

MACHINE LEARNING

AYUSH KUMAR SA

Data mining :- Applying ML techniques to dig deep into large amounts of data to help discover patterns that were not immediately apparent.

Why Machine learning is important?

Machine learning is great for

- * Problems for which existing soln require a lot of hand tuning or long lists of rules.
ex → Spam Mail Detection

- * Complex problems for which no good soln is present
using traditional approach
ex → Speech Recognition

- * Fluctuating techniques environments

ex → Change in approach of spammers who send spam mail

- * getting insights about complex problem and large amounts of data, eg - clustering

Labels → Solⁿ given with the features. Only sol^{n^{value}} are known as labels.

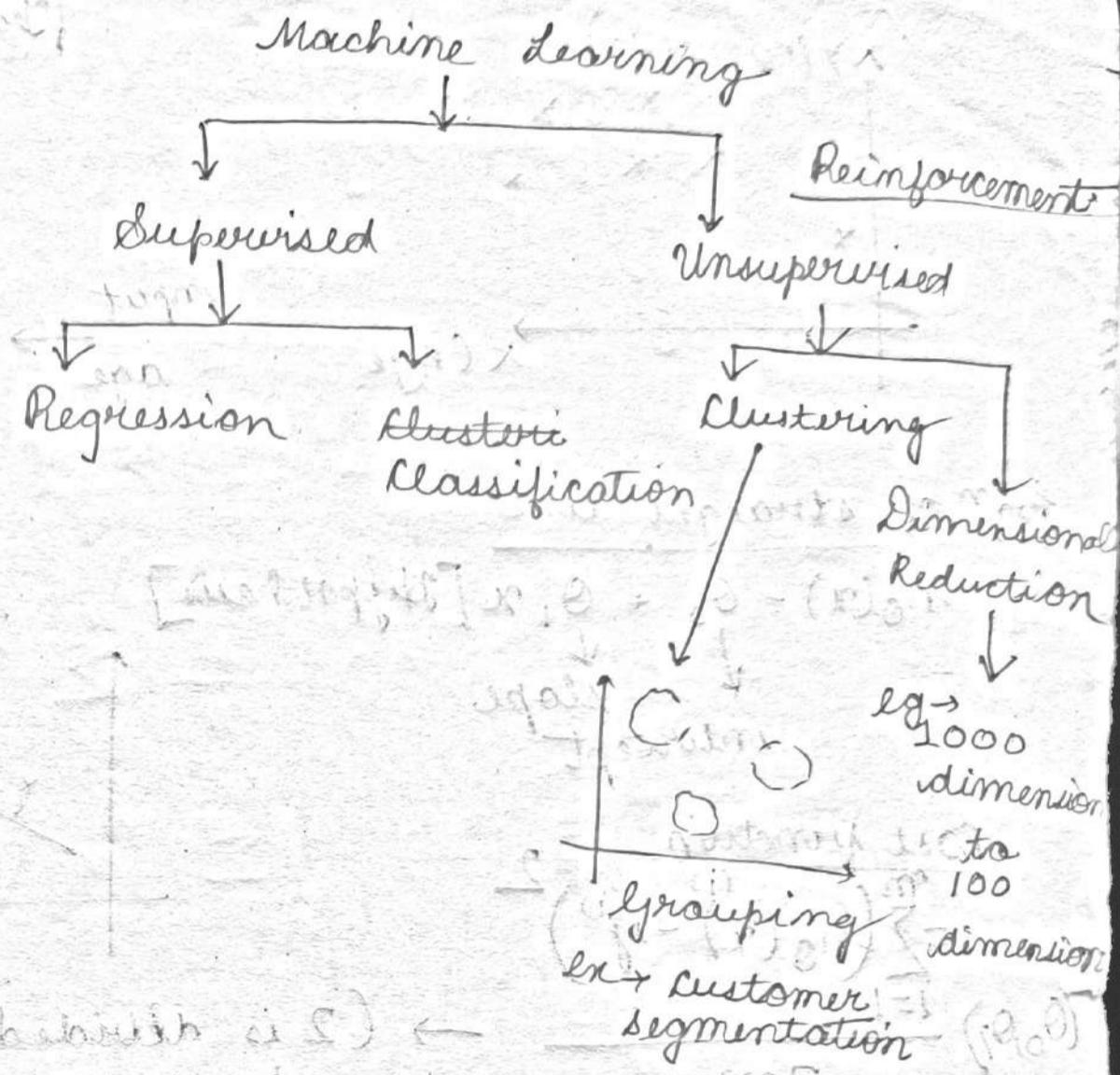
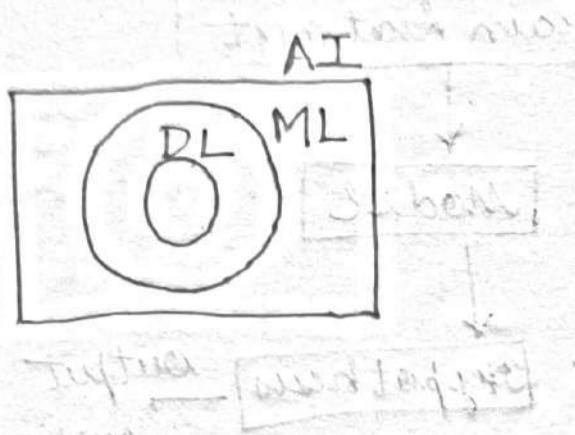
Attribute vs features (used interchangeably)

↓ ↓
eg: mileage eg: mileage is 5,000

Unsupervised learning also help in learning association

rule learning (to ~~not~~ identify relation between attributes).

