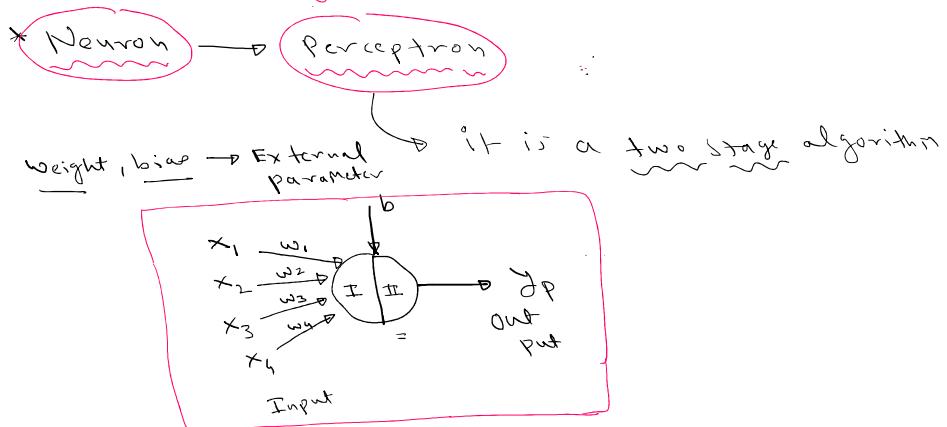


Deep Learning Day - 2



Stage I → Summation function (z)

$$z = \sum_{i=1}^n w_i x_i + b$$

Stage II - Activation function (Sigmoid)

$$\sigma(z) = \frac{1}{1+e^{-z}} = y_p$$

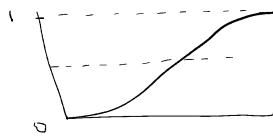
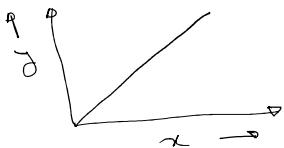
This whole process called as perceptron and perceptron if say act as a Neuron in Neural Network.

* What is Sigmoid ?

→ it is used to add Nonlinearity.

$$y = mx + c$$

$$\sigma(y) = \frac{1}{1+e^{-y}}$$



→ Sigmoid convert continuous values into probability distribution.

→ Range → 0 to 1

→ it always give probabilistic value.

→ it always give probability value of any one class.



$$P(x \in \text{Class } 1) = P(1) = 0.7$$

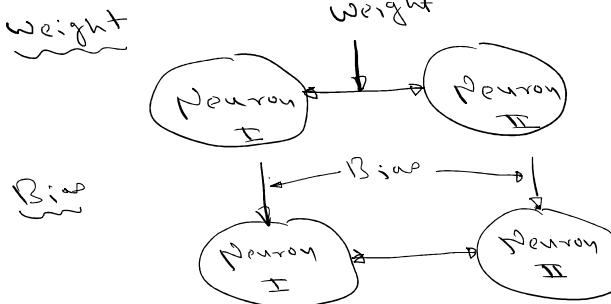
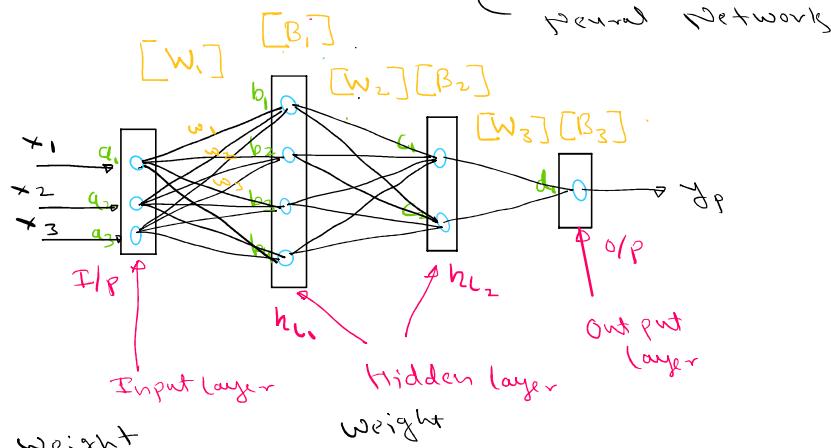
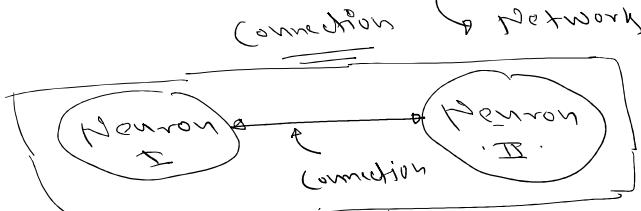
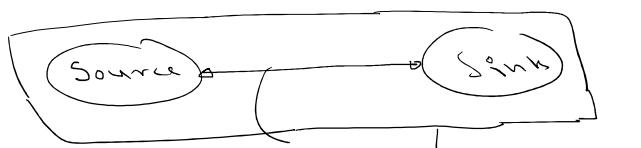
$$P(x \in \text{Class } 0) = P(0) = 1 - P(1)$$

$$\frac{1}{0.7} = \frac{1}{0.3}$$

$$= 0.3$$

threshold → 0.5

* Structure of Neural Network :-



$$[W_1] = [w_1, w_2, w_3, w_n, \dots, w_{12}]$$

$$[B_1] = [b_1, b_2, b_3, b_n] \dots$$

* Input layer → It hold the input and feed them to hidden layer.

