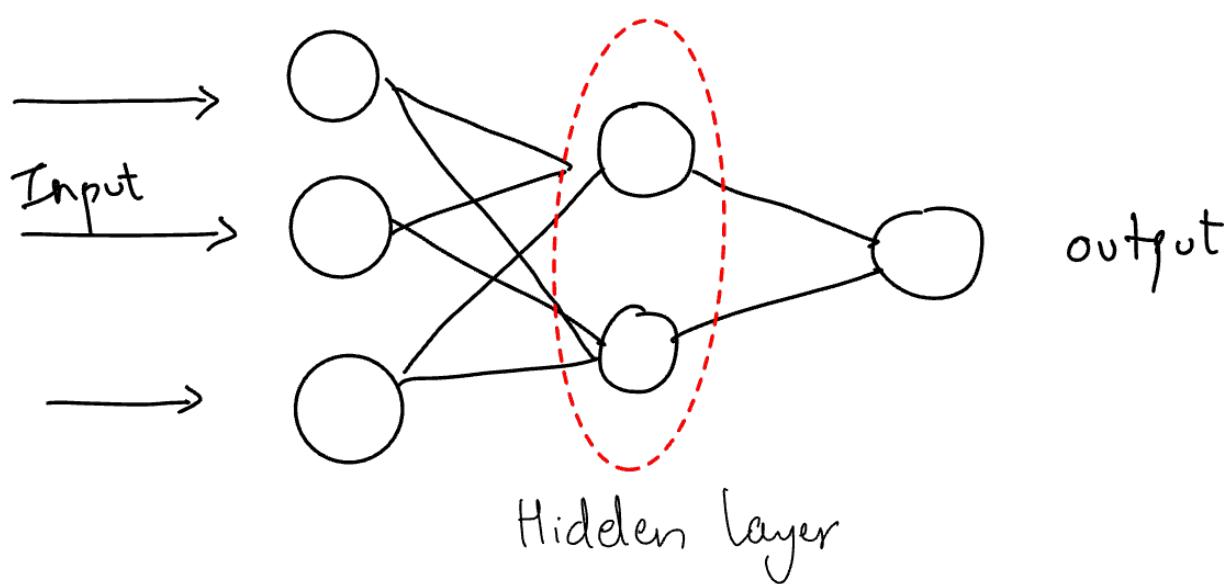
What's the problem with perceptron?  
 → It only works on linear data and not on non-linear data

# Multi-Layer Perceptron : MLP

Multi-layer perceptrons, or MLPs, are a powerful member of the Artificial Neural Networks family that can be used to solve complex problems that a single perceptron alone cannot. At their core, MLPs are complex and it has a collection of interconnected single perceptrons, also known as neurons or nodes, working together to process and analyze data.



Here as we can see, outputs of perceptrons are again fed as an input to another perceptron. This layering is called MLP



Each perceptron is called nodes when stacked in a neural network architecture.

There are four ways to improve the outputs.

1. Add more nodes to the Hidden Layer
2. Add more nodes to the input layer  
(This is when there are more input cols, more nodes are present by default)
3. Adding Nodes in Output
4. Adding More Hidden Layers  
↓  
Deep Neural Network

# MLP Notation

$x_{ii}$

Output dendition

$(w, b) \rightarrow$  notation

Denoting weights:  $w_{ij}^k$   
where  $k \rightarrow$  layer your  
weight is  
ending

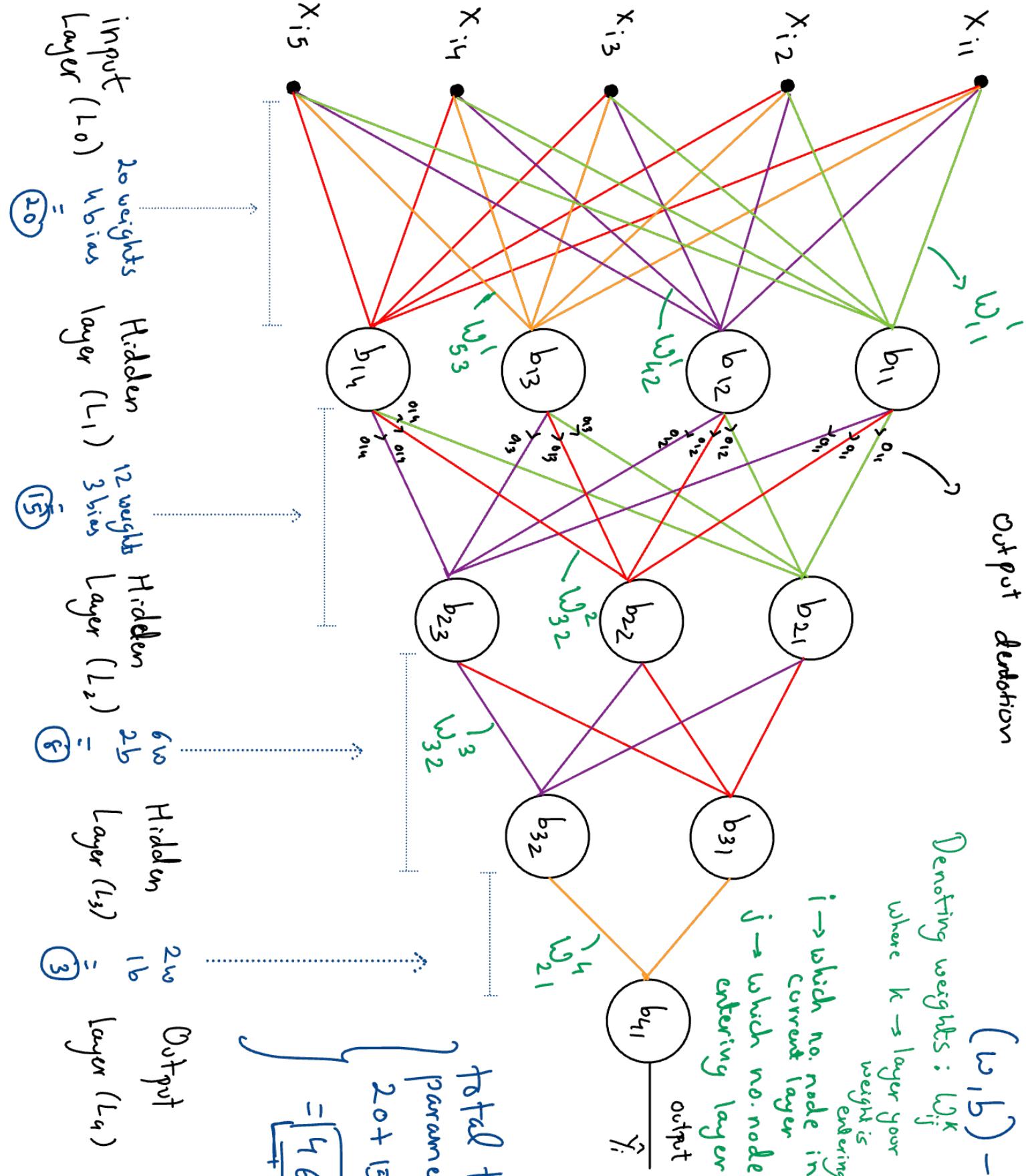
$i \rightarrow$  which no. node in  
current layer  
 $j \rightarrow$  which no. node in  
entering layer

