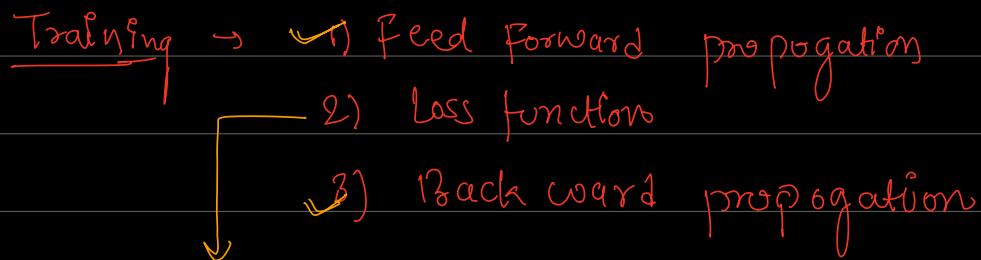


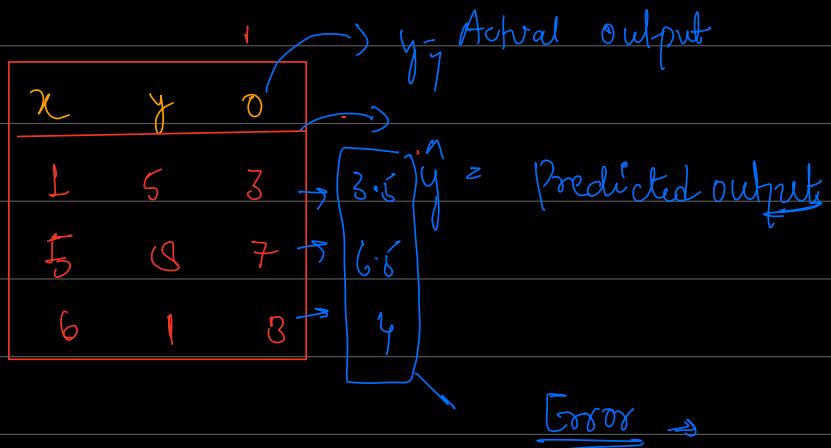
Loss Function

Backpropagation Mathematical Approach -



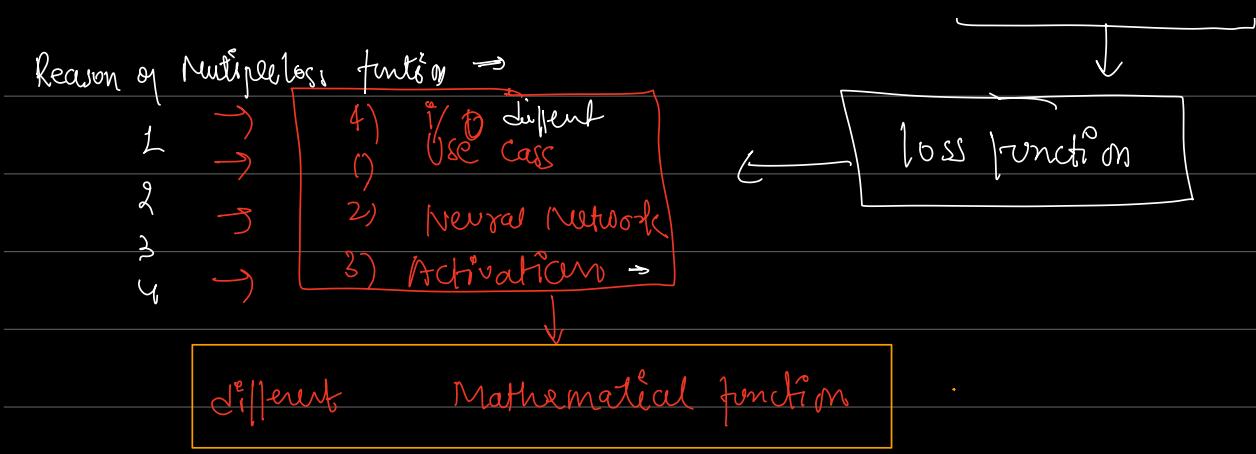
① Loss function →

Def ⇒ It is a method for evaluating the algo's performance →



$$\text{Error} = \boxed{\text{Actual} - \text{predicted}}$$

Neutral = $\boxed{\text{Mines}} \rightarrow \boxed{\text{Error}} \rightarrow \boxed{\text{Mathematical}}$