* Hidden layer → Done all the processing/learning

* Output layer → It shows the output

* Input layer → Should be one

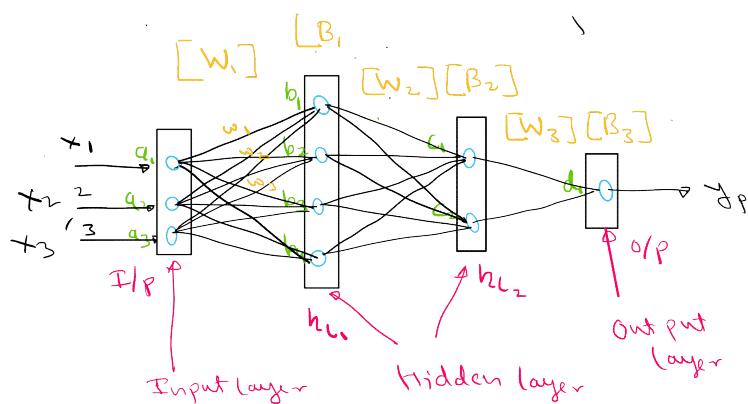
* Hidden layer → One or more than one

* Hidden layer → One or More than one

↓
Shallow Neural Network Deep Neural Network

* Output layer → Should be one

* How data flow in Neural Network



* No of Neuron in Input layer = No of input features

* No of Neuron in out put layer = decided by the activation fun and No of event.

Case I → Let say we are dealing with binary classification and in our out put layer we have Sigmoid activation fun.

No of Neuron in out put layer = one.

* Case II → Let's say we are dealing with multi class classification and in our out put layer we have Softmax as a activation function. and No. of event is 5.

No of Neuron in out put layer = 5

--- → i.e. if we are dealing

out put larger

- * Case III \rightarrow Let's say we are dealing with regression problem and in our output layer we are using as an activation.

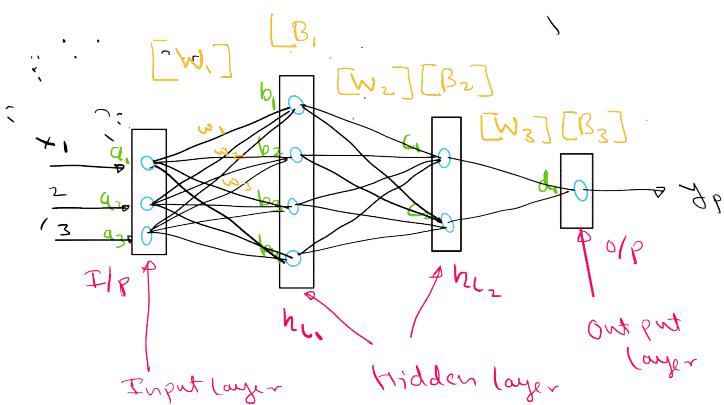
No of Neuron
in out put layer = one

- * No of Neuron in hidden layer = Not fixed

General there's rule (Manually build \rightarrow ANN)

The diagram illustrates a neural network structure:

- Input Layer:** Labeled "input layer".
- Hidden Layer:** Labeled "Hidden layer".
- Output Layer:** Labeled "out put layer".
- Connections:** Arrows indicate connections from the input layer to the hidden layer, and from the hidden layer to the output layer.
- Neuron Labels:** "1st Neuron" is labeled above the first neuron in the hidden layer, and "3rd Neuron" is labeled above the third neuron in the output layer.
- Range Labels:** "larger" is written vertically next to the hidden layer, and "between 3 to 15" is written below the output layer.
- Start and End:** "Start" is at the beginning of the input layer, and "end" is at the end of the output layer.



$$\text{Shape of Input} = (1, 3) \quad \left. \begin{array}{l} \\ \text{layer} \\ (4, 3). \end{array} \right\} \text{Stage I}$$

$$\text{Shape of hidden} = \begin{pmatrix} 1, u \\ 2, u \end{pmatrix} \xrightarrow{\text{larger 1}} \text{Stage II}$$

$$\text{Shape of Hidden Layer 2} = (1, 2) \quad (1, 2) \quad \text{Stage III}$$

$$\text{shape of out put} = (1,1)$$

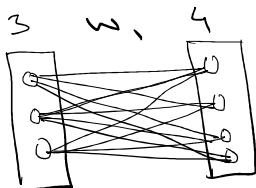
Shape of w in layers

→ How to control shape of layers in neural network.

→ Let's consider w , (weight matrix)

We start from 3 neuron and

We end with 4 neuron



$w \rightarrow 12$ weights.

