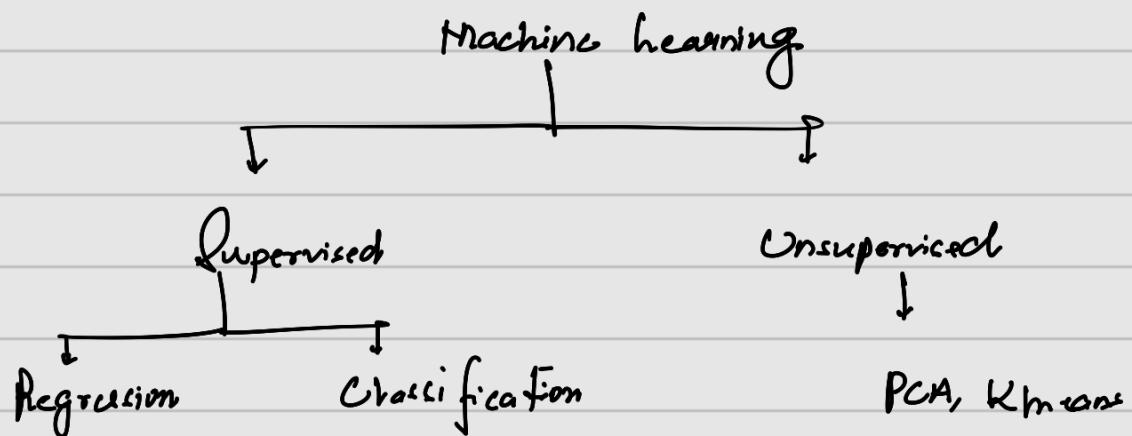


## Introduction to Deep learning:



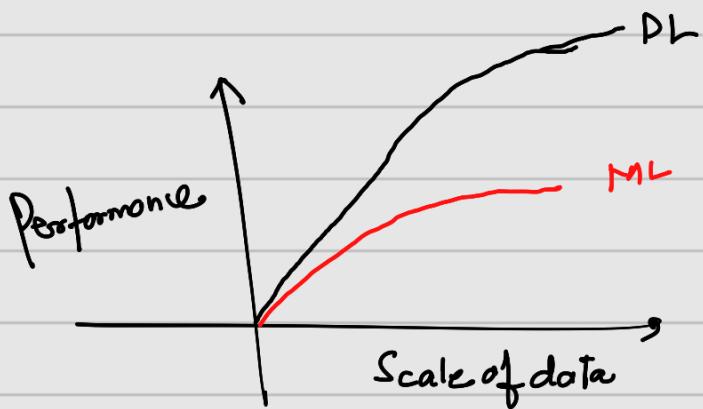
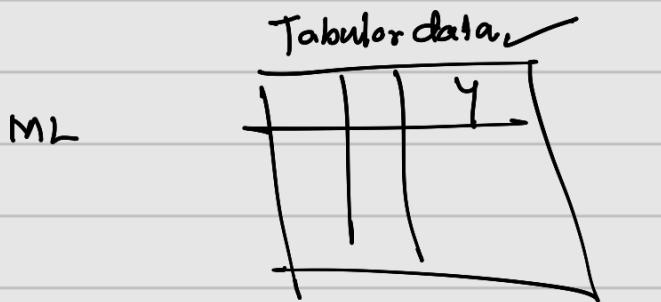
Linear Reg  
Log Reg  
DT, RF, Boosting

Mh → Parametric → equation  $y = f(x)$   
Non-Parametric process DT, RF, XGBoost

Deep learning → Parametric form.

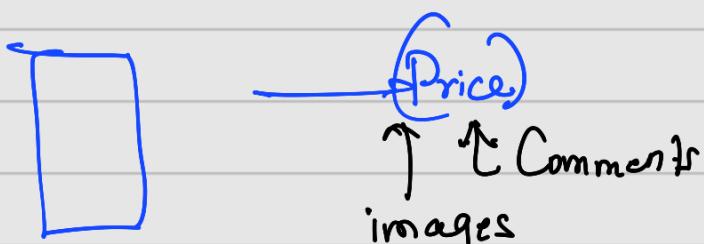
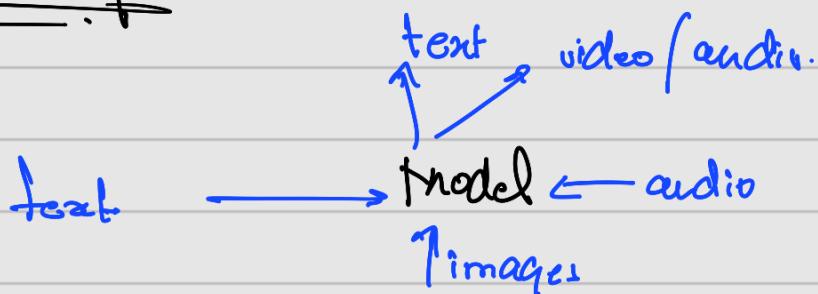
$$y \sim f(x)$$

