# Deep Learning
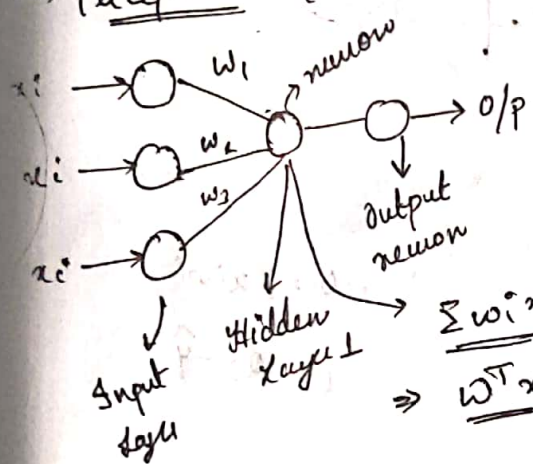
① Forward and Backward propogation
② Loss Function
③ Activation functions
④ Optimizers .

---

→ Deep learning tries to mimic human brain — Multi Layered neural network

→ Deep Learning is becoming popular because of huge amount of data &

    ⓐ Hardware advancement (GPUs)

    ⓑ

→ Perception:- if single Layer neural network &



$\left.\begin{array}{l} w_1 \\ w_2 \\ w_3 \end{array}\right\}$ weights

→ weights say that how much a neuron should get activated or deactivated

$$\sum w_i x_i \}\ \text{step 1}$$
$$\Rightarrow W^T x$$

step 2:- Pass this to a activation function

→ A Bias is added (constant) as if all intercepts weights are 0, then also the training proceeds.

Sigmoid activation function $\Rightarrow \dfrac{1}{1 + e^{-y}} \Rightarrow \dfrac{1}{1 + e^{-(\sum x_i w_i + b)}}$

$\Rightarrow Act(y) \longrightarrow$ then the output

→ This entire process is called forward propogation

→ $(y - \hat{y})$ should be equal to almost 0 . $(y - \hat{y})$ is also called as loss function. If difference is huge we do back propogation.

→ Optimizers will help to update the weights while back propogation.

→ Best known example of optimizer is Gradient descent.

---

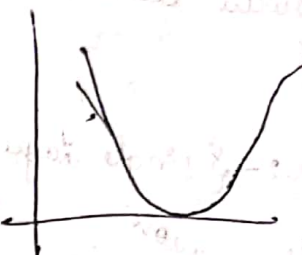① Vanishing gradient problem
② Loss function

## Backpropogation
① weight updation formula
② chain Rule of Differntation

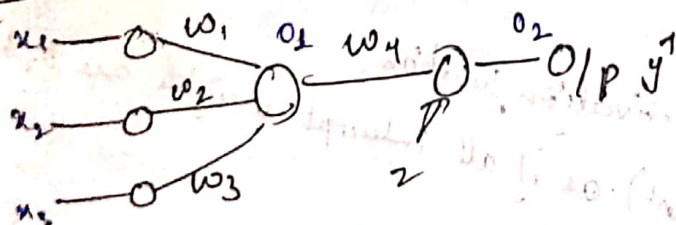$$W_{new} = W_{old} - \eta \frac{\partial L}{\partial W_{old}}$$

↳ Learning Rate

$$\frac{\partial L}{\partial W_{old}} = slope \qquad \eta = 0.01$$

## Chain Rule of Differntiation



$$\frac{\partial L}{\partial W_{4 old}} = \frac{\partial L}{\partial O_2} \times \frac{\partial O_2}{\partial W_{4 old}}$$
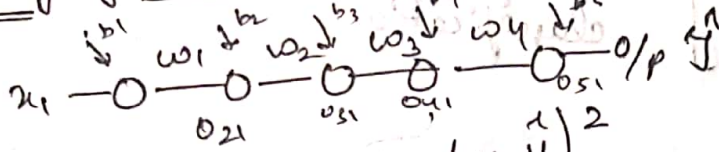
$$W_{4 new} = W_{4 old} - \eta \frac{\partial L}{\partial W_{4 old}}$$

$$b_{2 new} = b_{2 old} - \eta \frac{\partial L}{\partial b_{2 old}}$$

$$\frac{\partial L}{\partial W_{1 old}} = \frac{\partial L}{\partial O_2} \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1 old}$$