Deep learning and Machine Learning



Supervised learning Algorithms

- 1) Linear Regression
- 2) Ridge and Lasso
- 3) Logistic Regression
- 4) Decision Tree
- 5) AdaBoost
- 6) Random Forest
- 7) Gradient Boosting
- 8) Xgboost
- 9) Naive Bayes
- 10) SVM
- 11) KNN

Unsupervised learning algorithms

- 1) K Means
- 2) DB Scan
- 3) Hierarchical
- 4) K Means & Nearest Neighbours clustering
- 5) PCA
- 6) LDA

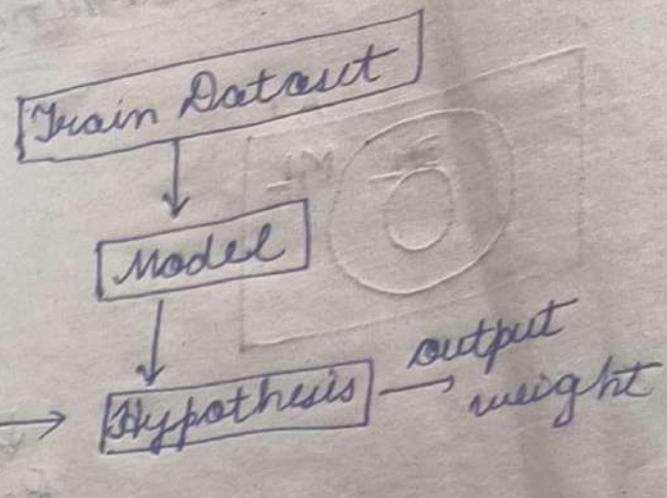
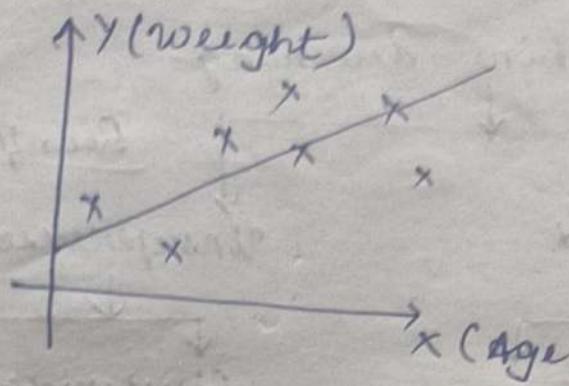
$$E_{\text{PCA}} = \frac{1}{2} \|X - \hat{X}\|_F^2 + \lambda \|W\|_F^2$$

minimum $(10, 0)$

$$\theta = \theta_0 + (\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n)$$

① Linear Regression

ALGORITHMS



Eqn of straight line

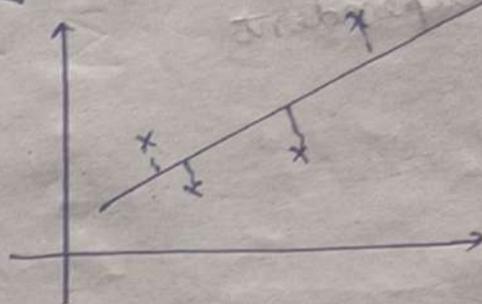
$$h_0(x) = \theta_0 + \theta_1 x \quad [\text{Hypothesis}]$$

↓ ↓
slope intercept

Cost function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

→ (2 is divided so that its derivative will be simpler.)



$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

Squared Error Function

What we need to solve

$$\min_{(\theta_0, \theta_1)} \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$$

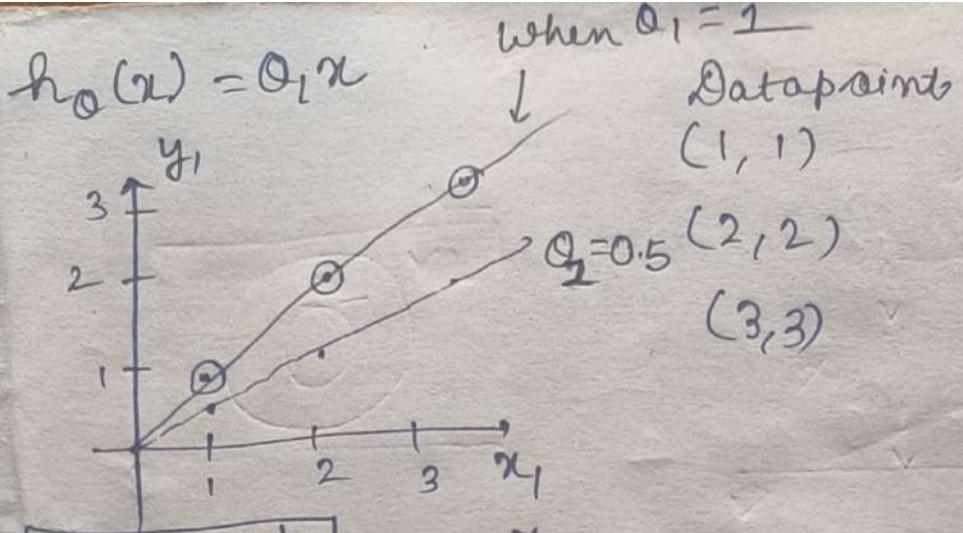
or

$$\min_{(\theta_0, \theta_1)} J(\theta_0, \theta_1)$$

$$\Rightarrow h_0(x) = \theta_0 + \theta_1 x \quad \text{if } \theta_0 = 0$$

If $\theta_0 = 0$

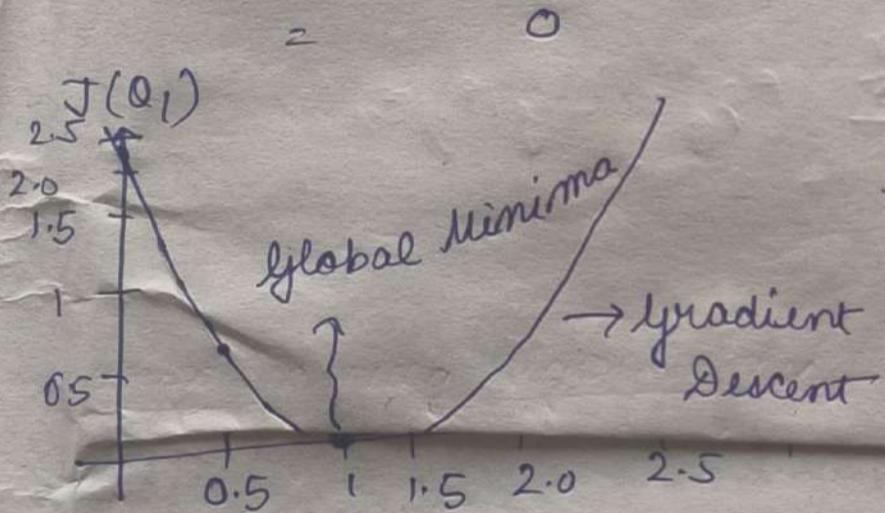
$$h_0(x) = \theta_1 x$$



When $\theta_1 = 1$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} [(1-1)^2 + (2-2)^2 + (3-3)^2]$$