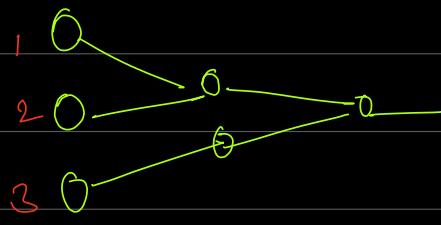


② Types of Loss function

③ Cost function & Loss function →



single iteration



y	x	0
1	2	3
2	3	4

Forward propo

1) Loss function

Error

A-Error

Loss function

Cost function

L

Normal Error

Average Error

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

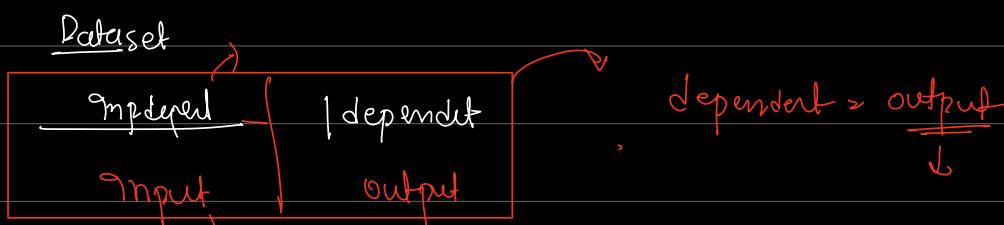
Actual predicted

Types of loss Function

Regression

Classification

- | | |
|--------------------------------|-------------------------------|
| 1) MSE (Mean square Error) | 1) Binary Cross Entropy |
| 2) MAE = (Mean Absolute Error) | 2) Categorical Cross Entropy |
| 3) Huber Loss | 3) Sparse Categorical Entropy |