↓  
equation

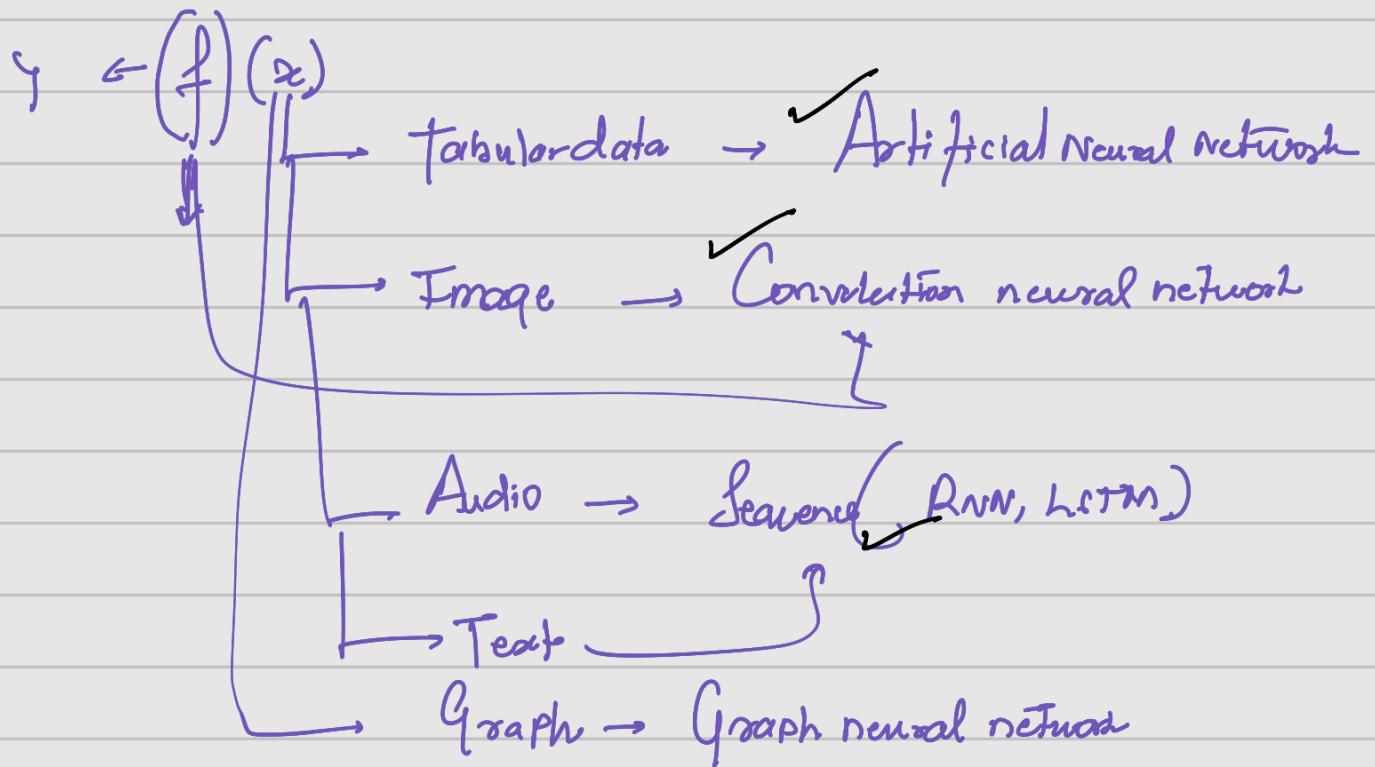


Automated feature engineering → Images, Videos, text,  
Audio, Graph, etc.