When $\theta_1 = 0.5$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m [0.5-1]^2 + [1-2]^2 + [3-1.5]^2$$

$$= \frac{1}{2 \times 3} [0.25 + 1 + 2.25]$$

$$\approx 0.58$$

When $\theta_1 = 0$

$$J(\theta_1) = \frac{1}{2m} [1^2 + 9^2 + 4^2] = \frac{1}{6} [14] = 2.3$$

Convergence Algorithm

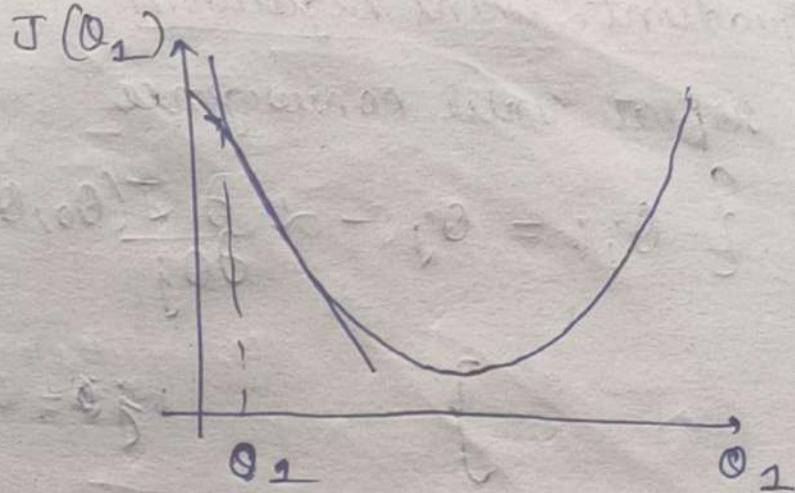
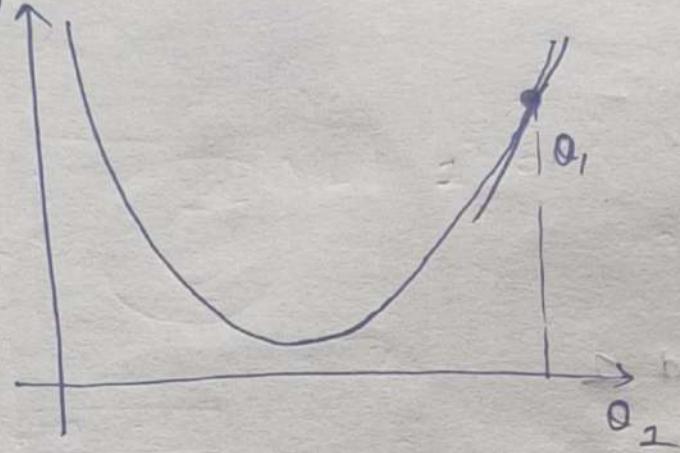
Repeat until convergence

$$\left\{ \theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \right\}$$

learning rate

Derivative (slope)

$\alpha = 0.01$



~~for~~

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta, \theta_i)}{\partial \theta_j}$$

In this case $\frac{\partial J(\theta, \theta_i)}{\partial \theta_j} = +ve$

$$\therefore \theta_j = \theta_j - \alpha (+ve)$$

\therefore so, we should reach global minima

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta, \theta_i)}{\partial \theta_j}$$

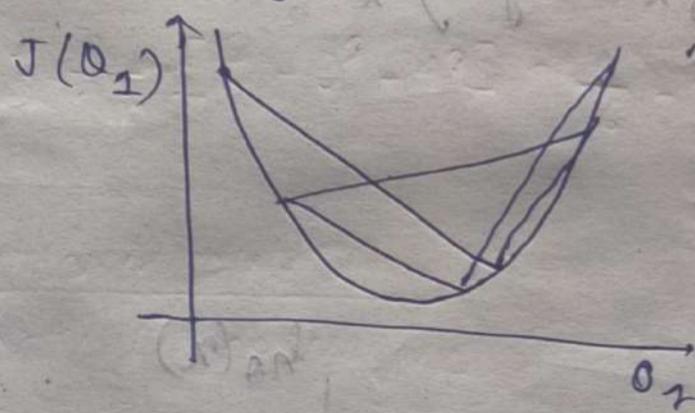
In this case $\frac{\partial J(\theta, \theta_i)}{\partial \theta_j} = -ve$

$$\therefore \theta_j = \theta_j - \alpha (-ve)$$

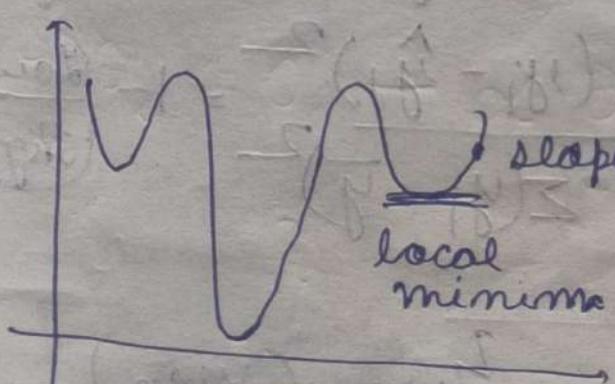
\therefore so, we should reach global minima.