- ① Numerical = Regression
- ② Categorical = Classification

A	B	C	O
1	90	70	1
2	10	60	3

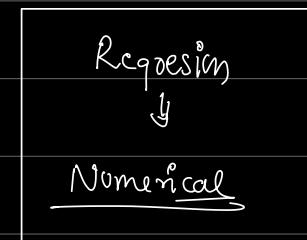
A	B	C	O
1	60	20	Yes
2	50	10	No

3	1	80	80	5
---	---	----	----	---

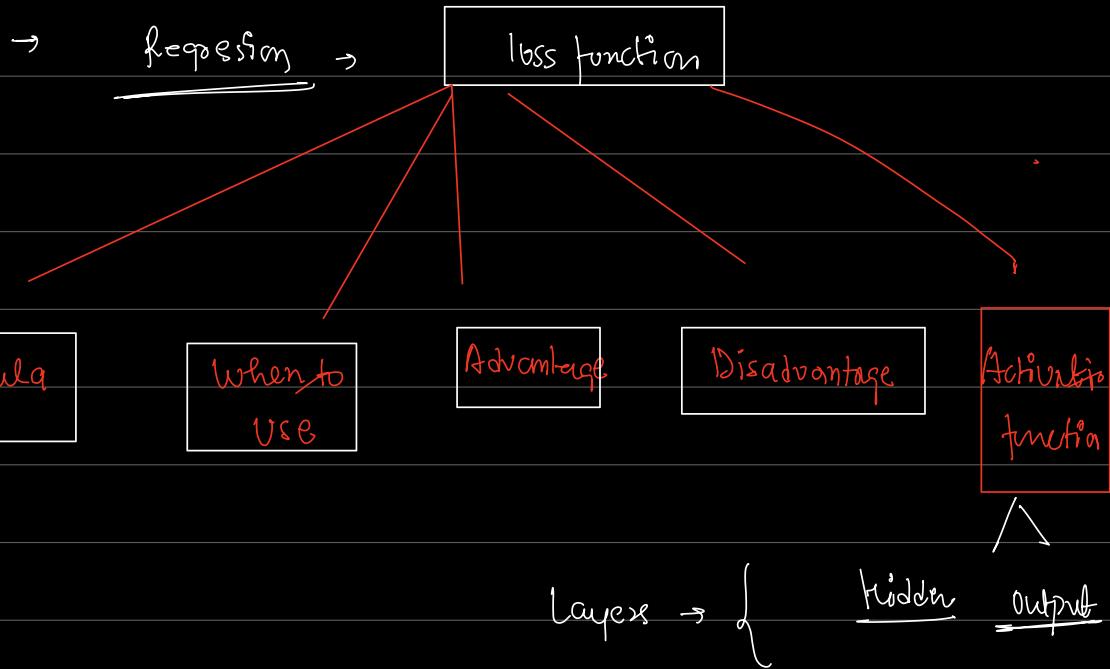
↑

3	1	80	32	Yes	0
---	---	----	----	-----	---

↑



Category



① What is global & local minima ?



$\vec{r} \rightarrow \vec{s}$

\downarrow
global

Loss \rightarrow Convergence = minima error

Global point where error = minima

④ Equation of loss function = Quadratic - global minima

① MSE :

Mean Squared Error \Rightarrow 1) MSE (Mean Square Error)

2) Squared loss

3) L^2 loss

$$\text{MSE} = \frac{1}{N} \sum_i^N (Y_i - \hat{Y}_i)^2$$

↗ ↓ ↓
 No. of value Actual Predicted

Q) Where to use Regression \rightarrow 1) Regression

2) No outlier present



$$\begin{array}{c}
 x \quad y \\
 \boxed{0.2} \quad \boxed{0.4} \\
 (5^1) \quad (0^0) \\
 1 \quad (5) \rightarrow (25) \xrightarrow{\text{expand}} \\
 6 \quad (8) \rightarrow (64)
 \end{array}$$



(2). Differential easily $\Rightarrow \frac{d}{dx}$ \Rightarrow Differentiable



- Advantages \Rightarrow
 - (1) Easy interpretation
 - (2) Always Differentiable
 - (3) Global minima

Disadvantages \therefore

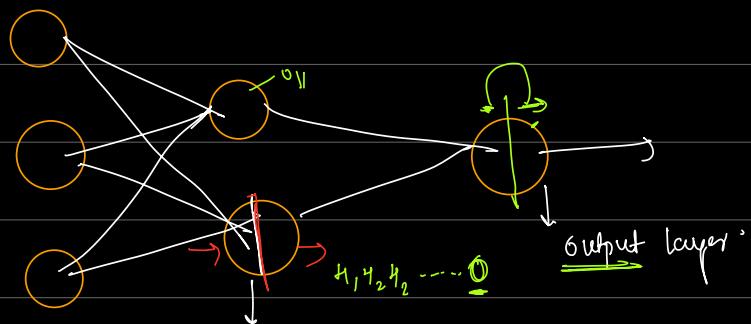
- (1) Not Robust to outliers

(not able to handle)

(2). Sensitive to Overfitting.

Activation function \rightarrow Hidden layer: ReLU

output layer: Linear Activ. funt.



Hidden layers

↓

ReLU \rightarrow No add Non-linearity

↓

to learn Complex pattern

Output linear \rightarrow We use linear because we want the output =

2. MAF

1) Mean Absolute Error

2) LL loss function

$$\text{Cost Function} \rightarrow \text{MAE} = \frac{1}{N} \sum_{i=1}^N |Y_i - \hat{Y}_i|$$