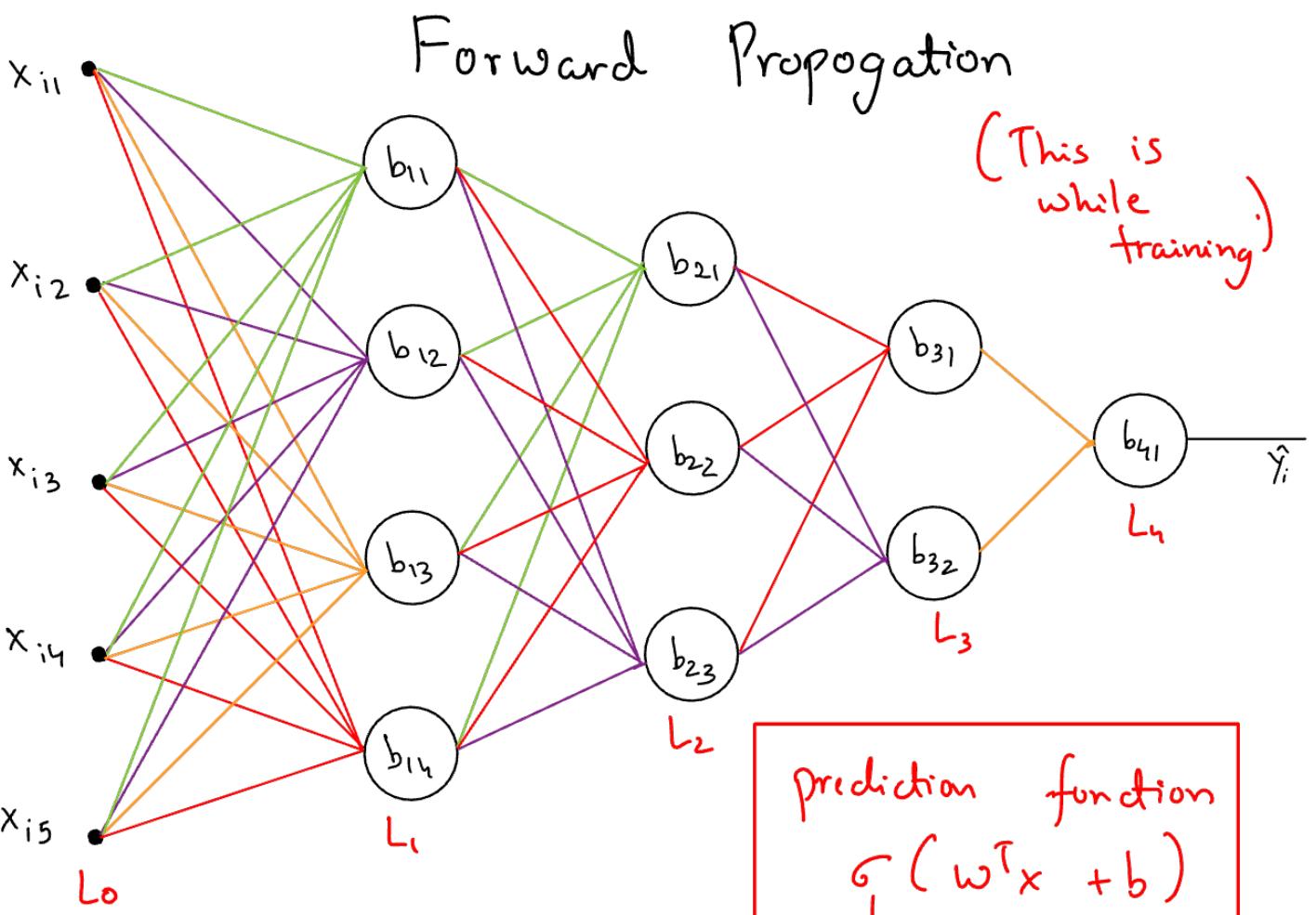
$b_{ii}$   
 $\hat{y}_i$   
output

total trainable  
parameters:

$$20 + 15 + 8 + 3$$

$$= \boxed{46}$$





(first row is passed)