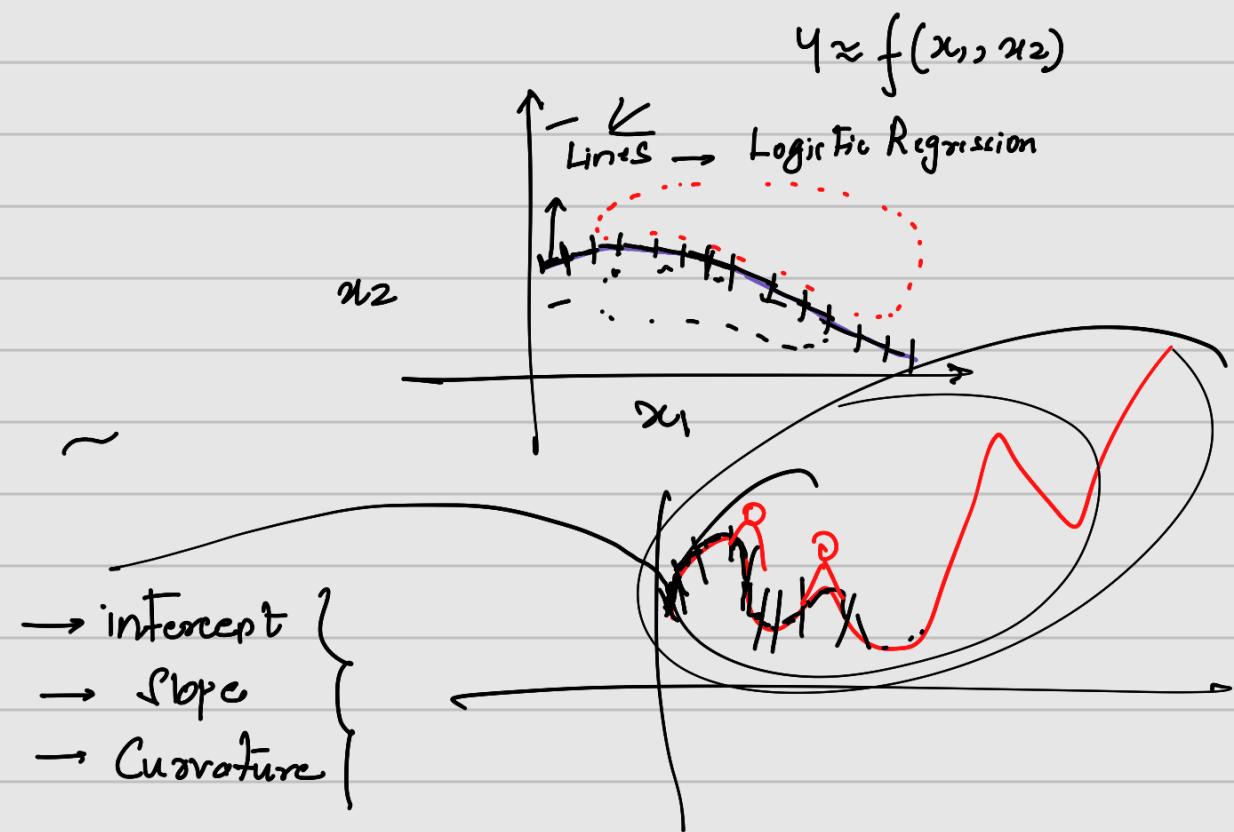
### Multimodality

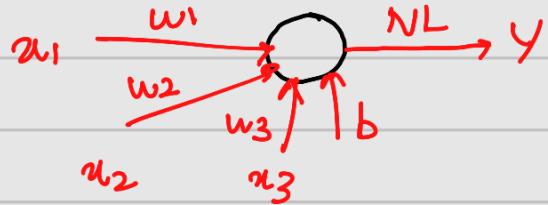


Deep learning a framework wherein we create a parametric model to predict  $y$  as a function of  $x$ .