α (learning rate) should be less, so updation of θ_j should be in steps, (like $\theta \leftarrow \theta - 0.01$)

If $\alpha = 100$ eg,



we will not get any local minima



In deep learning techniques we may get such gradient Descent graph which can solved by various optimizers, but in linear regression, we will ~~not~~ not get such graph

Gradient Descent Algorithm

repeat until convergence

$$\left\{ \theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \right.$$

$j = 0 \text{ and } 1$

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$j = 0 \Rightarrow \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

$$j = 1 \Rightarrow \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

repeat until convergence

$\alpha = 0.01$ (Learning rate)

$$\left\{ \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \right. \quad \begin{matrix} \text{(at each iteration)} \\ \text{move towards zero} \end{matrix}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)} \quad \begin{matrix} \text{move towards zero} \\ \text{move towards zero} \end{matrix}$$

minimum local

Performance metrics

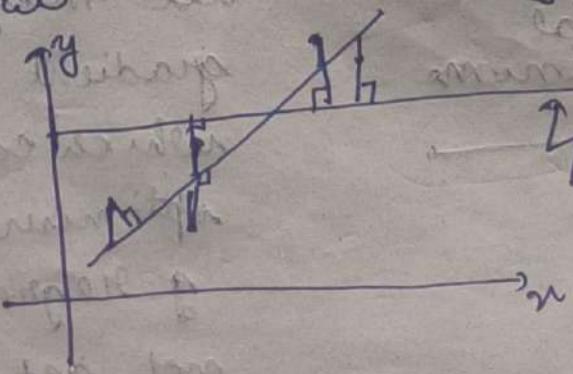
R^2 and Adjusted R^2

$$R^2 = 1 - \frac{SS_{res}}{SS_{total}}$$

$$1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\text{Error}}{\text{Total}}$$

Small
number
of points
High
error

* R^2 can also be negative



Worse (Mean)

* R^2 also increases when attributes nowhere related to the output are added. So we use adjusted R^2 .

Adjusted R^2

$$(R^2 \text{ (adjusted)} = 1 - \frac{(1-R^2)(N-1)}{N-p-1})$$

$$\boxed{\text{Adjusted } R^2 \leq R^2}$$

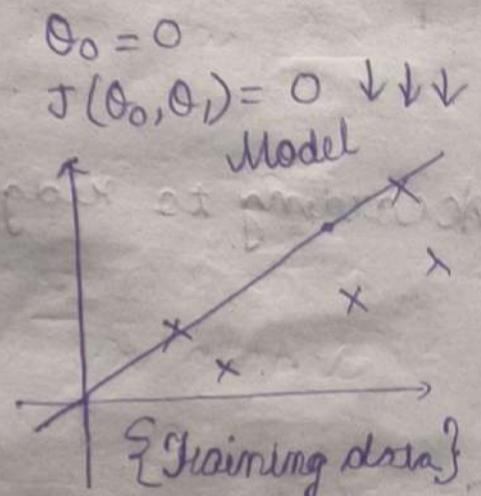
* N = no. of data points

* P = no. of predictors

Ridge and Lasso Regression

Overfitting

- ① Model performs well \rightarrow Training data (low bias)
- ② Fail to perform well \rightarrow Test Data (high variance)



Underfitting { High bias, high variance }

- ① Model accuracy is bad with training data.
- ② Model accuracy is also bad with test data.

ex →

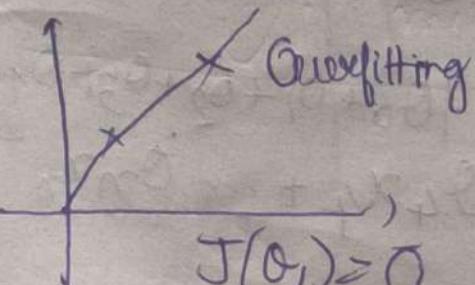
Model 1

Training Acc = 90%

Test acc = 80%

Overfitting

Low bias, high variance



Model 2

Training acc = 92%

Test acc = 91%

Generalised

Low bias, low variance

Model 3

Training acc = 70%

Test acc = 65%

Underfitting

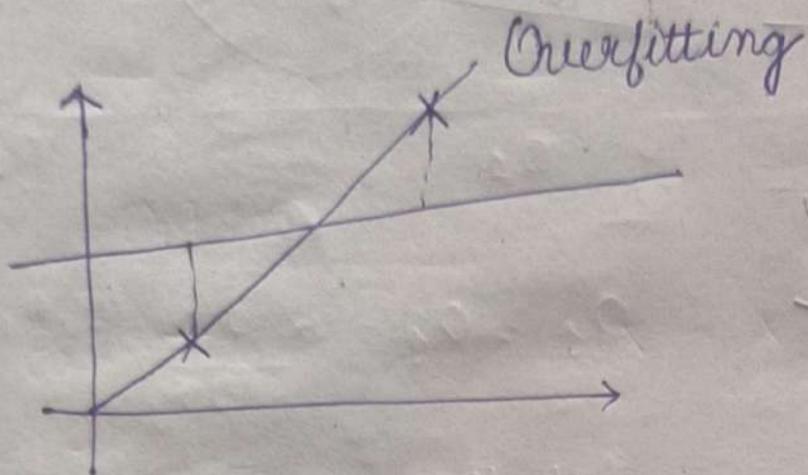
High bias, high variance

$$J(\theta_1) = 0$$

$$= \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Ridge (L2 Regularization)

$$J(\theta_1) = (\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope})^2$$



(i-11) (-2-1)

i-9-11

$$J(\theta_1) = 0$$

$$= \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$$= (\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope})^2$$

According to ridge we add $[\lambda (\text{slope})^2]$

$$= 0 + 1(2)^2 = 4$$

$$= (\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope})^2$$

$$\left\{ \begin{array}{l} \text{Prevents Overfitting} \\ \end{array} \right\} = \text{small value} + 1(1.5)^2$$

$$= \text{small value} + 2.25$$

iterations { hyperparameter }

R^2 , adjusted R^2

$\lambda \rightarrow$ hyperparameter

convergence algorithm

Lasso (L1 Regularization)

$$J(\theta_1) = (\hat{y}^{(i)} - y^{(i)})^2 + \lambda |\text{slope}| \quad \xrightarrow{\text{feature selection}} |0_0 + \theta_1 + \theta_2 + \theta_3 + \theta_4 + \dots + \theta_m|$$

$$h_0(x) = \hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \dots + \theta_n x_n$$