$\downarrow \quad \downarrow$
Actual Predicted

Why ? Mode - () \rightarrow

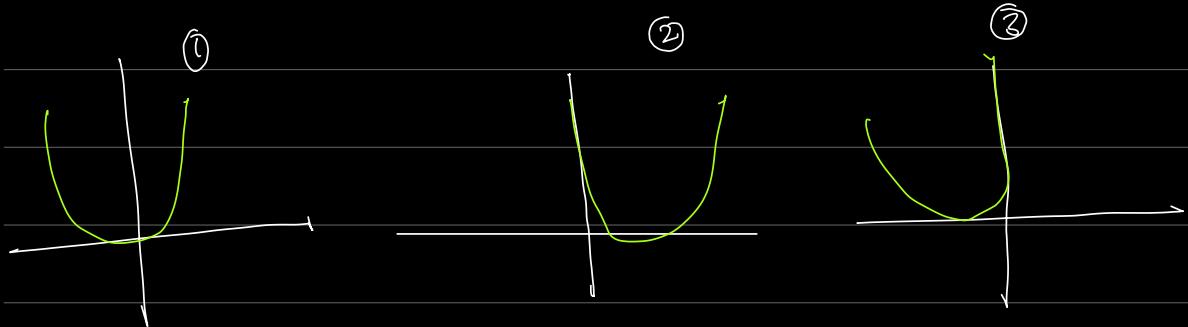
$$\begin{array}{c}
 \text{→} \quad \hat{y} \quad \hat{y} \\
 \left| \begin{array}{cc|c}
 & 3-2 & \Rightarrow L_{\text{wg}} \\
 \text{Model} \quad \leftarrow & 1-4 & -3 \text{ map} \\
 & 5-2 & 3 \text{ map} \\
 \end{array} \right. \quad \text{Actual Error} \\
 \text{---} \\
 1-3 \rightarrow (3) \\
 \left| \begin{array}{c|c}
 1 & \Rightarrow L + 3 + 3 \rightarrow 7 \\
 3 & \\
 3 & \\
 \end{array} \right.
 \end{array}$$

1

Where to use → 1) Dataset containing outlier

(σ) Activation

(3) It is not differentiable



- Advantages → 1) Robust to outliers
2) Easy to interpret

Disadvantage → 1) It is not differentiable

2) Less sensitive to large error

(3) → Not able find ten direction

of graph

Activation function \rightarrow Hidden \rightarrow ReLU (Non-linearity)

Output \Rightarrow Linear

(3) Huber loss \Rightarrow Reason: \Rightarrow

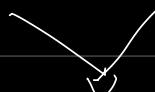
Regression

outlier
 $\xrightarrow{L_{\text{outlier}}}$

MSE \neq MAE

Not Handle

Handle



less sensitive to outliers \Leftarrow [loss function] = Huber loss

- 1) global
- 2) differential
- 3) outlier

$$1) \text{MSE} = \bigcup q_{\text{global}}$$

$$3) \text{MAE} \rightarrow \text{Outlier Handle}$$

Huber loss

Not
When to prefer which \Rightarrow

$\times \rightarrow$ 1) Large dataset \rightarrow
 $\times \rightarrow$ 2) Slow processing \curvearrowright (It is slow)

Formal \Rightarrow

$$\text{Huber} = \begin{cases} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (y_i - \hat{y}_i)^2 & |y_i - \hat{y}_i| \leq \delta \\ \frac{1}{n} \sum_{i=1}^n \delta \left(|y_i - \hat{y}_i| - \frac{1}{2} \delta \right) & |y_i - \hat{y}_i| > \delta \end{cases}$$

if $\leq \delta$
if $> \delta$

n = number of data points

y = Actual Value

\hat{y} = Predicted Value

$\delta =$ Define point where Huber loss
transitions from Quadratic to Linear

Condition = threshold

(1) Advantages \rightarrow 1) It handles the outlier dataset
as well

(2) - It lies b/w MAE and MSE



(3) \rightarrow GT is differentiable \rightarrow global minima

Disadvantages

- 1) GT is slow processing.
- 2) GT is not good for large dataset.
- 3) GT is complex.
- 4) GT need hyperparameter δ

Activation Function \rightarrow Hidden Layer \rightarrow Relu

Output = Linear AF

Classification

- 1) Binary
- 2) Multiclass
- 3) Sparse

① Binary Cross Entropy \Rightarrow ② Log loss function

Where \rightarrow 1) Classification

2) Binary classification

x	y	y
x	y	N



Q)

Cost -

$$\text{Log Loss} = -\frac{1}{N} \sum_{i=1}^N y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)$$