## Lecture-2 What is a Neuron?

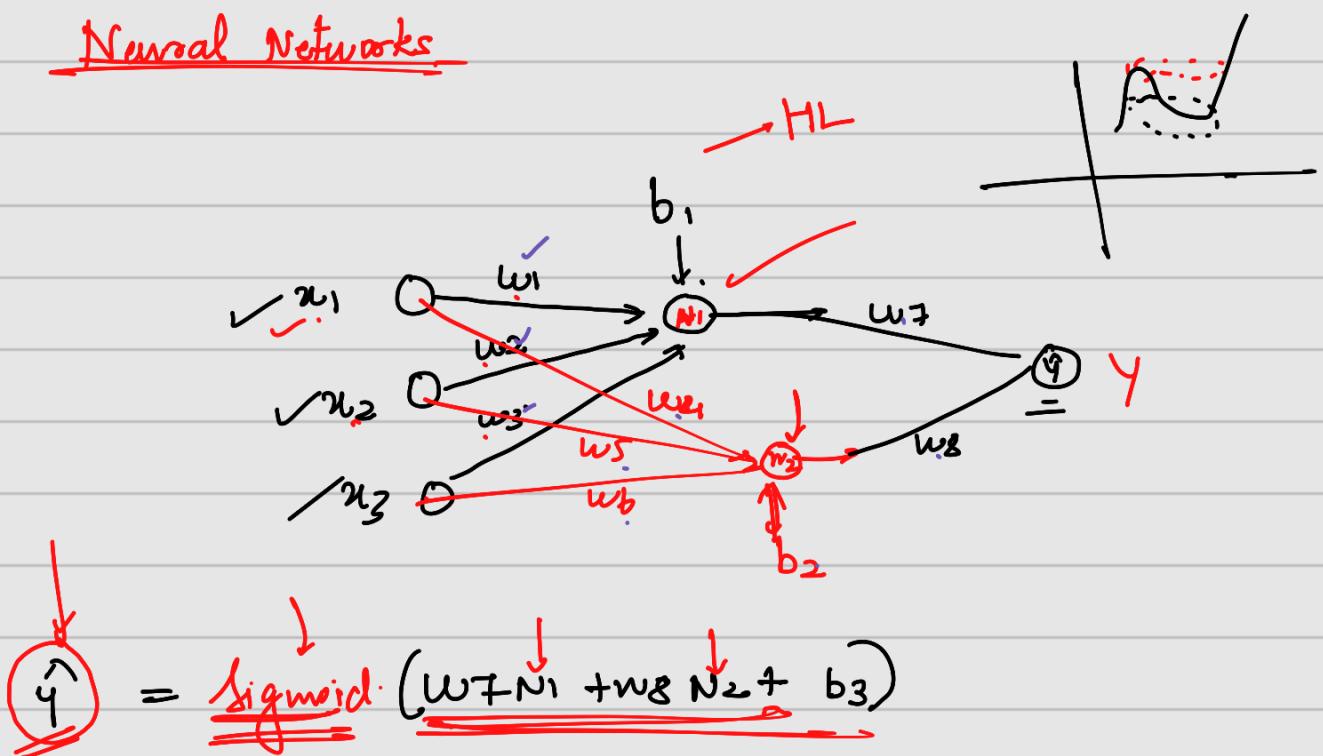




$$y = \text{NL}(w_1x_1 + w_2x_2 + w_3x_3 + b)$$

activation functions

## Neural Networks



$x_1$	$x_2$	$x_3$	$y$
-	-	-	-
-	-	-	-
-	-	-	-

If I want my  $\hat{y}$  to go as close as possible to  $y$ , modifying the weights and biases.

① Define the architecture.

Weights and biases  $\rightarrow$  randomly initialize.

$$(\hat{y} \rightarrow y) \rightarrow \text{loss}$$

② Reduce the loss by modifying the weights and biases.

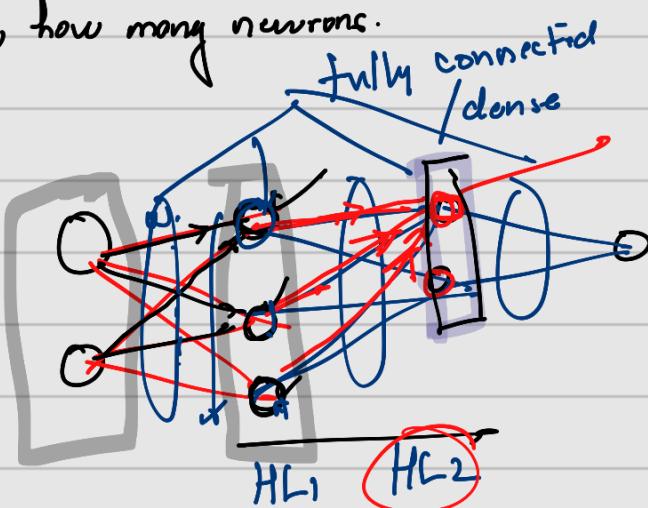
### Gradient Descent

How to train Neural networks:

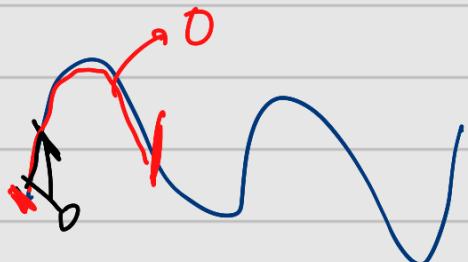
① Define my neural net architecture.

How many HL, how many neurons.

$x_1$	$x_2$	$y$



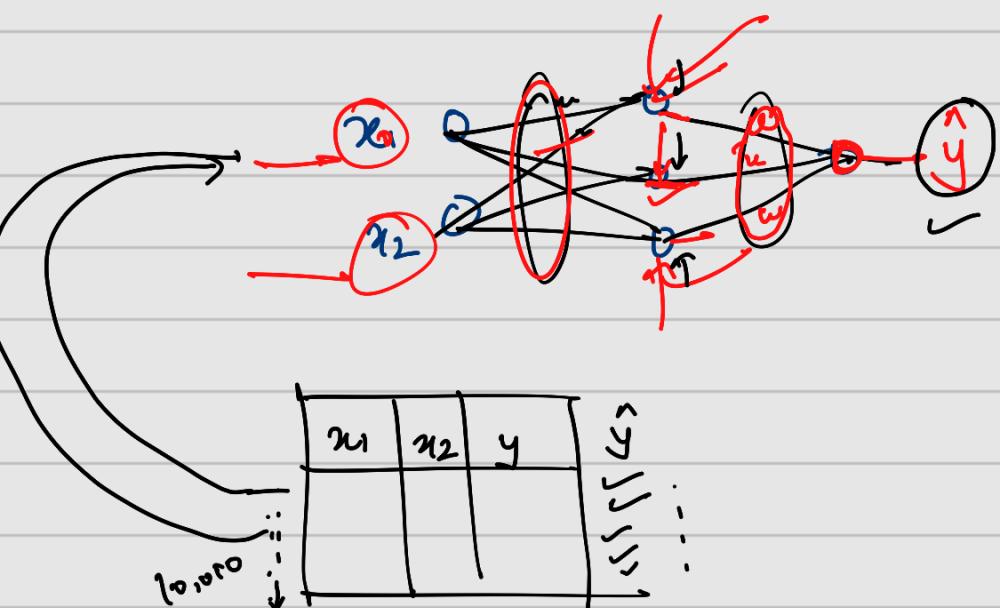
$\xrightarrow{\quad}$  add more degree of non-linearity.