Layer 1:

$$\begin{bmatrix} w_{11}^1 & w_{12}^1 & w_{13}^1 & w_{14}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 & w_{24}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 & w_{34}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 & w_{44}^1 \\ w_{51}^1 & w_{52}^1 & w_{53}^1 & w_{54}^1 \end{bmatrix}_{5 \times 4} \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \\ x_{i5} \end{bmatrix}_{5 \times 1} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \\ b_{14} \end{bmatrix}_{4 \times 1}$$

$w$

$x$

$b$

Sigmoid activation func

Transposing  $w \rightarrow w^T$

$$\begin{bmatrix} w_{11}^1 & w_{21}^1 & w_{31}^1 & w_{41}^1 & w_{51}^1 \\ w_{12}^1 & w_{22}^1 & w_{32}^1 & w_{42}^1 & w_{52}^1 \\ w_{13}^1 & w_{23}^1 & w_{33}^1 & w_{43}^1 & w_{53}^1 \\ w_{14}^1 & w_{24}^1 & w_{34}^1 & w_{44}^1 & w_{54}^1 \end{bmatrix}_{4 \times 5}$$

Multiplying  $W^T$  and  $x \rightarrow W^T x$   
 $[4 \times 3] [5 \times 1] \rightarrow [4 \times 1]$

$$\begin{bmatrix} w_{11}'x_{i1} + w_{21}'x_{i2} + w_{31}'x_{i3} + w_{41}'x_{i4} + w_{51}'x_{i5} \\ w_{12}'x_{i1} + w_{22}'x_{i2} + w_{32}'x_{i3} + w_{42}'x_{i4} + w_{52}'x_{i5} \\ w_{13}'x_{i1} + w_{23}'x_{i2} + w_{33}'x_{i3} + w_{43}'x_{i4} + w_{53}'x_{i5} \\ w_{14}'x_{i1} + w_{24}'x_{i2} + w_{34}'x_{i3} + w_{44}'x_{i4} + w_{54}'x_{i5} \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} \quad 4 \times 1 \quad 4 \times 1$$

$$\begin{bmatrix} w_{11}'x_{i1} + w_{21}'x_{i2} + w_{31}'x_{i3} + w_{41}'x_{i4} + w_{51}'x_{i5} + b_1 \\ w_{12}'x_{i1} + w_{22}'x_{i2} + w_{32}'x_{i3} + w_{42}'x_{i4} + w_{52}'x_{i5} + b_2 \\ w_{13}'x_{i1} + w_{23}'x_{i2} + w_{33}'x_{i3} + w_{43}'x_{i4} + w_{53}'x_{i5} + b_3 \\ w_{14}'x_{i1} + w_{24}'x_{i2} + w_{34}'x_{i3} + w_{44}'x_{i4} + w_{54}'x_{i5} + b_4 \end{bmatrix}$$

$w^T x + b \rightarrow$  Now apply Sigmoid  
activation function  
 $\sigma(w^T x + b)$

6

$$\begin{bmatrix} w_{11}'x_{i1} + w_{21}'x_{i2} + w_{31}'x_{i3} + w_{41}'x_{i4} + w_{51}'x_{i5} + b_1 \\ w_{12}'x_{i1} + w_{22}'x_{i2} + w_{32}'x_{i3} + w_{42}'x_{i4} + w_{52}'x_{i5} + b_2 \\ w_{13}'x_{i1} + w_{23}'x_{i2} + w_{33}'x_{i3} + w_{43}'x_{i4} + w_{53}'x_{i5} + b_3 \\ w_{14}'x_{i1} + w_{24}'x_{i2} + w_{34}'x_{i3} + w_{44}'x_{i4} + w_{54}'x_{i5} + b_4 \end{bmatrix}$$

$$\begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \\ o_{14} \end{bmatrix} \quad \text{outputs}$$

Layer 2 :

$$\begin{bmatrix} w_{11}^2 & w_{12}^2 & w_{13}^2 \\ w_{21}^2 & w_{22}^2 & w_{23}^2 \\ w_{31}^2 & w_{32}^2 & w_{33}^2 \\ w_{41}^2 & w_{42}^2 & w_{43}^2 \end{bmatrix}_{4 \times 3} \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \\ o_{14} \end{bmatrix}_{4 \times 1} + \begin{bmatrix} b_{21} \\ b_{22} \\ b_{23} \end{bmatrix}_{3 \times 1}$$

Transposing  $w$   $\rightarrow w^T$

$$\begin{bmatrix} w_{11}^2 & w_{21}^2 & w_{31}^2 & w_{41}^2 \\ w_{12}^2 & w_{22}^2 & w_{32}^2 & w_{42}^2 \\ w_{13}^2 & w_{23}^2 & w_{33}^2 & w_{43}^2 \end{bmatrix}_{3 \times 4}$$

Multiplying with  $X$  (outputs) and adding bias and then  $G$

$$G \left( \begin{bmatrix} w_{11}^2 o_{11} + w_{21}^2 o_{12} + w_{31}^2 o_{13} + w_{41}^2 o_{14} + b_{21} \\ w_{12}^2 o_{11} + w_{22}^2 o_{12} + w_{32}^2 o_{13} + w_{42}^2 o_{14} + b_{22} \\ w_{13}^2 o_{11} + w_{23}^2 o_{12} + w_{33}^2 o_{13} + w_{43}^2 o_{14} + b_{23} \end{bmatrix} \right)$$



$$\begin{bmatrix} o_{21} \\ o_{22} \\ o_{23} \end{bmatrix}$$

Layer 3:

$$\begin{bmatrix} w_{11}^3 & w_{12}^3 \\ w_{21}^3 & w_{22}^3 \\ w_{31}^3 & w_{32}^3 \end{bmatrix} \begin{bmatrix} o_{21} \\ o_{22} \\ o_{23} \end{bmatrix} + \begin{bmatrix} b_{31} \\ b_{32} \end{bmatrix}$$