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

Actual Predicted

Why Log ? \rightarrow ① Probabilistic Answers

② Penality for incorrect prediction

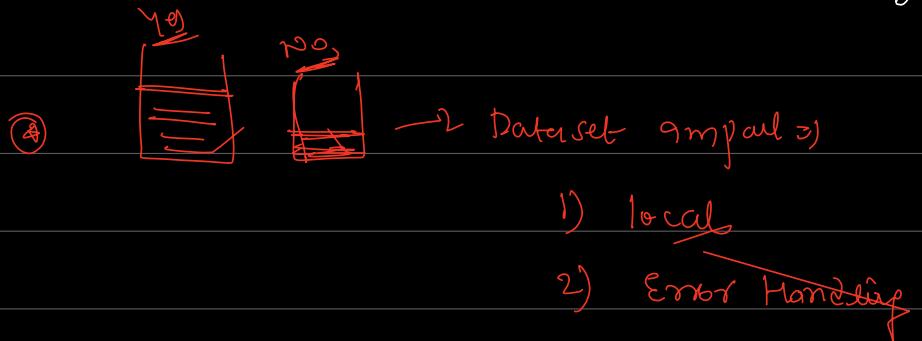
③ Smoother gradient

④ Advantages \rightarrow 1) It penalize incorrect prediction

2) Differentiable

3) for probabilistic interpretation

- * Disadvantage \Rightarrow
 - 1) Problem of local Minima
 - 2) Not Optimal
 - 3) \cdot It depends on class imbalance
(Sensitive to class imbalance)

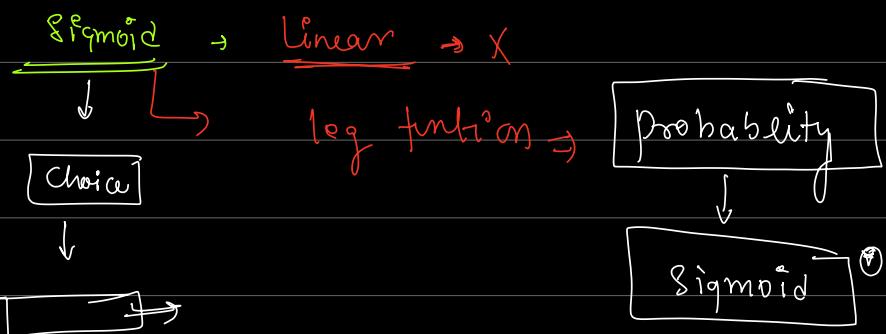


Activation functions \Rightarrow

Hidden layers \rightarrow ReLU \Rightarrow

Output layer \rightarrow Sigmoid

ReLU \rightarrow to add non-linearity.



Q

Classification

Log loss

↳ 1) Binary Cross Entropy \star

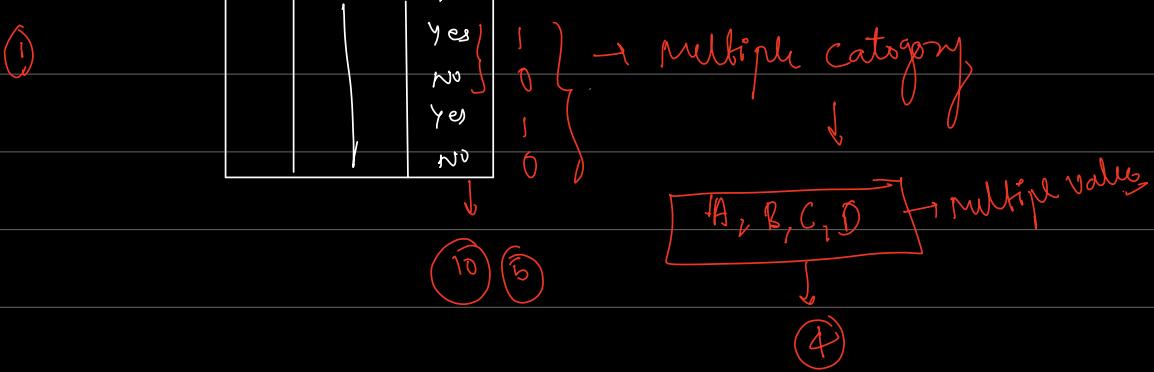
→ 2) Categorical Cross Entropy

3) Sparse Categorical Cross Entropy

Categorical Cross Entropy

1)

1 hot column \rightarrow One hot Encoding



②

Loss = $-\sum_{j=1}^k y_j \log(\hat{y}_j)$

where k is number of classes in the data

Actual
Predicted value
multi

↓
Iteration

④

Sample
no of iteration
class
 $\Rightarrow 2$

Cost = $\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k [y_{ij} \log(\hat{y}_{ij})]$