Shape of w = $(4, 3)$

$$z = \sum_{i=1}^3 w_{i*} x_i + b$$

Start $\rightarrow (1, 3)$

$$(1, 3) \cdot (4, 3) = (1, 4)$$

end $\rightarrow (1, n)$

Transpose

$$(1, 3) (3, 4) = (1, 4)$$

$$(1, 4) = (1, n)$$

$$z = \sum_{i=1}^n w_{i*} x_i + b$$

* Importance of Bias

Bias \rightarrow one dimensional array that consist values of bias.

$$z = \sum_{i=1}^n w_{i*} x_i + b$$

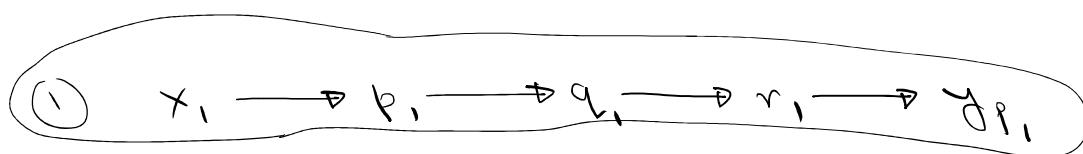
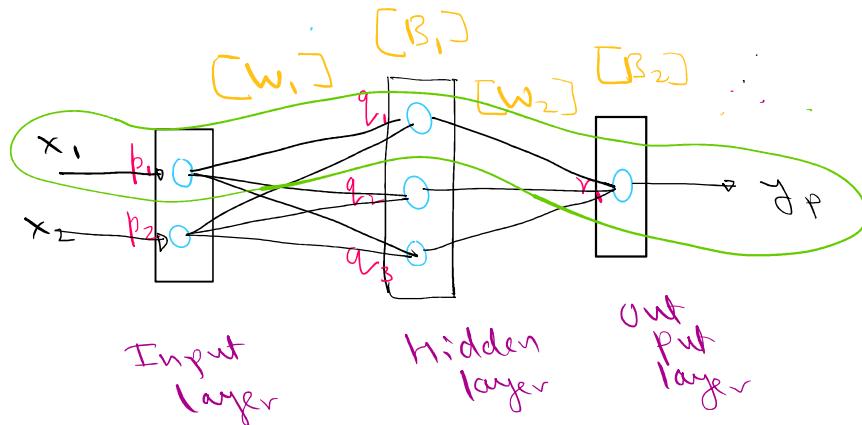
	x_1	x_2	x_3	y
(1)	100	0	300	cat

② $200 \quad 3 \quad 400 \quad \text{Day}$

$$\begin{aligned} z &= w_1x_1 + b_1 + w_2x_2 + b_2 + w_3x_3 + b_3 \\ &= 100w_1 + b_1 + 0 + b_2 + 300w_3 + b_3 \end{aligned}$$

* How Neural Network learn

(main Rule)



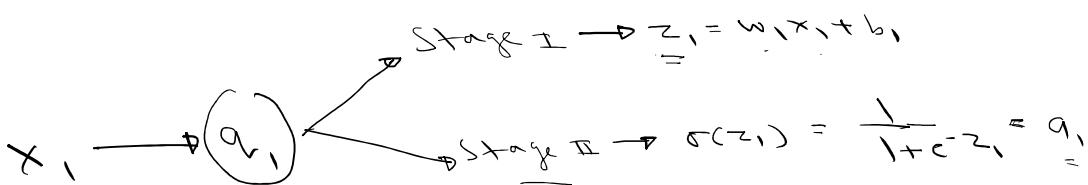
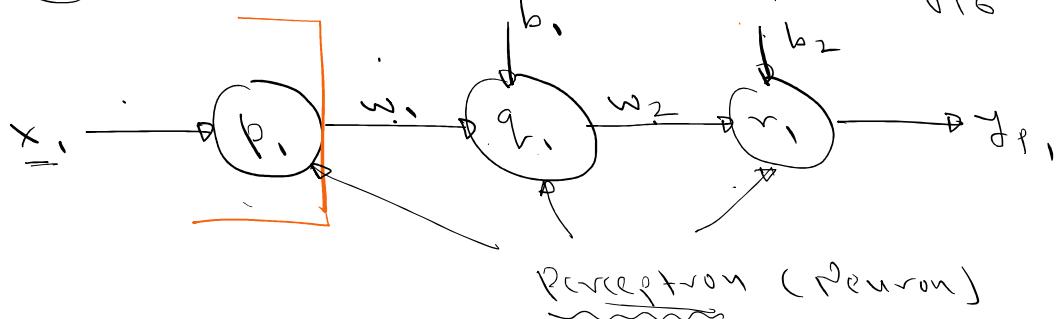
② $x_1 \rightarrow p_1 \rightarrow q_2 \rightarrow r_1 \rightarrow y_{P2}$

③ $x_1 \rightarrow p_1 \rightarrow q_3 \rightarrow r_1 \rightarrow y_{P3}$

④ $x_2 \rightarrow p_2 \rightarrow q_1 \rightarrow r_1 \rightarrow y_{P4}$

⑤ $x_2 \rightarrow p_2 \rightarrow q_2 \rightarrow r_1 \rightarrow y_{P5}$

⑥ $x_2 \rightarrow p_2 \rightarrow q_3 \rightarrow r_1 \rightarrow y_{P6}$



$$a_1 \rightarrow r_1 \xrightarrow{\text{Stage I} \rightarrow z_2 = w_2 a_1 + b_2} \xrightarrow{\text{Stage II} \rightarrow \sigma(z_2) = \frac{1}{1+e^{-z_2}} = y_p,}$$