$w^T$

$$\begin{bmatrix} w_{11}^3 & w_{21}^3 & w_{31}^3 \\ w_{12}^3 & w_{22}^3 & w_{32}^3 \end{bmatrix} \begin{bmatrix} o_{21} \\ o_{22} \\ o_{23} \end{bmatrix} + \begin{bmatrix} b_{31} \\ b_{32} \end{bmatrix}$$

$$\sigma(w^T x + b)$$

$$g \left( \begin{bmatrix} w_{11}^3 o_{21} + w_{21}^3 o_{22} + w_{31}^3 o_{23} + b_{31} \\ w_{12}^3 o_{21} + w_{22}^3 o_{22} + w_{32}^3 o_{23} + b_{32} \end{bmatrix} \right)$$



$$\begin{bmatrix} o_{31} \\ o_{32} \end{bmatrix}$$

$$\text{Layer } h : \begin{bmatrix} w_{11}^h \\ w_{21}^h \end{bmatrix} \begin{bmatrix} o_{31} \\ o_{32} \end{bmatrix} + \begin{bmatrix} b_{n1} \end{bmatrix}$$

$$g(w^T x + b)$$

$$g \left( \left[ w_{11}^h o_{31} + w_{21}^h o_{32} + b_{n1} \right] \right)$$



$$[o_{n1}]$$

final output.

The outputs of current layer are given as inputs to the next layer's perceptrons.

This process keeps on following until we end up with final output layer output.

This is how training is done.

# Back Propagation

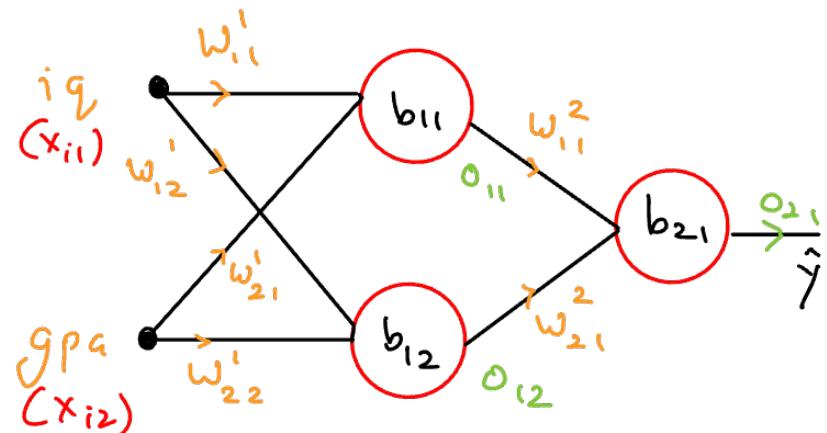
Back propagation (backward propagation of Errors) is the core algorithm for training neural networks. It works by performing a forward pass (forward propagation) to compute predictions calculating the loss between predicted and actual outputs, and then executing a backward pass to compute gradients using the chain rule of calculus. These gradients are used with optimization methods like gradient descent to update weights and biases, reducing prediction error over multiple epochs.

## Key steps in the process:

1. **Forward Pass** – Inputs are multiplied by weights, biases are added, and activation functions (e.g., Sigmoid, ReLU, Softmax) produce outputs layer by layer.
2. **Loss Calculation** – A loss function (e.g., Mean Squared Error, Cross-Entropy) measures prediction error.
3. **Backward Pass** – Gradients of the loss with respect to each parameter are computed from output to input layers.
4. **Weight Update** – Parameters are adjusted in the opposite direction of the gradient, scaled by the learning rate.

# Step by Step implementation

iq	gpa	Salary (lakh)
80	3.5	1.8
60	4	2.5
70	3.2	2



1. Initialize the values for  $w, b$

There are different initialization techniques:

- Random values for  $w, b$
- $w \rightarrow 1, b \rightarrow 0$
- And many more

We will take  $w \rightarrow 1$  and  $b \rightarrow 0$  for simplicity.

Here, we have used "linear" as the activation function.

Cuz this is a regression problem.

(Remember, linear is not mse, because this is activation function not loss function)

2. Select a point (from row of dataset) (student).

Now we use the variables and predict the salary of 1<sup>st</sup> student. This is forward propagation. (ex, we got  $\hat{y} = 5$ )

3. Choose a loss function (MSE)

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

$$\begin{aligned} L &= (\text{ground truth} - \text{predicted value})^2 \\ &= (1.8 - 5)^2 \\ &= 10.2 \end{aligned}$$

(This is very wrong).  
error

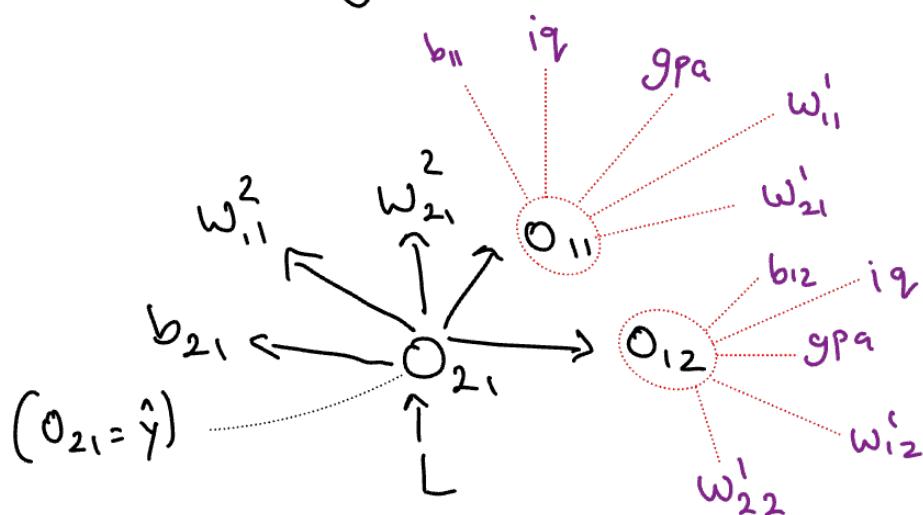
So, now we try to change the value of  $\hat{y}$  (we cannot change  $y$  because it is the ground truth)

To change the value of  $\hat{y}$ , we have to change all the values that resulted in the current value. Those are:

$$O_{21} = w_{11}^2 O_{11} + w_{21}^2 O_{12} + b_{21}$$

(note that there's no activation function at the end (final output)  
that's why we are not applying any!)