⑤ One - Hot - Encoding

Yes $\rightarrow \begin{cases} 1 \\ 0 \end{cases}$

One Hot Encoding

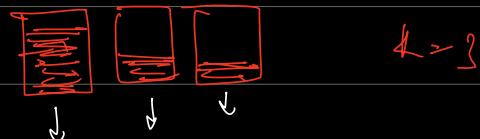
Categorical →

$$\left\{ \begin{array}{l} (1) \rightarrow \text{Formula} = \underline{k} = \text{Class} \\ (2) \rightarrow \text{One Hot Encoding} \rightarrow \end{array} \right.$$

③ Advantages → 1) classify Multiple class

(4) Disadvantages → 1) One-Hot Encoding for the Train

2) sensitive Class Imbalance



Class Imbalance = Accuracy ↓

Activation function → Hidden L → ReLU (Non-linearity)

Output layer → (Softmax) ○

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$



Probability

(1) Sigmoid → Binary classification

(2) Multiple class = Softmax AF

[Sparse] Categorical Cross Entropy ↗

Why sparse? = To remove Empty space
 ⇒ Purpose → Same for Multi class classification

[Sparse Matrix → Empty]

1	0	1	0
1	0	0	1
0	1	0	0

Sparse = Value
Empty = empty

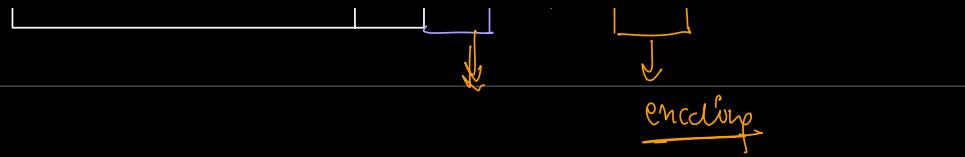
	1	2	3	4
y_g	1			
No	0			
R_{ij}	0			
HNo	1			

sparse

Endless

sparse

sparse



Basic \rightarrow Way of doing Encoding

which \rightarrow Sparse ✓

1 Formula = Same

Advantage \Rightarrow - (1) \rightarrow No one hot Encoding
 (2) Computation efficient
 than Categorical

Disadvantage \rightarrow 1) It is less intuitive as compare
 to categorical

(2) Sensitive \rightarrow Class imbalance

Which to prefer \rightarrow Intuitive \rightarrow Categorical

Computation power \Rightarrow Sparse

Gradient Descent

→ Optimization
Algorithm

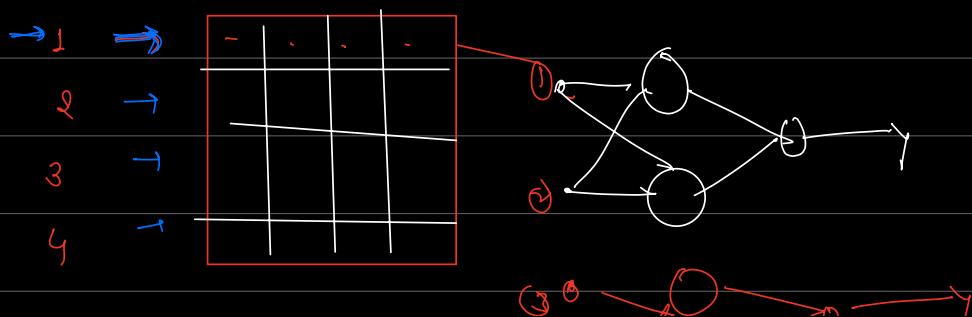
- 1) Term and Terminology ↗
- 2) Why we use gradient ↗
- 3) Working ↗
- 4) Types →
 - 1) Batch Gradient Descent ↗
 - 2) Stochastic Gradient Descent ↗