$$y_a : y_p,$$

$$\log(L_1) = \tilde{y}_a - \tilde{y}_p,$$

Linear Regression

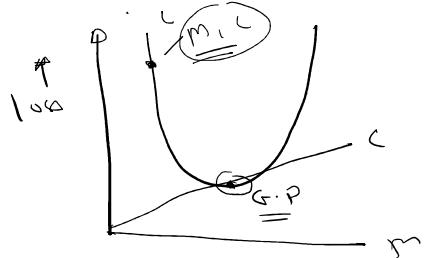
$$y = mx + c$$

Best value of m and c

Gradient Descent algorithm

$$m_{\text{new}} = m_{\text{old}} - \alpha \frac{\partial L}{\partial m}$$

$$c_{\text{new}} = c_{\text{old}} - \alpha \frac{\partial L}{\partial c}$$



$$w_i' = w_i - \alpha \frac{\partial L_i}{\partial w}$$

$$b_i' = b_i - \alpha \frac{\partial L_i}{\partial b}$$



$$\min_{w, b} [L_1, L_2] \text{ s.t. } L_3, L_4, L_5, L_6 \quad \text{G.P.}$$

minimum loss

ANN

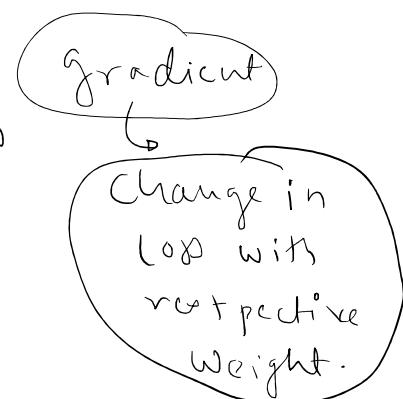
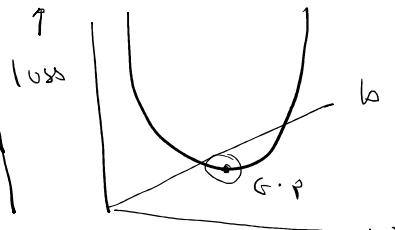
$$z = \sum_{j=0}^n (w_j x_j + b)$$

Best value of w and b

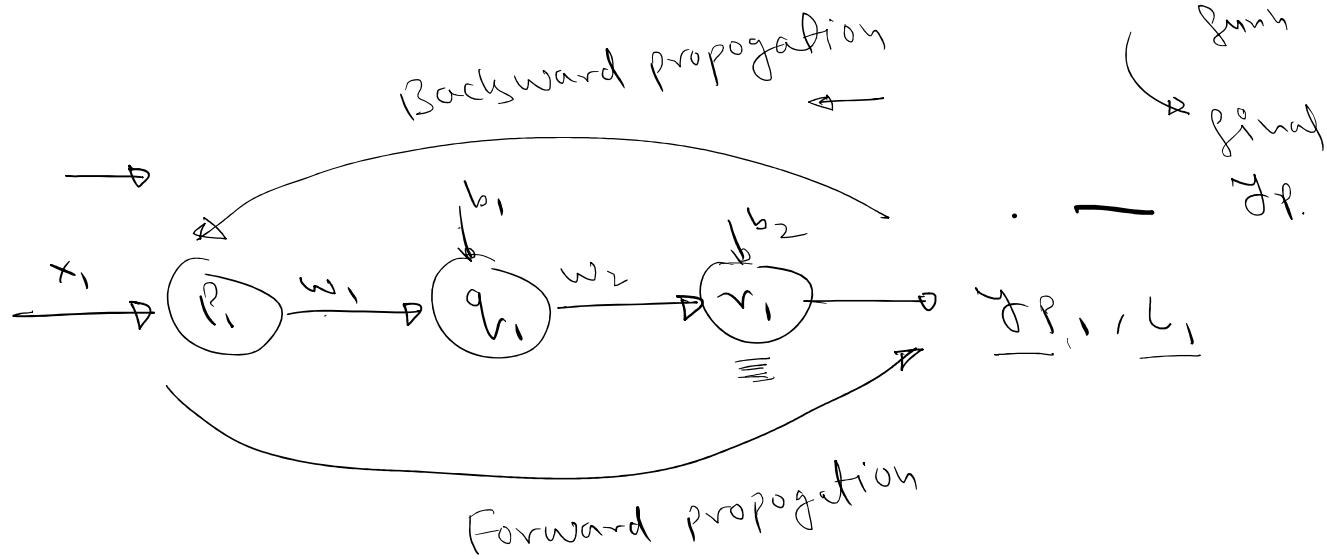
Gradient Descent algorithm.