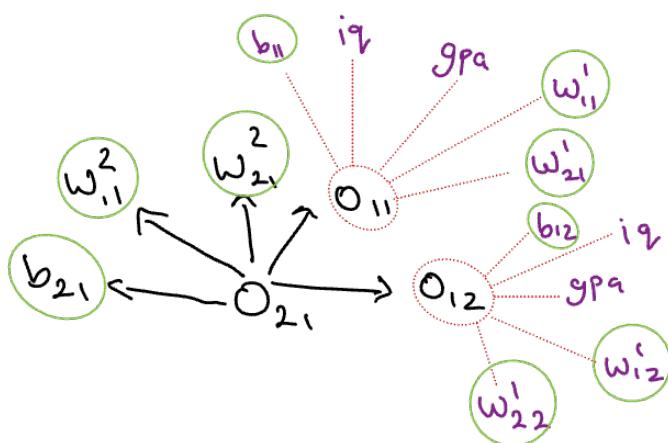
Now again,  $O_{11}$  and  $O_{12}$ 's output depends on other weights too.



↳ Final step : Updating Weights and Bias  
Using Gradient Descent  
↓ formulas

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

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



Here, we have a total of 9 trainable parameters that we can change.  
(ofc! we can't change the actual data like gpa, iq.)

Example formula for one of the trainable parameters

$$b_{21,\text{new}} = b_{21,\text{old}} - \eta \frac{\partial L}{\partial b_{21,\text{old}}}$$

We have to apply this same formula for each of those 9 trainable parameters.

We already have value of  $b_{21}$  (we took 0) and if its weight  $w$ , we take  $w=1$  initially.

We also have the value of  $\eta$  (learning rate) (0.1)  
But we don't have the value for  $\frac{\partial L}{\partial b_{21,\text{old}}}$ .

So we have to calculate the values of all parameters w.r.t. Loss.

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

$$\frac{\partial L}{\partial w_{11}^1}, \frac{\partial L}{\partial w_{21}^1}, \frac{\partial L}{\partial b_{11}}$$

$$\frac{\partial L}{\partial w_{12}^1}, \frac{\partial L}{\partial w_{22}^1}, \frac{\partial L}{\partial b_{12}}$$

Let's calculate the first set's 1<sup>st</sup> parameter:

$\frac{\partial L}{\partial w_{11}^2}$ : There might be a question that why are we calculating Loss function wrt  $w_{11}^2$  when it's not even directly connected?



The reason is that we want to see what will be the change in value of  $O_{21}$  if we change  $w_{11}^2$  and then based on that, what will be getting the change of  $L$  (loss) when  $O_{21}$  gets changed.

We will be using chain rule for differentiation to break down  $\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_{11}^2}$

( $\hat{y}$  is  $O_{21}$ )

$$\textcircled{1} \quad \frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_{11}^2}$$

We need to find value of  $\frac{\partial L}{\partial \hat{y}}$  and  $\frac{\partial \hat{y}}{\partial w_{11}^2}$

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

$$\frac{\partial \hat{y}}{\partial w_{11}^2} = \frac{\partial}{\partial w_{11}^2} (O_{11}w_{11}^2 + O_{21}w_{21}^2 + b_{21}) \rightarrow \boxed{O_{11}}$$

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

$$\textcircled{2} \quad \frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial \hat{y}} \times \underbrace{\frac{\partial \hat{y}}{\partial w_{21}^2}}$$

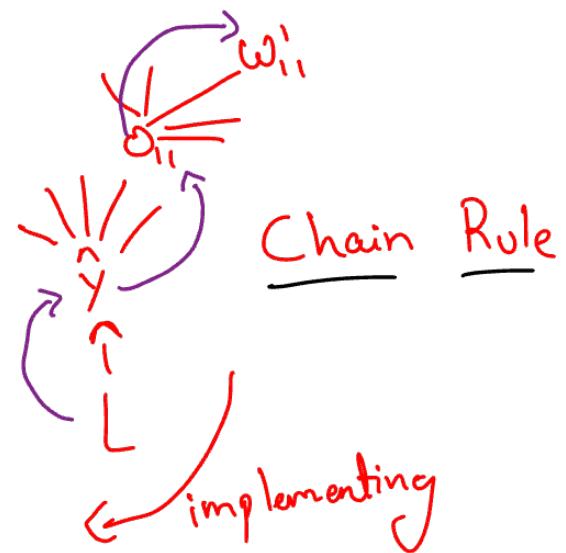
$$\frac{\partial \hat{y}}{\partial w_{21}^2} = \frac{\partial}{\partial w_{21}^2} (O_{11}w_{11}^2 + O_{12}w_{21}^2 + b_{21}) \\ = O_{12}$$

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

$$\textcircled{3} \quad \frac{\partial L}{\partial b_{21}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial b_{21}}$$

$$\frac{\partial \hat{y}}{\partial b_{21}} = 1$$

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



Next Set  $\rightarrow$

$$\textcircled{4} \quad \frac{\partial L}{\partial w_{11}^i} = \underbrace{\frac{\partial L}{\partial \hat{y}}}_{\text{already have it from above}} \underbrace{\frac{\partial \hat{y}}{\partial O_{11}}}_{\text{from above}} \underbrace{\frac{\partial O_{11}}{\partial w_{11}^i}}$$

$$\frac{\partial \hat{y}}{\partial O_{11}} = \frac{\partial (w_{11}^i O_{11} + w_{21}^i O_{21} + b_{21})}{\partial O_{11}} = w_{11}^i$$