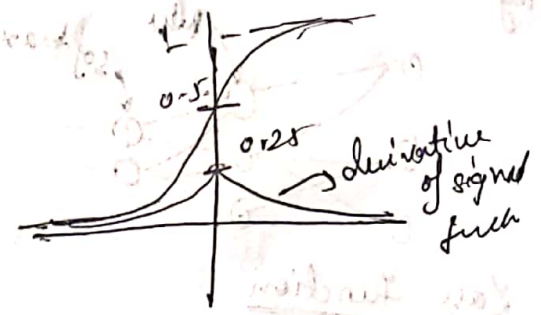$$O_1 = \Sigma x_i w_i$$

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

$$Z = \sigma(O_1 w_4 + b)$$

add both paths

## Vanishing gradient problem

$$x_1 \xrightarrow{\omega_1} \underset{O_{21}}{O} \xrightarrow{\omega_2} \underset{O_{31}}{O} \xrightarrow{\omega_3} \underset{O_{41}}{O} \xrightarrow{\omega_4} \underset{O_{51}}{O} \longrightarrow o/p \; \hat{y}$$

$$LOSS = \frac{1}{2}M(y-\hat{y})^2$$
(Mean squared error)   (MSE)

$$\frac{\partial h}{\partial w_{1new}} = \frac{\partial h}{\partial O_{51}} \times \frac{\partial O_{51}}{\partial O_{41}} \times \frac{\partial O_{41}}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial w_{1new}}$$

$$w_{1new} = w_{1old} - \eta \times \frac{\partial L}{\partial w_{1new}}$$

$$O_{51} = \sigma(O_{41} \times w_4 + b)$$
↳ Activation fn

when   $w_{new} \approx w_{old}$ ⟹ vanishing gradient problem

## Activation functions

① Sigmoid  → o/p layer
② Tanh $\tanh x = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$  → binary classification
③ Relu → for hidden layer
④ Leky Relu   $\max(0,x)$   $d(\tanh x) =$
⑤ Part relu

$\max(0.01 \times x)$

$\max(0.01 \times x)$

Technique which Activation $f^n$ we should use

In hidden layer - Relu ⟶ Pre Relu or Elu

✓ O/p layer - Sigmoid } Binary classification

define



Relu

softmax } Multi class classification

Loss function

Deep Learning          Loss $f^n$          Cost function

① Regression

→ Mean Square Error $\frac{1}{2}(y-\hat{y})^2$          $\frac{1}{2n}\sum_{i=1}^{n}(y-\hat{y})^2$

→ Mean Absolute Error $\frac{1}{2}|y-\hat{y}|$          $\frac{1}{2n}\sum_{i=1}^{n}|y-\hat{y}|$

→ Huber Loss

Advantages
① Differentiable
② only one global/local minima
③ It converges faster

Disadvantages
① Not Robust to outliers

→ Robust to outliers

Huber Loss

when not punit - outliers

$Loss = \begin{cases} \frac{1}{2}(y-\hat{y})^2, & \text{if } (y-\hat{y}) \leq \delta \\ \delta|y-\hat{y}| - \frac{1}{2}\delta^2, & \text{otherwise} \end{cases}$

classification

└→ cross Entropy ───┐
                ├→ Binary cross Entropy
                └→ Catagorical cross entropy