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

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



$= - \int$ minimum loss
 ↓ best value of
 w and b
 ↓ best summation
 sum
 ↓ activation

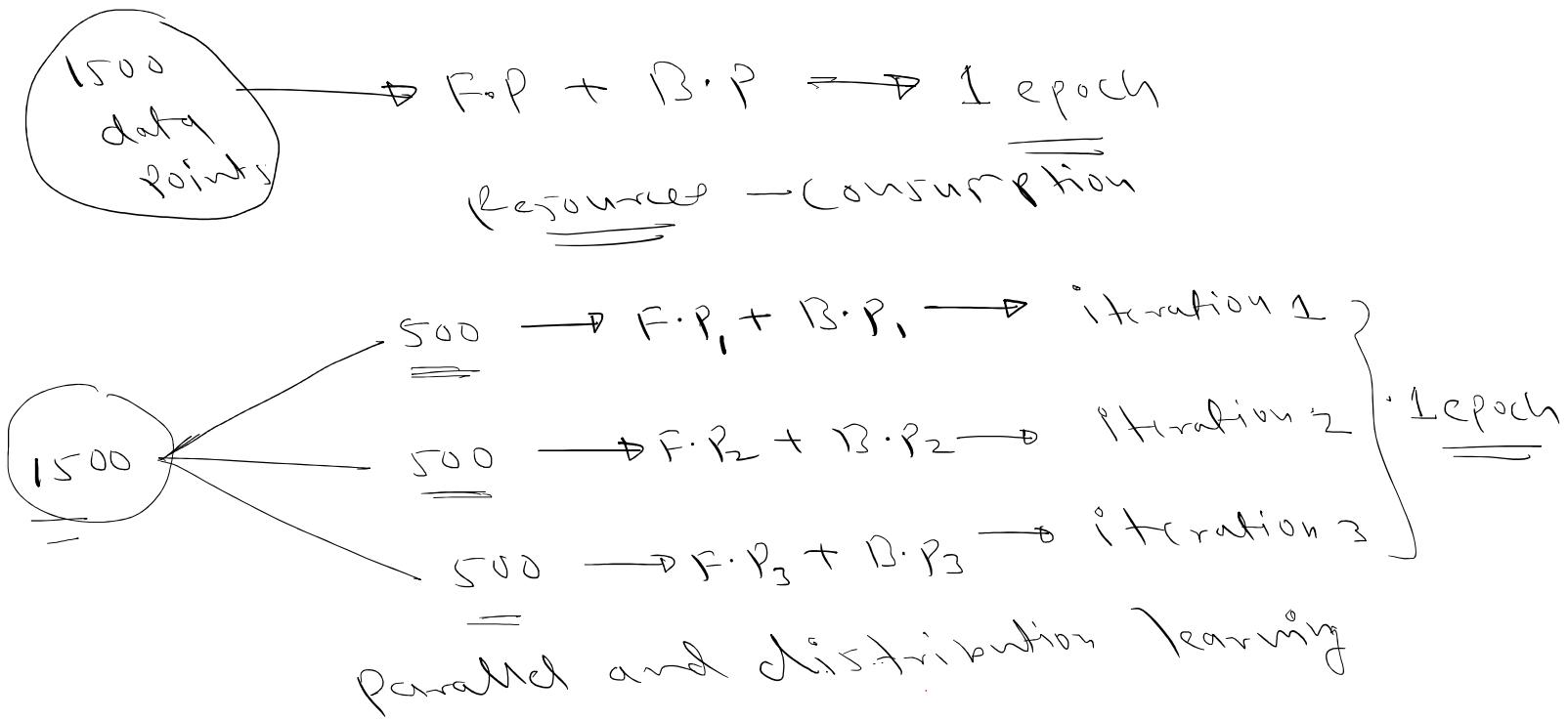


Forward Propagation $\rightarrow y_{f, l_1}$

Backward propagation \rightarrow to update weight and bias.

* Epoch / Iteration

Let's say we have 1500 data points



* : Type of function use in Neural Network

① Optimization func

→ it is used to minimize or optimize the error

→ it is used to find best value for weight and bias

→ it helps to achieve convergence or global minima point.

② Activation func

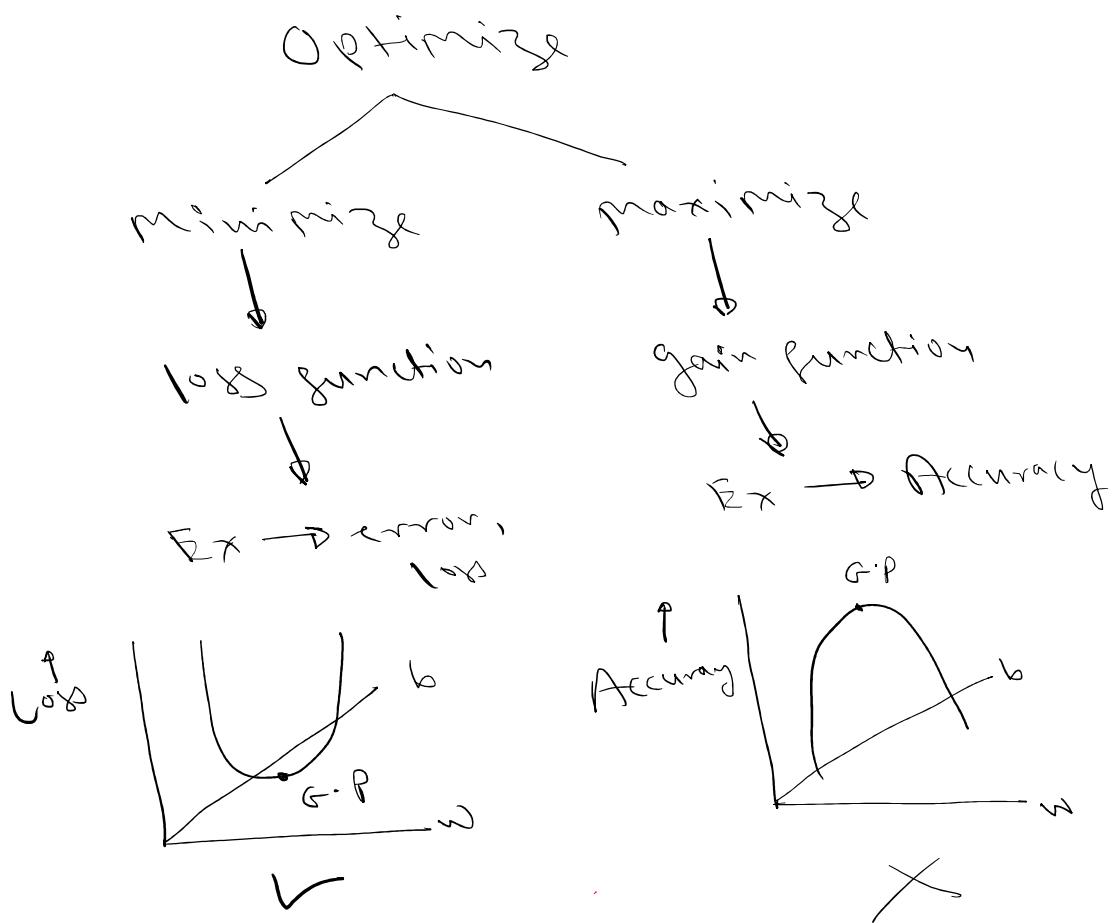
→ Governing the way a neuron behaves.