Using Lasso

- ⇒ Preventing overfitting } $\rightarrow L_1$ regularization
- ⇒ Feature selection } $\rightarrow L_2$ regularization
- {⇒ cross-validation}

Summary

Ridge regression (L2 norm)

$$\text{Cost function} = \frac{1}{2m} \left(\sum ((h_\theta(x)) - y^{(i)})^2 + \lambda (\text{slope})^2 \right)$$

Purpose: Preventing overfitting

Lasso regression (L1 norm)

$$\text{Cost function} = \frac{1}{2m} \left(\sum ((h_\theta(x)) - y^{(i)})^2 + \lambda |\text{slope}| \right)$$

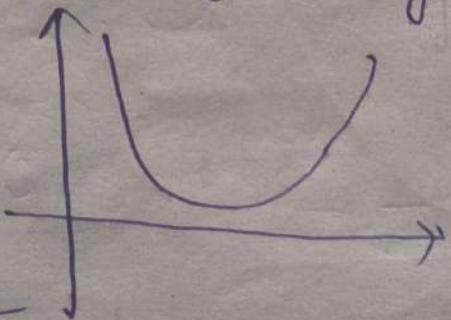
Purpose:

- 1) Prevent overfitting
- 2) Feature selection

Assumption of Linear Regression

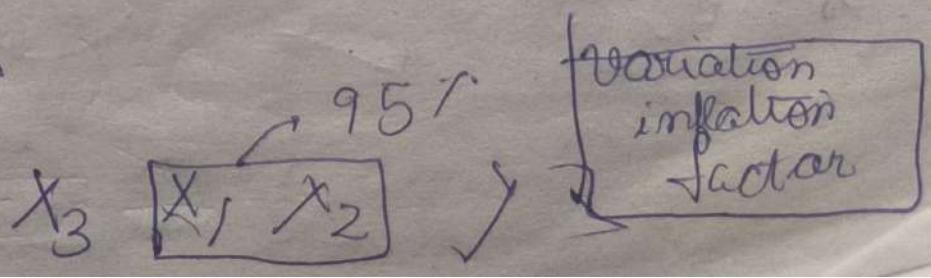
⇒ Normal (Gaussian Distribution \rightarrow model will get trained well)

⇒ Standardization {Scaling data} \rightarrow Z score $\mu=0, \sigma=1$



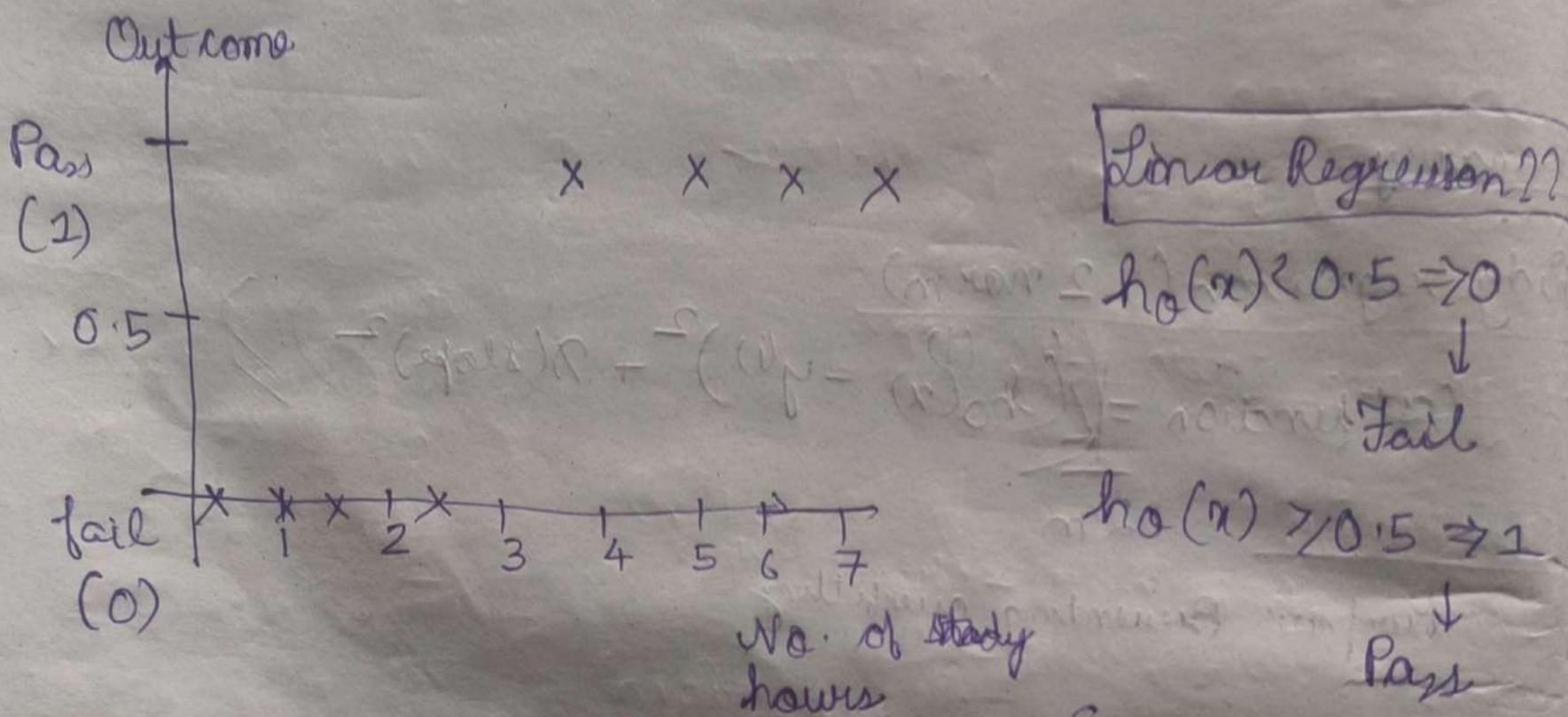
* Homoscedasticity

⇒ Linearity



⇒ Multi-collinearity

Logistic Regression (Classification algorithm) \rightarrow Binary classification
 or
 No. of study hours No. of play P/P
 — —
 — —
 — —
 P
 F
 (Also can solve multiclass classification)



* Why we don't use linear regression?