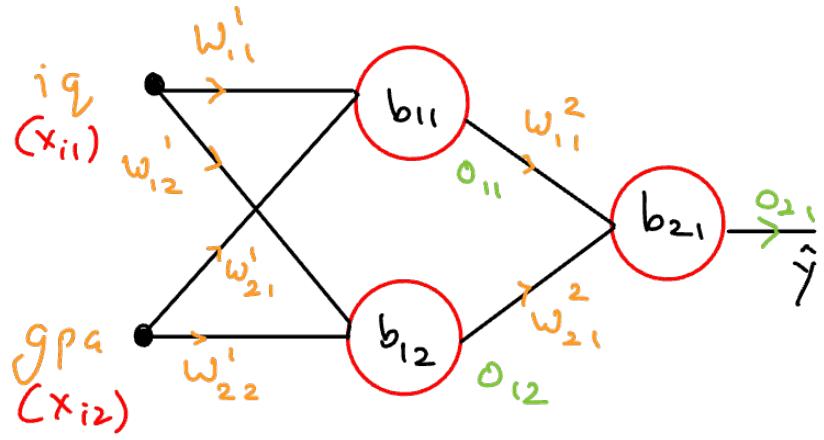
$$\begin{aligned} \frac{\partial O_{11}}{\partial w_{11}^i} &= \frac{\partial (w_{11}^i x_{i1} + w_{21}^i x_{i2} + b_{11})}{\partial w_{11}^i} \\ &= x_{i1} \end{aligned}$$

$$\therefore \boxed{\frac{\partial L}{\partial w_{11}^i} = -2(y - \hat{y}) w_{11}^i x_{i1}}$$

$$\textcircled{5} \quad \frac{\partial L}{\partial w_{21}^i} = \frac{\partial L}{\partial \hat{y}} \underbrace{\frac{\partial \hat{y}}{\partial O_{11}}}_{\text{from above}} \underbrace{\frac{\partial O_{11}}{\partial w_{21}^i}}$$

$$\begin{aligned} \frac{\partial O_{11}}{\partial w_{21}^i} &= \frac{\partial (w_{11}^i x_{i1} + w_{21}^i x_{i2} + b_{11})}{\partial w_{21}^i} \\ &= x_{i2} \end{aligned}$$

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

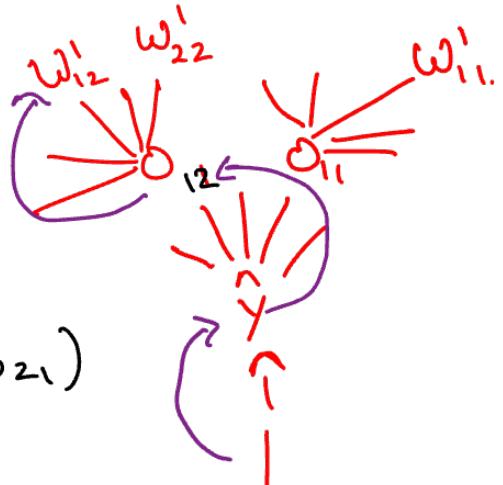


$$⑥ \frac{\partial L}{\partial b_{11}} = \underbrace{\frac{\partial L}{\partial \hat{y}}}_{\text{L}} \underbrace{\frac{\partial \hat{y}}{\partial o_{11}}}_{\text{L}} \underbrace{\frac{\partial o_{11}}{\partial b_{11}}}_{\text{L}}$$

$$\frac{\partial o_{11}}{\partial b_{11}} = \frac{\partial (w_{11}^1 x_{i1} + w_{21}^1 x_{i2} + b_{11})}{\partial b_{11}} = 1$$

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

$$⑦ \frac{\partial L}{\partial w_{12}^1} = \underbrace{\frac{\partial L}{\partial \hat{y}}}_{\text{L}} \underbrace{\frac{\partial \hat{y}}{\partial o_{12}}}_{\text{L}} \underbrace{\frac{\partial o_{12}}{\partial w_{12}^1}}_{\text{L}}$$



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

$$\frac{\partial o_{12}}{\partial w_{12}^1} = \frac{\partial (w_{12}^1 x_{i1} + w_{22}^1 x_{i2} + b_{12})}{\partial w_{12}^1} = x_{i1}$$

$$\therefore \frac{\partial L}{\partial w_{i2}^1} = -2(\gamma - \hat{y}) w_{21}^2 x_{i1}$$

$$\textcircled{8} \quad \frac{\partial L}{\partial w_{22}^1} = \underbrace{\frac{\partial L}{\partial \hat{y}}}_{\text{we have it}} \underbrace{\frac{\partial \hat{y}}{\partial O_{12}}}_{\text{we have it}} \underbrace{\frac{\partial O_{12}}{\partial w_{22}^1}}$$

$$\begin{aligned} \frac{\partial O_{12}}{\partial w_{22}^1} &= \frac{\partial}{\partial w_{22}^1} (\omega_{12}^1 x_{i1} + \omega_{22}^1 x_{i2} + b_{12}) \\ &= x_{i2} \end{aligned}$$

$$\therefore \frac{\partial L}{\partial w_{22}^1} = -2(\gamma - \hat{y}) w_{21}^2 x_{i2}$$

$$\textcircled{9} \quad \frac{\partial L}{\partial b_{12}} = \underbrace{\frac{\partial L}{\partial \hat{y}}}_{\text{we have it}} \underbrace{\frac{\partial \hat{y}}{\partial O_{12}}}_{\text{we have it}} \underbrace{\frac{\partial O_{12}}{\partial b_{12}}}_{\text{we have it}}$$

$$\begin{aligned} \frac{\partial O_{12}}{\partial b_{12}} &= \frac{\partial}{\partial b_{12}} (\omega_{12}^1 x_{i1} + \omega_{22}^1 x_{i2} + b_{12}) \\ &= 1 \end{aligned}$$

$$\therefore \frac{\partial L}{\partial b_{12}} = -2(\gamma - \hat{y}) w_{21}^2$$