→ It controls the output of any neuron.

Neuron.

- (3) Loss function :- error ↑ error ↓
- It gives a number which represent the goodness model.
 - It helps to update the neural network parameters.

* Optimization function



* Type of optimization func

1. weight +

2. ...

3. ...

right +

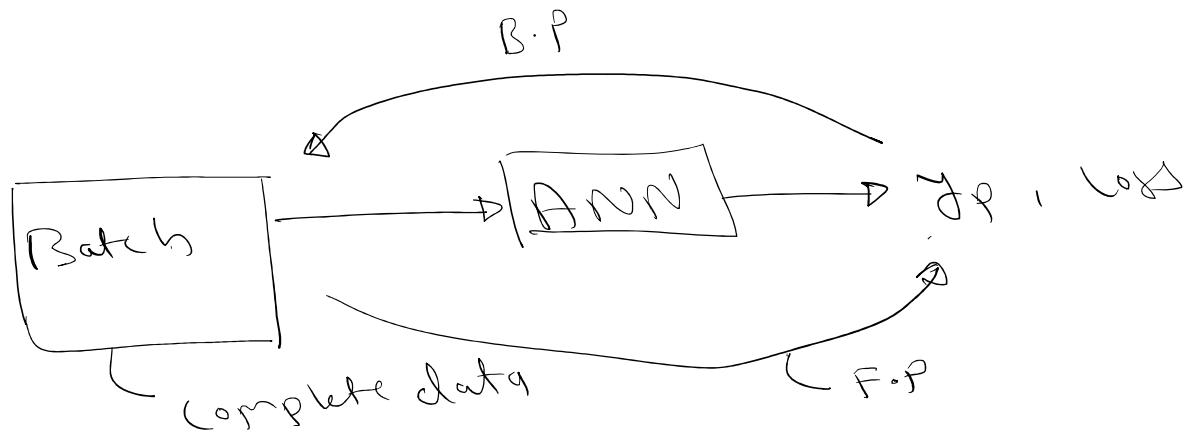
... & third

weight +
 Bias
 (dynamically
 update)

learning
 rate

weight + bias
 + learning

① gradient Decent algorithm :-
 Batch gradient descent - algorithm
 (BGDA)

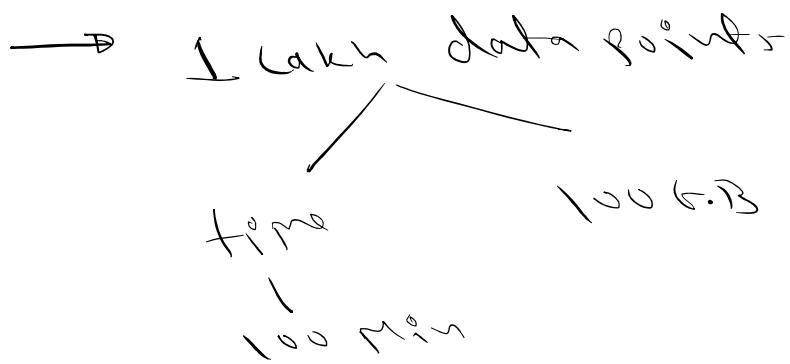


→ in BGDA epoch get perform on whole data , no matter how big our data set.

→ let's say we have 1000 data points and we apply BGDA , and we achieve convergence in 1 min and memory consumption is 1 GB.