3) Mini Batch Gradient Descent ↗

- 5) Compare ↗ Speed →
Convergence stability
- 6) When to use which one ↗

(1) Terms & Terminology →

- 1) epoch →



④ θ

\leftarrow \rightarrow

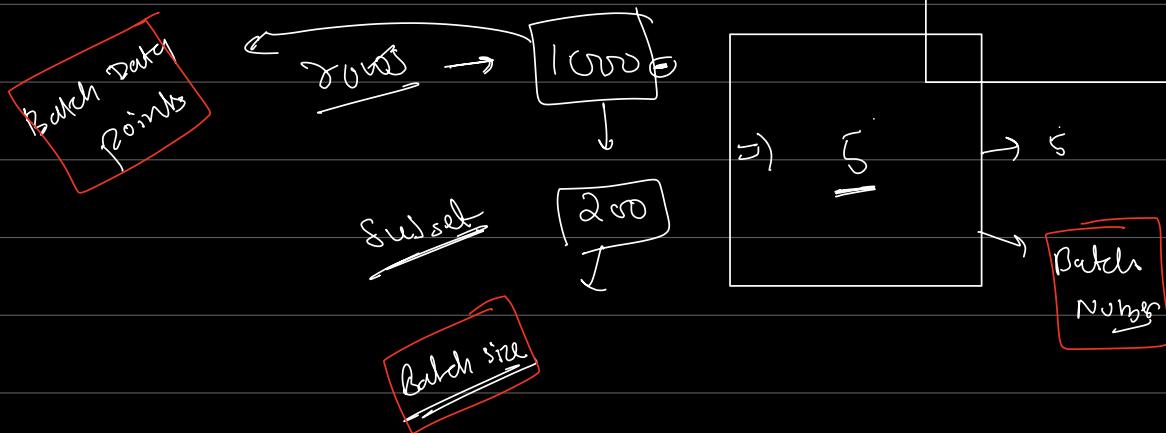
① \rightarrow For 1 Row

② $F_p + B_p \rightarrow$ step function

② \rightarrow Complete dataset $\Rightarrow F_p + B_p \rightarrow$ Epoch

③ \rightarrow Batch \Rightarrow collection of Iteration \Rightarrow

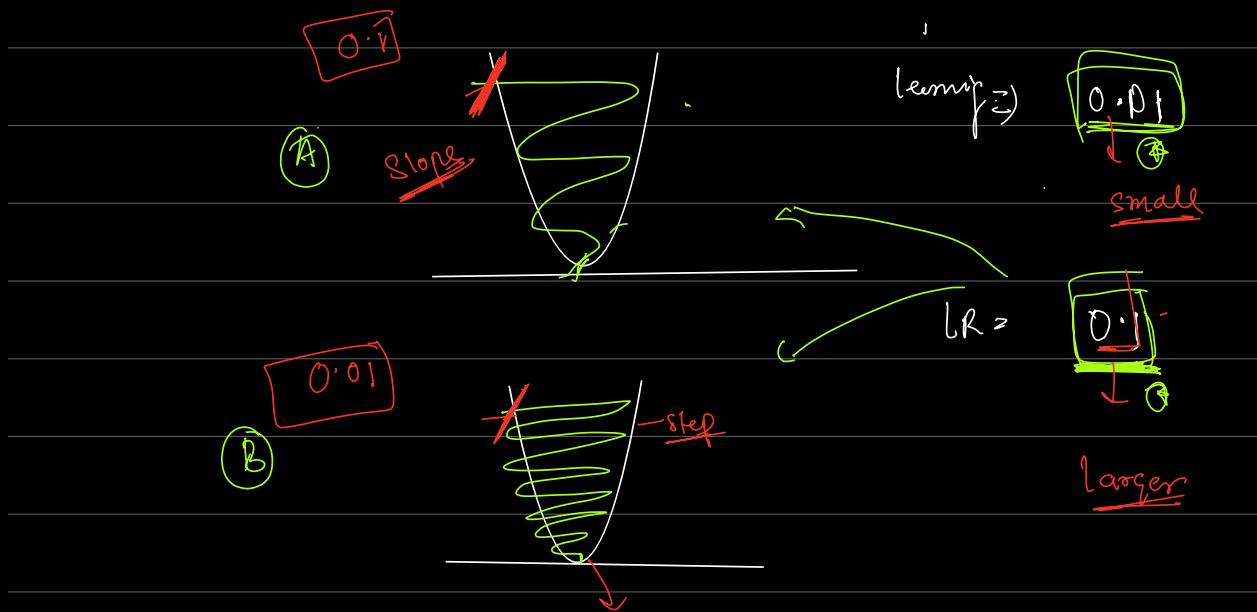
dataset \Rightarrow 100 Data point \Rightarrow 5
20