Going over the steps once again

o) weight and bias value initialize  
for i in epochs: (usually 100-1000)

for i in range(size):

i) 1 student  $\rightarrow$  forward propagation  $\rightarrow$  predict

ii) Loss calculate (mse)

iii) Adjust all weight and biases

↓

$$w_{\text{new}} = w_{\text{old}} - \frac{n \delta L}{\sum \Delta w_{\text{old}}}$$

iv) Calculate the average loss for epochs

Now, let's apply this algorithm to a real world dataset from scratch.  
We already have the gradient descent formulas:

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

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

$$\frac{\partial L}{\partial w_{21}^2} = -2(y - \hat{y}) \theta_{12}$$

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

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

$$\frac{\partial L}{\partial w_{12}^1} = -2(y - \hat{y}) w_{21}^2 x_{i1}$$

$$\frac{\partial L}{\partial w_{12}^2} = -2(y - \hat{y}) w_{11}^2 x_{i1}$$

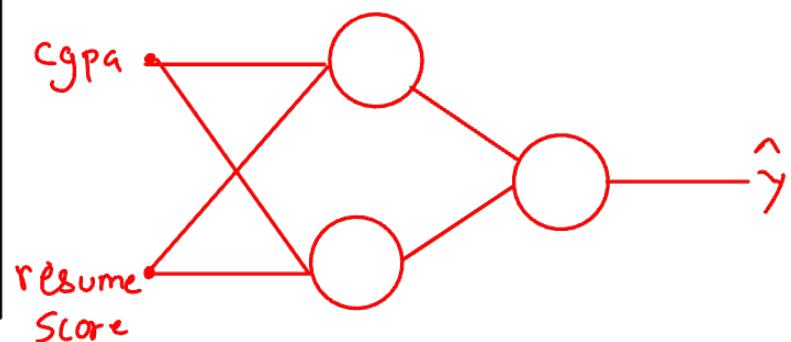
$$\frac{\partial L}{\partial w_{22}^1} = -2(y - \hat{y}) w_{21}^2 x_{i2}$$

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

# Regression Problem

Applying the algorithm from scratch to custom dataset.

cgpa	resume score	salary (1 pa)
8	8	4
7	9	5
6	10	6
5	12	7



This dataset is applied and coded in one of the files in my repo.

Same goes for classification.

fun  
fact

We take derivation as  $\frac{dy}{dx}$  when the derivative has one variable that changes w.r.t to other but when we have more than one, we take  $\frac{\partial z}{\partial y}$  and this is simply gradient.

# Back propagation Intuition (Last time)

We know that in back propagation the main formula that updates the values of weights and biases is

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

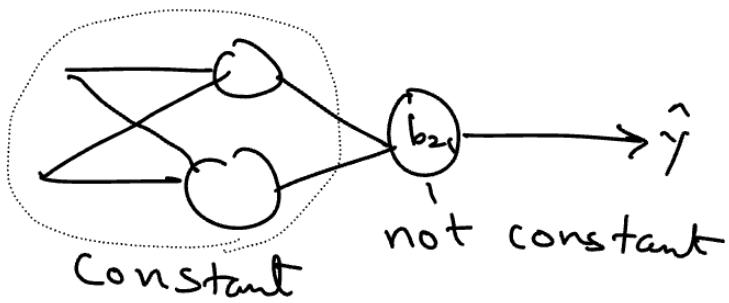
But why does this formula work?

→ for now, assume  $\eta = 1$



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

And for now, just assume that all weights and biases are constant except the bias of  $\hat{y}$  (just for explanation purpose)



$$\therefore L(b_{21})$$

$b_{21} = b_{21} - \frac{\partial L}{\partial b_{21}}$  → Why are we subtracting  
 this value  
 and why is it  
 subtraction ?

So, the -ve sign is there because:

1) We know, if we change value of  $b_{21}$ , then value of  $L$  will also change.  $\frac{\partial L}{\partial b_{21}}$