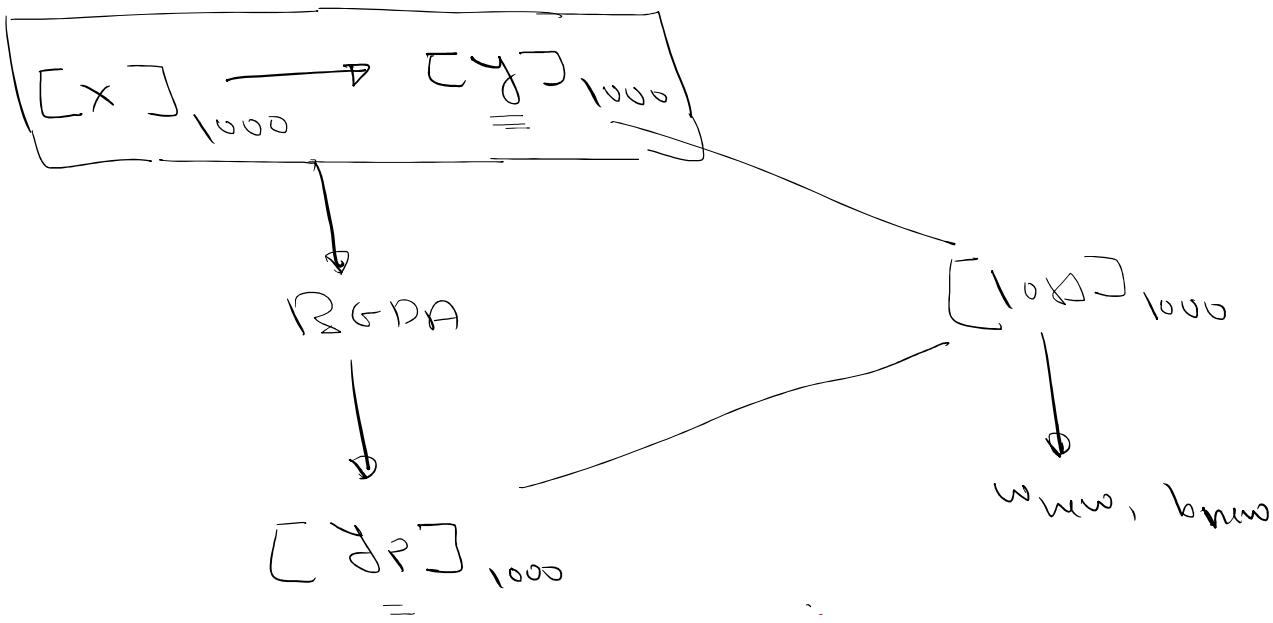
Consumption is 1GB.

1000 - 1

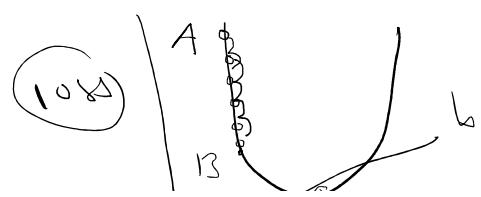


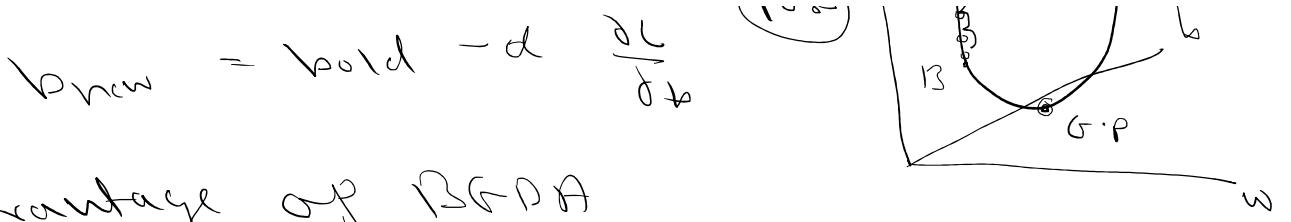
1,00,000 - 100

- * Draw back of BGDA
 - We cannot use BGDA in case of Big data because epoch get perform on whole data.
- * How BGDA work in Back end :-



$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial L}{\partial w}$$
$$b_{\text{new}} = b_{\text{old}} - \alpha \frac{\partial L}{\partial b}$$





* Advantage of BGD

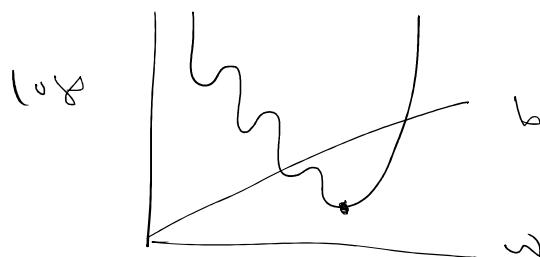
→ We get smooth learning curve in case of BGD because we perform epoch on whole data.

$$\textcircled{1} [2, 4, 6, 8, 10, 12, 14, 16, 18, 20]$$

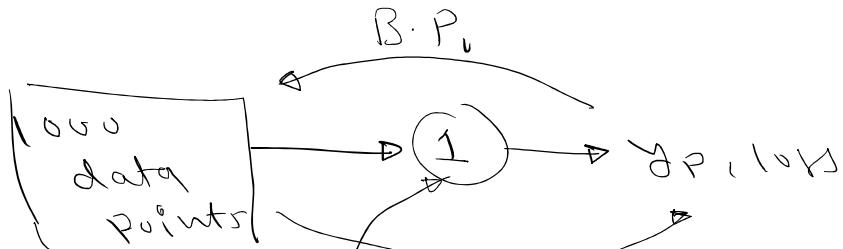
$$\textcircled{2} \underbrace{[2, 4, 6, 8]}_T, \underbrace{[10, 12, 14, 16]}_{\text{learning step}}$$