→ all weights and biases will be randomly initialized.

→ given these weights,  $f$  will multiply the weights to the inputs  
add the biases



$x_1$	$x_2$	$y$	$\hat{y}$
1.0	0.0		

$\text{loss} = \sum_{i=1}^n \frac{1}{2} (y_i - \hat{y}_i)^2$

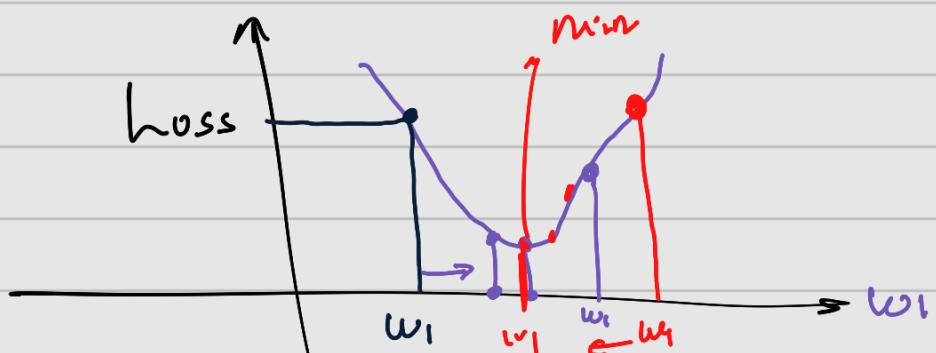
Regression

Forward Propagation

Based on the loss, we find out better set of weights and biases.

$$\text{loss} = \sum_{i=1}^n \frac{1}{n} (y_i - \hat{y}_i)^2$$

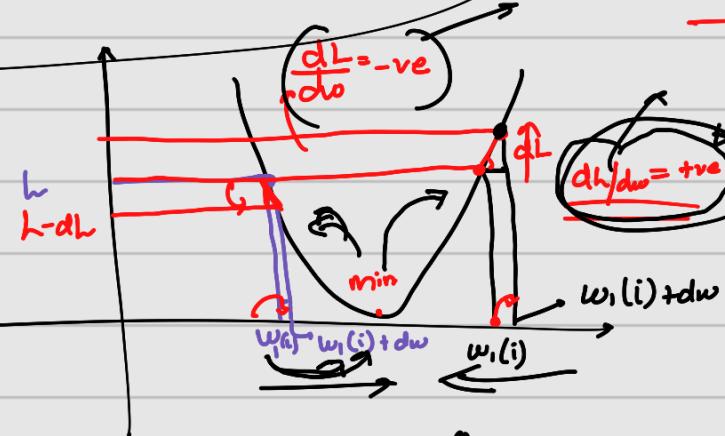
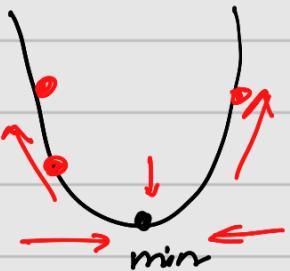
$$\text{loss} = f(\underline{w}, b)$$



for all weights and biases.