→ because, in linear regression, values less than 0 and more than 1 also occur due to which it does not become a sigmoid function (values should be between 0 to 1)

* Why don't

→ Because of outliers, the predicted values are far off from actual values.

Decision boundary logistic regression

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

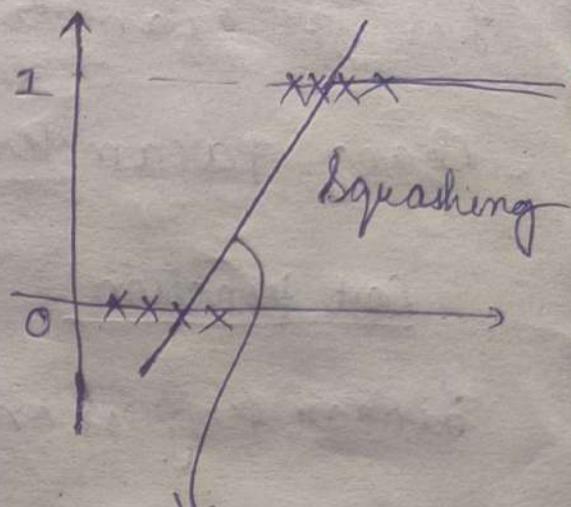
$\boxed{h_0(x) = \theta^T x}$

$$\text{Let } z = \theta_0 + \theta_1 x$$

$$h_0(x) = g(z)$$

$$h_0(x) = \frac{1}{1 + e^{-z}}$$

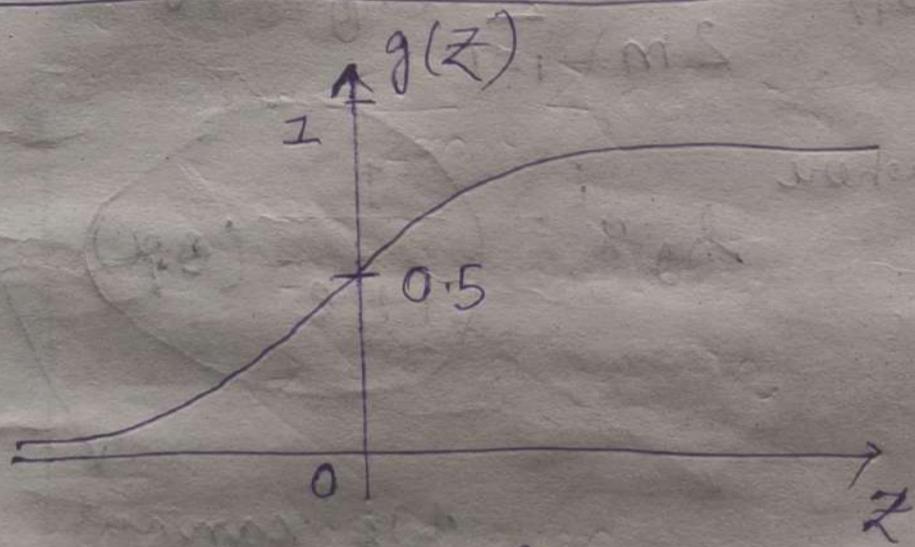
→ sigmoid
or
logistic function



$$h_0(x) = \theta_0 + \theta_1 x$$

$$h_0(x) = g(\theta_0 + \theta_1 x)$$

$$\boxed{h_0(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}}$$



$$\left. \begin{array}{l} g(z) \geq 0.5 \\ \text{when } z \geq 0 \end{array} \right\}$$

Training set

$$\{(x^1, y^1), (x^2, y^2), (x^3, y^3), \dots, (x^n, y^n)\}$$

$$y \in \{0, 1\} \rightarrow \underline{\underline{2 \%}}$$

$$h_0(z) = \frac{1}{1+e^{-z}}$$

$$z = \theta_0 + \theta_1 x$$

Let for example $\theta_0 = 0$, $z = \theta_1 x$

Change parameter θ_1 ?

Cost function

$$\text{Linear regression } J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0 x^i - y^i)^2$$

$$\text{Logistic regression, } h_0 x = \frac{1}{1+e^{-(\theta_1 x)}}$$

Cost function,

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0 x^i - y^i)^2$$

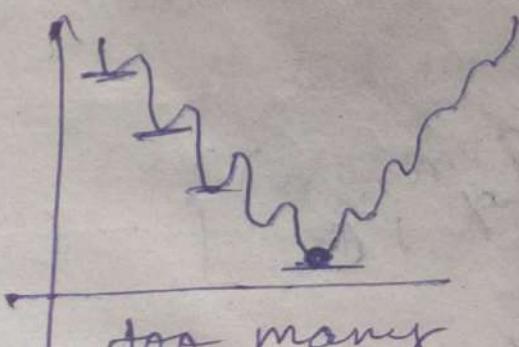
where

$$h_0 x^i = \frac{1}{1+e^{-(\theta_1 x^i)}}$$

non-convex
 f^n .

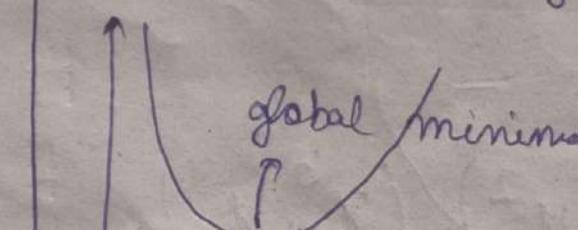
We cannot
use this cost
 f^n for logistic
regression.