④ Learning Rate \Rightarrow if is step size Optimal value

Weird & bias

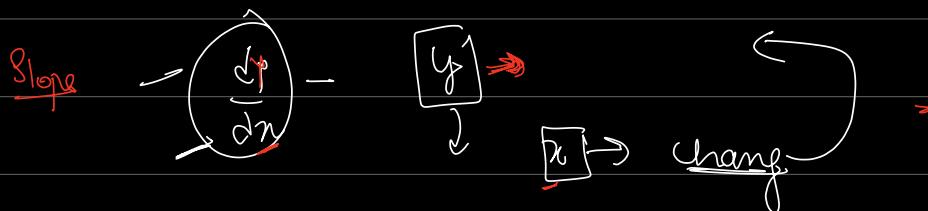
How fast & slow Model learn.



\Rightarrow gradient descent \Rightarrow optimization
algo

Slope \Rightarrow 1) direction

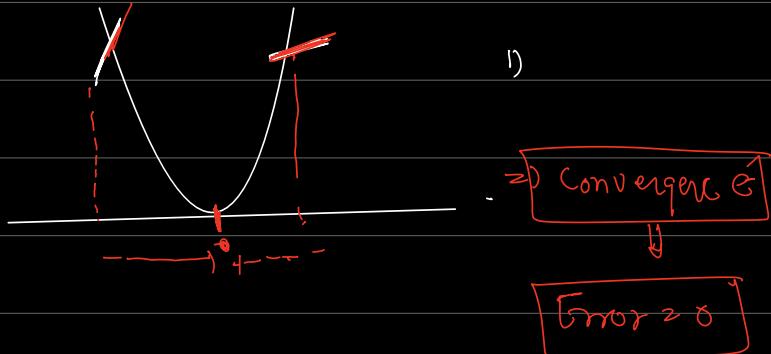
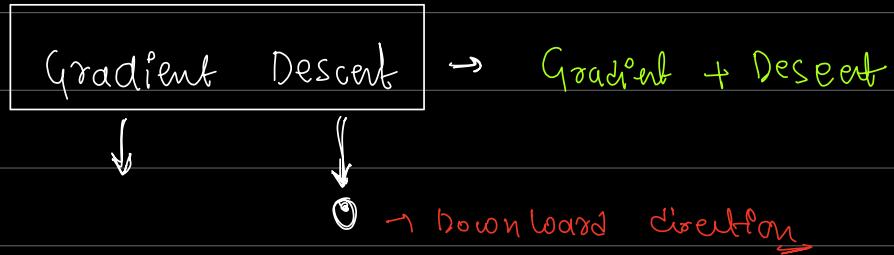
$$\text{Slope} \rightarrow 2) \frac{dy}{dx} = \frac{y_2 - y_1}{x_2 - x_1}$$



$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial y_i} \times \frac{\partial y_i}{\partial w_i} \quad \left. \right\} \Rightarrow \frac{\text{slope}}{\parallel} \quad \boxed{\text{gradient}}$$

gradient

$$(2) \rightarrow \text{Cost function} = \frac{1}{2} \sum ()$$



Why Gradient Descent :-

1) Convergence - ✓ 1) It helps to Convergence.

2) Optimization: ✓ 2) Optimization Algorithm → Try to find optimal value

④ Scalability → ④ 3) It can be used better for linear & Non-linear models.

efficiency - ④ → It is computationally efficient.

High dimensional
data set

Efficiency for large dataset.

④ Gradient Descent = ① optimization Algorithms
which is used to

minimize the cost function:

loss function ①

Working \Rightarrow

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

(1) Weight / Bias \rightarrow Randomly

(2) loss \rightarrow Error

(3) Back propagation \Rightarrow Weight / Bias

To reduce the Error

L/D

Adjust \rightarrow optimize

Algo \Rightarrow

Gradient descent

② Increases / Decreases the
Value of w / b

Weight Bias Update Formula

Optimization \Rightarrow Error Reduce

1) Gradient Descent - Types

2) —

3) —

4) —

5) —

6) —

Gradient Types

Reason

Batch size

Batch Gradient
Descent

Stochastic
Gradient Descent

Mini Batch
Gradient D.s.

epoch -

Complete Data

single datapoint

Collection
of
dataset