┌─ Same loss function as that of Logistic regression

└→ Do one hot encoding

$$L_{of}(x_i, y_i) = -\sum_{j=1}^{c} y_{ij} \times \ln(\hat{y}_{ij})$$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in class} \\ 0 & \text{othrw} \end{cases}$$
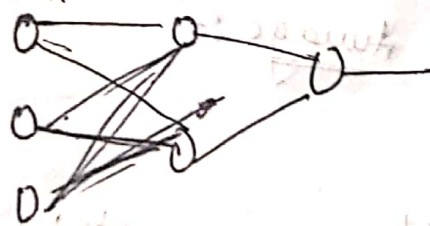
softmax activation

$$\sigma(z) = \frac{e^{z_i}}{\sum\limits_{j=1}^{k} e^{z_j}}$$

Optimizers

→ ① Gradient Descent

weight updation formula:

$$\boxed{W_{new} = W_{old} - \eta \left[\frac{\partial h}{\partial W_{old}}\right]}$$

I/p layer    Hiddenlayer

$$\text{Loss function} = \frac{1}{2n}\sum_{i=1}^{n}(y - \hat{y})^2$$

Disadvantage
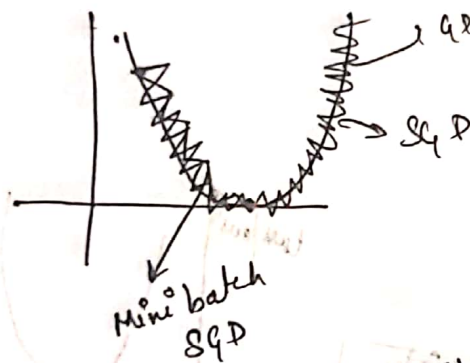
① Resource extensive task (Huge RAM)

Epoch

1 FP and BP → 1 epoch

② Stochastic Gradient Descent

→ We will pass datapoints in iterations.

→ Same like Gradient Descent

→ Ram required is less but the converge will be

→ very slow. High Time complexity

③ Mini Batch Stochastic Gradient Descent

→ In batches, one batch might have 1000 reads

→ Resource intensive

→ Converge rate will be latter

→ Time complexity will improve



Mini batch
SGD

④ SGD with Mini Batch SGD with Momentum

→ Momentum will help to reduce the noise

Exponential Moving/Weighted Average :-

$$w_t = w_{t-1} - \eta \frac{\partial h}{\partial w_{t-1}}$$

$$V_{t_1} = a_1$$

$$V_{t_2} = \beta \times V_{t_1} + (1-\beta) \times a_2$$

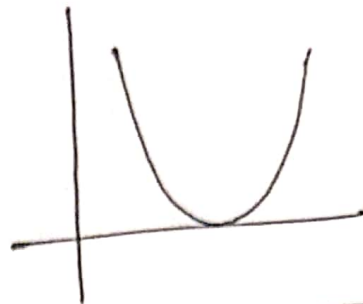$$V_{t_3} = \beta \times V_{t_2} + (1-\beta) \times a_3$$

$\beta = 0$ to $1$
↳ Hyperparameter

$$w_t = w_{t-1} - \eta \, V_{dw}$$

$$\boxed{V_{dw} = \beta \times V_{dw \, t-1} + (1-\beta) \times \frac{\partial L}{\partial w_{t-1}}}$$

⑤ Adagrad — Adaptive Gradient Descent

→ Learning Rate is not constant

$$\boxed{w_t = w_{t-1} - \eta' \frac{\partial L}{\partial w_{t-1}}}$$

$$\eta' = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

$\eta \rightarrow$ decrease it

↘ might become very less

$\boxed{\epsilon = \text{small number to avoid divide by } 0 \text{ condition}}$

$$\alpha_t = \sum_{i=1}^{t} \left(\frac{\partial L}{\partial w_t}\right)^2$$

→ we should see $\alpha_t$ is never so big

⑥ Adadelta and RMS prop

$$\eta' = \frac{\eta}{\sqrt{sdw + \epsilon}}$$

↗ Exponential weighted average

$$sdw_t = \beta \times sdw_{t-1} + (1-\beta)\left(\frac{\partial h}{\partial t}\right)^2$$

$$(sdw)_t = 0$$

⑦ Adam Optimizer

Momentum + RMS prop

$$V_{dw} = \beta \times V_{dw \, t-1} + (1-\beta)\frac{\partial L}{\partial w_{t-1}}$$

$$0. \; w_t = w_{t-1} - \eta' V_{dw}$$