Non convexⁿ



too many
local minima

Convexⁿ



global minimum

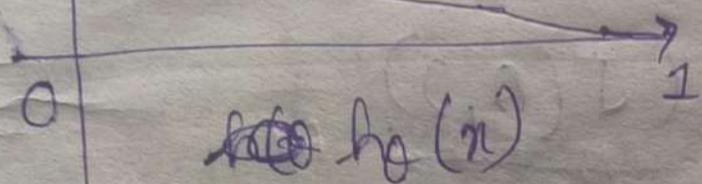
Logistic Regression cost function

$$J(\theta_1) = \begin{cases} -\log(h_\theta(x^{(i)})) & y=1 \\ -\log(1-h_\theta(x^{(i)})) & y=0 \end{cases}$$

$$h_\theta(x) = \frac{1}{1+e^{-(\theta_1 x)}}$$

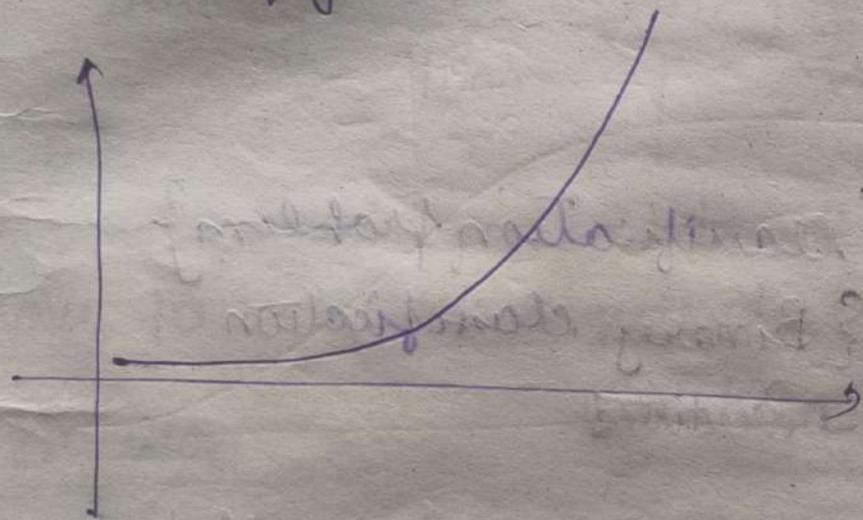
$J(\theta_1)$

if $y=1$ cost = 0 if $y=1$, $h_\theta x = 1$



~~$h_\theta(x)$~~

if $y=0$



$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

$\text{Cost}(h_\theta(x^i), y) = -y \log(h_\theta(x^i)) - (1-y) \log(1-h_\theta(x^i))$

if $y=1$

$$\text{Cost}(h_\theta(x^i), y) = -\log(h_\theta(x^i))$$

If $y = 0$

$$\text{cost}(h_0(x^i), y) = -\log(1 - h_0(x^i))$$

Cost function

$$J(\theta_1) = -\frac{1}{2m} \sum_{i=1}^m y^i \log(h_0(x^i)) + (1-y^i) \log(1-h_0(x^i))$$

$$h_0(x^i) = \frac{1}{1 + e^{-\theta_1 x}}$$

Repeat until convergence

$$\left\{ \theta_j = \theta_j - \alpha \frac{\partial J(\theta_1)}{\partial \theta_j} \right\}$$

Performance metrics {Classification problem}

{actual} {Binary classification}

$$x_1 \quad x_2 \quad \hat{y} \quad y$$

$\hat{y} \rightarrow$ Predicted

Confusion Matrix

		Actual	
		0	1
Predicted	0	3	2
	1	1	1

Predicted

		Actual	
		1	0
Predicted	1	TP	FP
	0	FN	TN

→ Confusion Matrix

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$= \frac{3+1}{3+2+1+1} = 4/7 = 0.57 = 57\%$$

0 → 900 } Unbalanced
1 → 100 } Dataset

0 → 600 } Balanced
1 → 400 } Data

$$\text{Accuracy} = \frac{900}{1000} = 90\%$$

but biased data

More parameters required for unbalanced data

① Precision

$$\frac{TP}{TP + FP}$$

② Recall

$$\frac{TP}{TP + FN}$$

③ F-Score

		Actual	
		1	0
Predict	1	TP	FP
	0	FN	TN

{ Spam classification }
↳ Precision
{ has CANCER OR NOT }
→ Recall

{ Tomorrow stock market will fall } \rightarrow f. Beta

$$\underline{F\text{-beta}} = \frac{(1+\beta^2) \text{ Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

$\checkmark \quad \beta = 1 \quad F_1 = \frac{2 (\text{Pre} \times \text{Recall})}{\text{Pre} + \text{Recall}} \quad [\text{HM}]$

both FP, FN

'imp'

β decrease \rightarrow more imp to FP than FN

$$\rightarrow 0.5 \quad F_{0.5}$$

more FP

dominated by

β increase \rightarrow more imp to FN than FP

$$\rightarrow 2 \quad F_2$$

$$F_2 = \frac{\text{Pre}}{\text{Pre} + \text{FN}}$$

stab based on
stab based on high frequency and low imp

stab 7

stab 8

NAIVE BAYES INTUITION

{Classification}

↓
{Bayes Theorem}

Independent event → rolling a dice

Dependent event → taking out marble without replacement

For dependent event,

$$P(A \text{ and } B) = P(A) * P(B/A)$$

$$P(B \text{ and } A) = P(A \text{ and } B)$$

$$P(A) * P(B/A) = P(B) * P(A/B)$$

$$P(y|A) = \frac{P(B) * P(A/B)}{P(A)} \rightarrow \text{Bayes Theorem}$$

Crux behind
Naive Bayes

ex on Dataset:-

→ Independent feature

$$x_1 \quad x_2 \quad x_3 \quad x_4 \quad \dots \quad x_m \quad \text{TREAT} \quad \boxed{y} \quad \text{output}$$

$$P(y/x_1 x_2 x_3 \dots x_n) = \frac{P(y) * P\left(\frac{x_1 x_2 \dots x_n}{y}\right)}{P(x_1 x_2 \dots x_n)}$$

$$= P(y) * P(x_1/y) * P(x_2/y) * P(x_3/y) \dots P(x_n/y)$$

$$P(x_1) * P(x_2) * P(x_3) \dots P(x_n)$$

ex →

$$\underline{x_1} \quad \underline{x_2} \quad \underline{x_3} \quad \underline{x_4} \quad \underline{y}$$

$$\boxed{}$$

$$P(y = \text{Yes}/x_i)$$

Yes

$$= P(\text{Yes}) * P(x_1/\text{Yes}) * P(x_2/\text{Yes}) \\ * P(x_3/\text{Yes}) * P(x_4/\text{Yes})$$

$$\boxed{}$$

No

$$P(x_1) * P(x_2) * P(x_3) * P(x_4)$$

Similarly for $P(y = \text{No}/x_i)$

$$P(y=\text{No}/x_1) = \frac{P(\text{no}) * P(x_1/\text{no}) * P(x_2/\text{no}) * P(x_3/\text{no}) * P(x_4/\text{no})}{P(x_1) * P(x_2) * P(x_3) * P(x_4)}$$