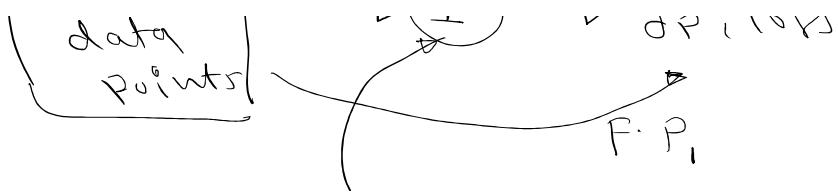
, learning Rate

learning
curve



$\textcircled{2}$ Stochastic gradient descent algorithm (SGD)
Random or Randomness





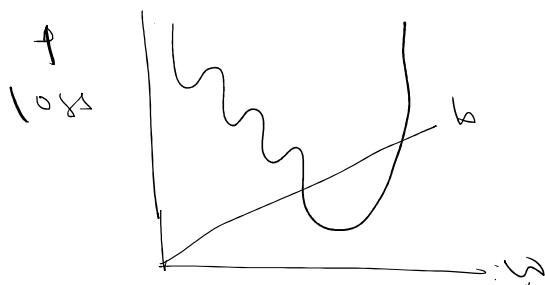
Here, we consider single
data point at a time randomly

→ In SGD → 1 epoch

↓
n → No of iteration

where n → No of data point

→ In SGD we learn error from single data
point at a time hence learning curve
become very noisy.



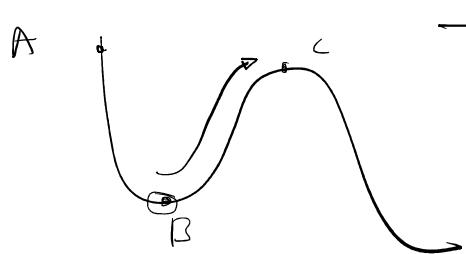
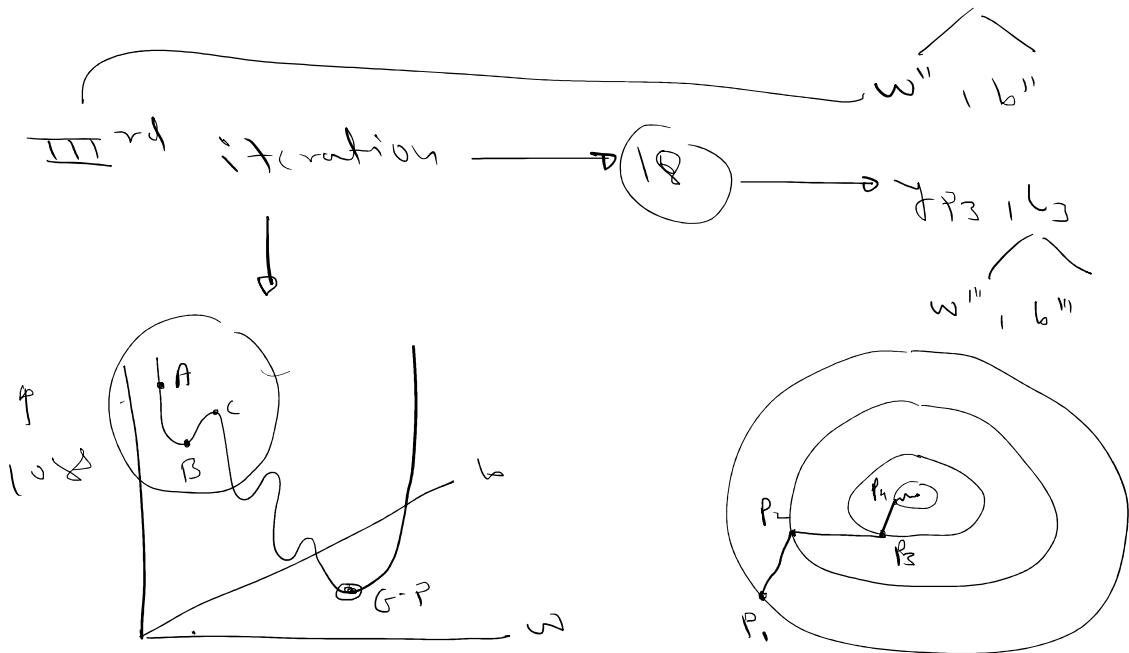
$$[2, 4, \textcircled{6}, 8, 10, \textcircled{12}, 14, 16, \textcircled{18}, 20]$$

w, b

1st iteration → $\textcircled{6} \rightarrow y_{p_1, l_1}$

\dots
 w', b'

2nd iteration → $\textcircled{16} \rightarrow y_{p_2, l_2}$



In SGD since we are learning from single data point at a time we don't have sufficient amount of learning to overcome local minimum point as a result after some time SGD consider local minimum as a global minimum and that is wrong.

* Advantage of SGD

→ Compare to BGD, SGD required less resources to achieve convergence

* Drawback of SGD :-

- * Draw back of SVM -
- It is very prone to local minima because in SGD we learn from single data point at a time hence we don't have sufficient amount of learning to overcome local minima.