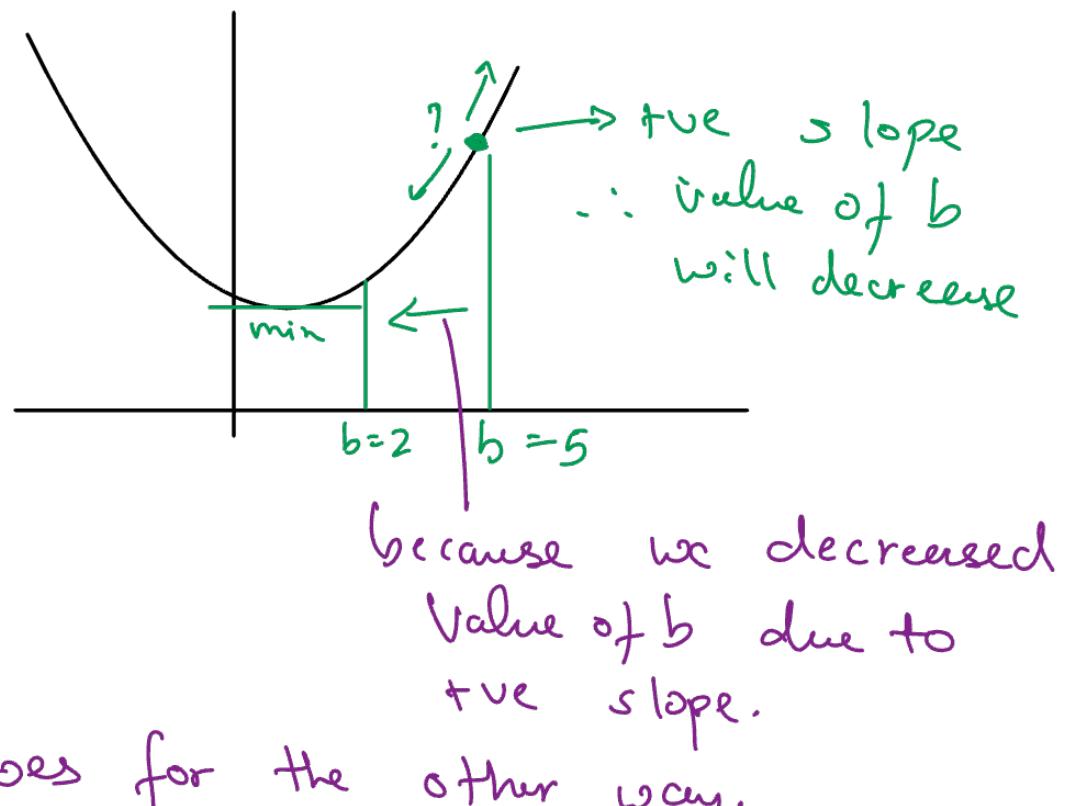
If value of  $\frac{\partial L}{\partial b_{21}} = +ve$ , then we increase the value of  $b_{21}$ ,

the value of  $L$  will also increase,  
 but we know the decrease the value of  $L$  (as much possible)  
 we have to  
 but if  $\frac{\partial L}{\partial b_{21}}$  is -ve, then when we  
 increase  $b_{21}$ , loss decreases, which is what we want!

So, the -ve sign in  $b_{21} - \frac{\partial L}{\partial b_{21}}$  is very smartly placed there that it handles both the cases.

Therefore, simply when the loss is at  $x$ , it knows either to move in -ve or +ve direction that ends up with minima of GD.

Our algorithm will simply calculate the derivative (slope).



Same goes for the other way.

That's why, gradient descent always goes towards the minima.

And, how much change should the value take? the same as the current value therefore, we decrease the value from its own :

$$w_{\text{new}} = w_{\text{old}} - n \frac{\partial L}{\partial w_0}$$

# What is Convergence?

Simply when we reach our minima,  
it's called convergence.

So, when  $w_n$  is roughly close to  $w_0$   
it's convergence.

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

The standard way is generally to  
run external loop which we refer  
to epochs.

We usually don't hard code it that  
when it reaches convergence it should  
stop because it's not always that  
 $w_{\text{new}}$  is exactly same as  $w_{\text{old}}$ .

# Memoization

Memoization is an optimization technique used primarily to speed up computer programs by storing the results of expensive function calls and returning the cached result when the same inputs occur again.

When a neural network has more than two hidden layer, the derivation that we calculate at every step becomes more and more complex, which significantly increases the time.

To solve this, Backpropagation uses memoization to store the results of differentiations that were done during the process.

In simple words, Back propagation is simply Chain rule (derivatives) + Memoization.

# Gradient Descent

We know that GD is an iterative optimization algorithm used to minimize a cost function by adjusting model's parameters in the steepest descent of the function's gradient.

There are three variants of GD, which differ in how much data we use to compute the gradient of the objective function. Depending on the amount of data, we make trade-off between the accuracy of the parameter update and the time it takes to perform an update.

**Batch GD**: updates parameters using the gradient (vanilla GD)  
default computed over the entire training dataset each step.

**Stochastic GD**: updates parameters using the gradient from a single randomly chosen example each step

**Mini-batch GD**: updates parameters using the gradient from a small random subset of examples each step, balancing speed and stability.

The default backpropagation uses stochastic GD where we update weights  $n$  times (rows)  $\times$  epochs.

Batch GD is faster than stochastic because in batch GD, we update the weight epoch times while in stochastic GD, we run for  $n \times$  epochs.

Stochastic GD is more faster to reach convergence since it updates the weights way more times than it Batch GD.

The Loss function's reduceness is very unstable in stochastic GD while very smooth/stable in batch GD. because we pick one random point in stochastic GD multiple times. The pro for this is that stochastic may sometimes move out of local minima to reach global but the con is that it does not converge at the exact solution. (it is approximate)

Batch GD is very faster because it uses vectorization (dot product) rather than running loops.

The only problem is that if the dataset is very big, it won't work because we have to load the entire dataset into the ram.

Mini-Batch GD is the middle ground.  
(generally MBGD is used in Keras)

Usually, in Keras when we use the batch size parameter in Keras, we generally use the  $2^n$  numbers as parameters.

2, 4, 8, 32, 64, 128, 256.

This is because in giving two multiple input research has proven that it is very effective way to handle the RAM.

# Vanishing Gradient Problem

The **vanishing gradient problem** is a significant challenge in training deep neural networks. It occurs when the gradients of the loss function with respect to the network's weights become exceedingly small during backpropagation. This issue is particularly prevalent in deep networks with many layers, where the gradients diminish as they propagate backward through the network.

When using activation function, specifically tanh or sigmoid, vanishing gradient problem can occur where the initial weights become too much small (after repetitive Back propagation) that the model never reaches convergence, in other words, model does not get trained.

There are two ways to identify if this problem is occurring or not.

1. Observing the Loss function's value.  
If it does not change for multiple epochs.
2. plotting a graph for weights.  
epoch VS value graph.

# Solution to handle Vanishing GrP

1. Reduce model complexity
    - Decrease the number of hidden layers.
  2. Using different activation functions
    - RELU is the most commonly used.
  3. Proper Weight Initialization.
  4. Batch Norm → layer
  5. Residual Network
- Covered  
in next  
topics!