$$w_1(i+1) = \underline{w_1(i)} - d$$

$$\frac{dh}{dw(i)}$$

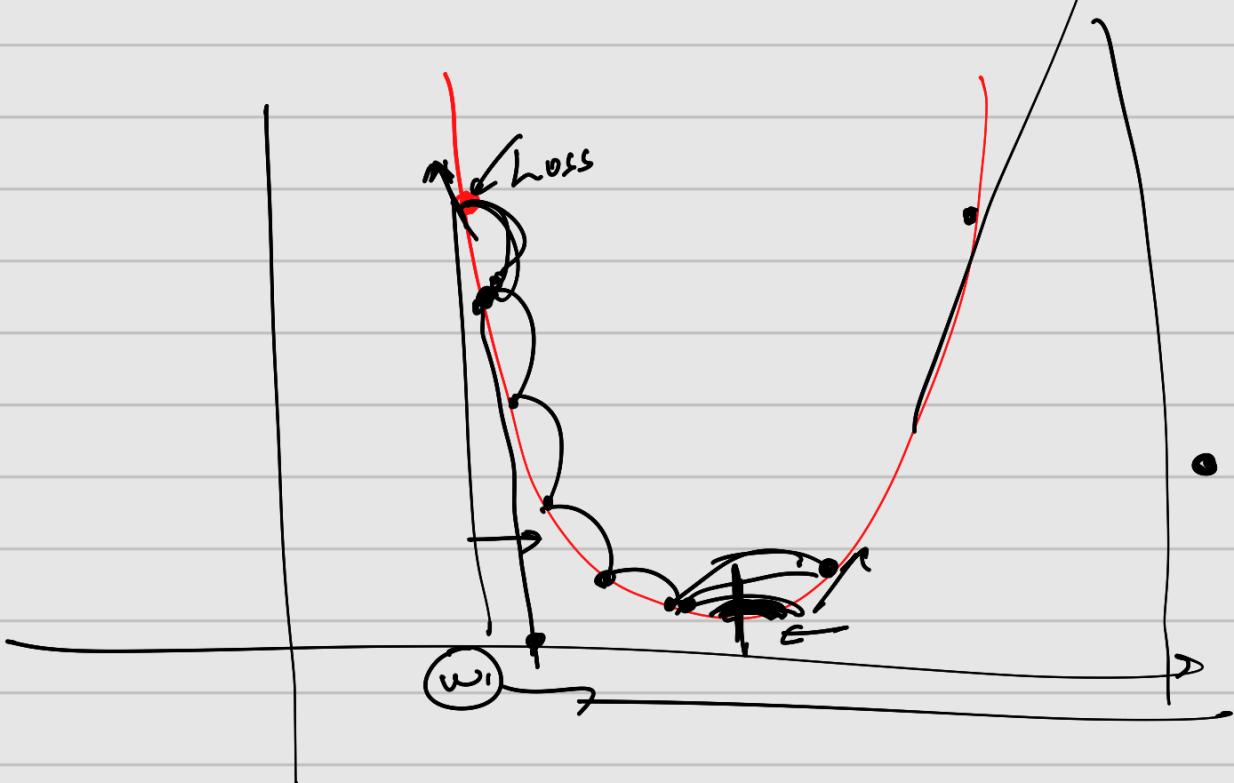


$$w_1(i+1) = \underline{w_1(i)} - \frac{dh}{dw(i)}$$

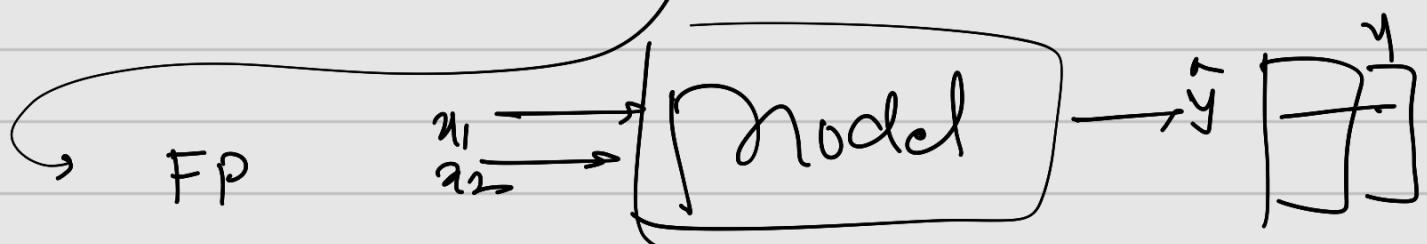
$w_1(i) + \text{amount}$

$$w_1(i+1) = \underline{w_1(i)} - \frac{dh}{dw(i)}$$

$w_1(i) - \text{amount}$

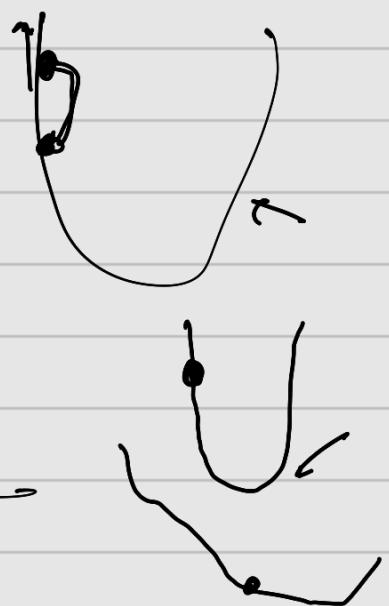


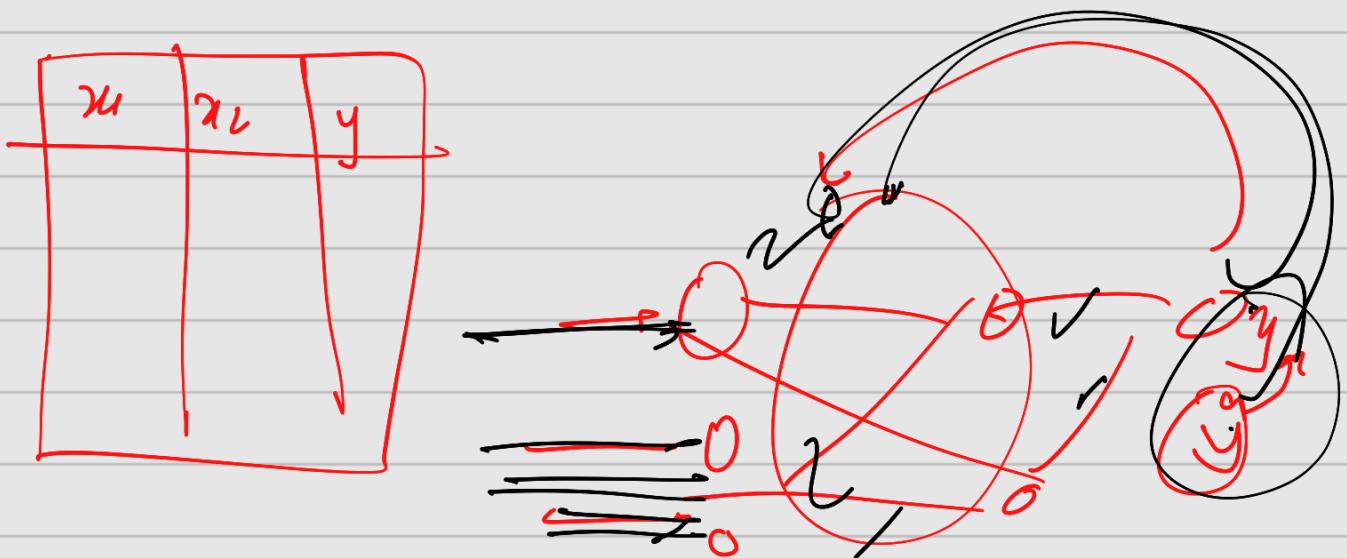
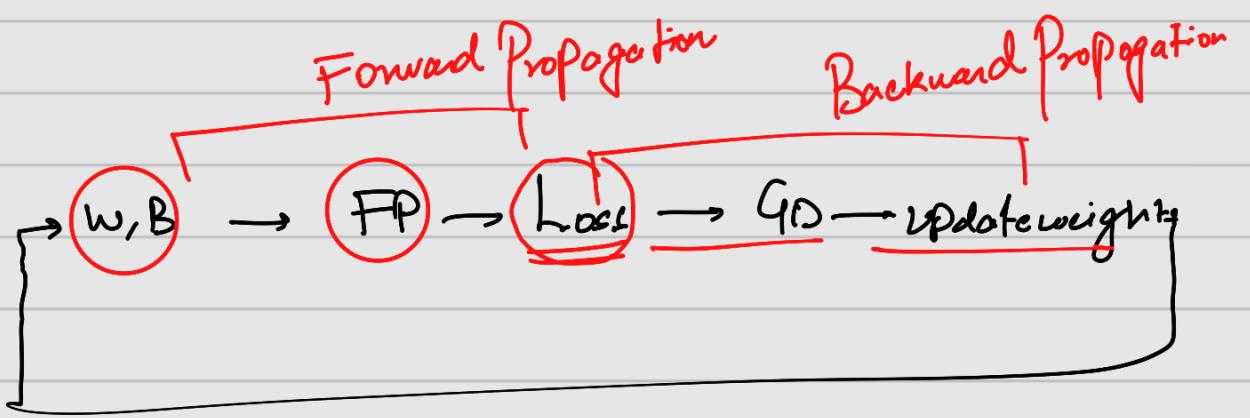
$\mathcal{J}$  will have new set of weights and biases



loss  
gradient descent

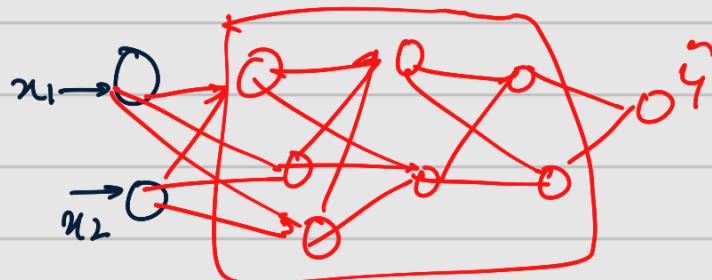
$w_j(i+1) = w_j(i) - \alpha \frac{\partial h}{\partial w_j(i)}$





IFP, IBP  $\rightarrow$  Cepoch.

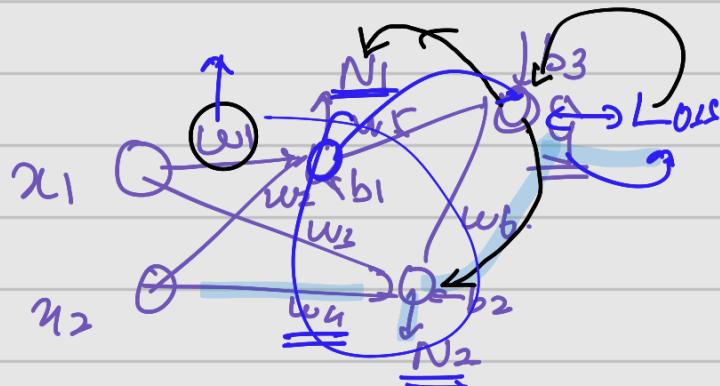
Chain Rule:



$$w_i(t+1) = w_i(t) - \frac{\partial L}{\partial w_i(t)} \quad \text{for all } i$$

$$\left( \frac{\partial h}{\partial w_1}, \frac{\partial L}{\partial w_1}, \frac{\partial h}{\partial w_2}, \frac{\partial L}{\partial w_2}, \dots \right) \left( \frac{\partial L}{\partial b_1}, \frac{\partial L}{\partial b_2}, \dots \right).$$

Regression



$$N_1 = w_1 x_1 + w_2 x_2 + b_1$$

$$N_2 = w_3 x_1 + w_4 x_2 + b_2$$

$$\hat{y} =$$

$$w_5 N_1 + w_6 N_2 + b_3$$

$$\frac{\partial L}{\partial w_4} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial w_4}$$

$$\frac{\partial L}{\partial \hat{y}} = \sum_{i=1}^n \frac{1}{n} \times (y_i - \hat{y}_i)^2$$

$$loss = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$L = f(w, \theta)$$

$$\frac{\partial L}{\partial w} = ?$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial N_1} \times \frac{\partial N_1}{\partial w_1}$$

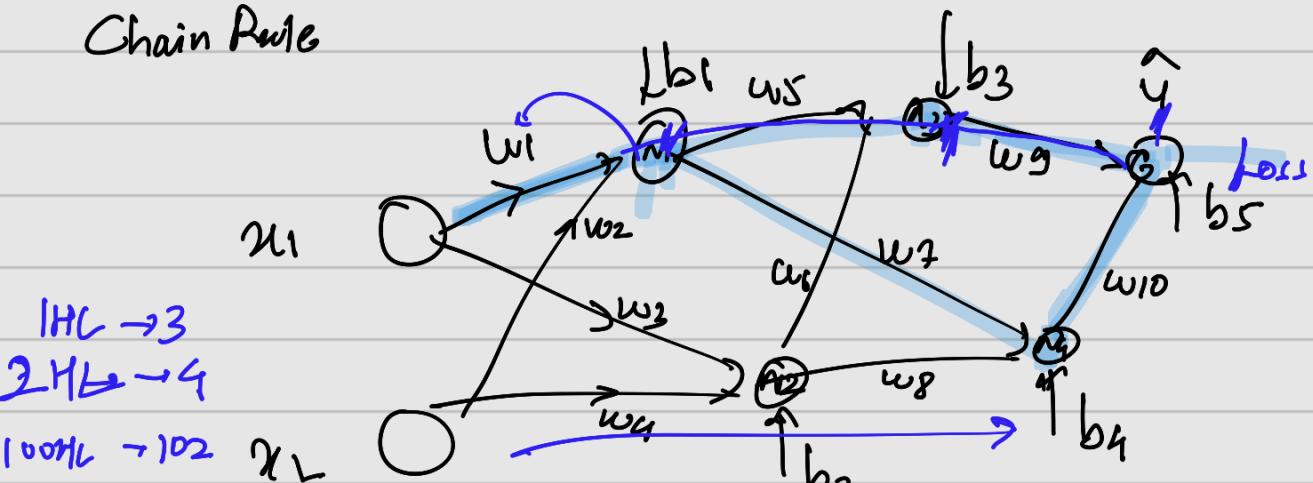
$$\begin{aligned} m &= z^2 \\ z &= m^2 \\ y &= z^2 \end{aligned}$$

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

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial z} \times \frac{\partial z}{\partial m} \times \frac{\partial m}{\partial x}$$

$$\frac{\partial y}{\partial x} = 2z \times 2m \times 2x$$

Chain Rule



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial N_1} \times \frac{\partial N_1}{\partial w_1} + \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial N_2} \times \frac{\partial N_2}{\partial w_1}$$

$$\begin{aligned} (0.99)^{365} \\ = 0 \dots \dots 0 \\ (1.01)^{365} - 1 \end{aligned}$$

Vanishing gradient  
Exploding gradient

## Activations