↙
#fixed

Let $P(\text{Yes}/x_1) = 0.13$ $P(\text{No}/x_1) = 0.05$

$\geq 0.5 \Rightarrow 1$

$< 0.5 \Rightarrow 0$

$$P(y_e/x_1) = \frac{0.13}{0.13 + 0.05} = 0.72 = 72\%$$

$$P(\text{No}/x_1) = 1 - 0.72 = 0.28 = 28\%$$

→ DATASET — ①

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Outlook (x_1)

	yes	No	$P(Y)$	$P(N)$
Sunny	2	3	$2/9$	$3/5$
Overcast	4	0	$4/9$	$0/5$
Rain	<u>3</u>	<u>2</u>	$3/9$	$2/5$
Total	9	5		

Temp (x_2)

	yes	No	$P(Y)$	$P(N)$
Hot	2	2	$2/9$	$2/5$
Mild	4	2	$4/9$	$2/5$
Cold	<u>3</u>	<u>1</u>	$3/9$	$1/5$
Total	9	5		

Play (x_3)

		$P(yes)$	$P(No)$
yes	9	$9/14$	$5/14$
No	<u>5</u>		
Total	14		

Then, for a test dataset, (sunny, Hot) \rightarrow Output?

$$\Rightarrow P(\text{yes} / \text{sunny, hot}) = \frac{P(\text{yes}) * P(\text{sunny} / \text{yes}) * P(\text{hot} / \text{yes})}{P(\text{sunny}) * P(\text{hot})}$$

$P(\text{No} / \text{sunny, hot})$

$$= \frac{1 - 9/14 * 3/5 * 2/5}{1} = 0.085$$

$$= 9/14 * 2/9 * 2/9 = 2/63 = 0.031$$

$$P(\text{Yes} / \text{Sunny, Hot}) = 0.031 = 1 - 0.73 = 27\%$$

$$P(\text{No} / \text{Sunny, Hot}) = \frac{0.085}{0.031 + 0.085} = 0.73 = 73\%$$

Normalized.

$\rightarrow (\text{Sunny, hot}) \rightarrow \text{Yes or } \boxed{\text{NO}}$

Assignment :-

(Overcast, Mild) \rightarrow Naive Bayes?

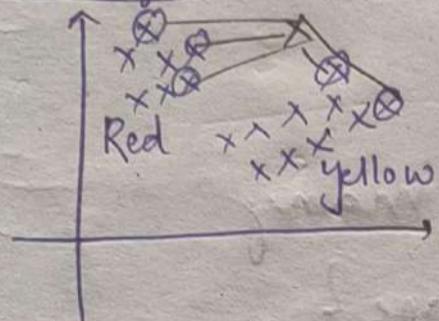
KNN Algorithm

{ K - Nearest Member }

Classification

Regression

① Classification



K-Nearest Neighbour

$$K = 5$$

Maximum No

$$\text{Red} = 3$$

$$\text{Yellow/white} = 2$$

Distances

Euclidean distance

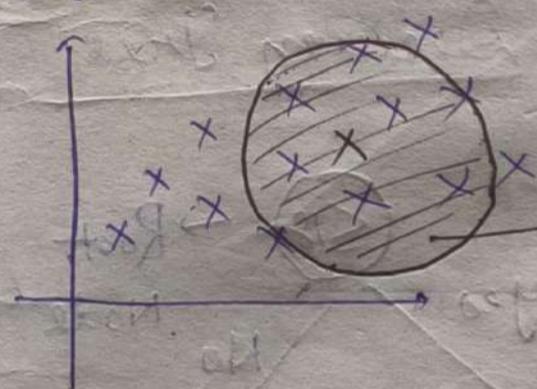
$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$(x_1, y_1)$$

Manhattan Distance

$$|(x_2 - x_1)| + |(y_2 - y_1)|$$

② Regression

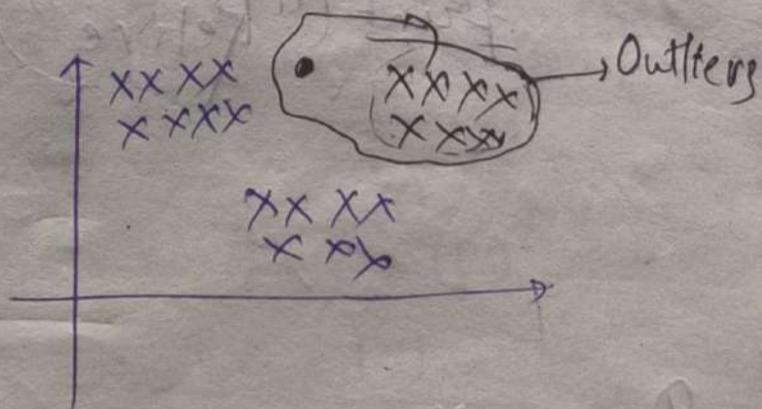


$$K = 5 \leftarrow \text{Hyperparameter}$$

K Nearest Works bad with

① Outlier

② Imbalanced Dataset



Minimise \leftarrow the function \leftarrow

Agenda

- ① Decision Tree Classification ✓
- ② Decision Tree Regression ✓
- ③ Practical Implementation ✓
- ④ Ensemble Techniques ✓

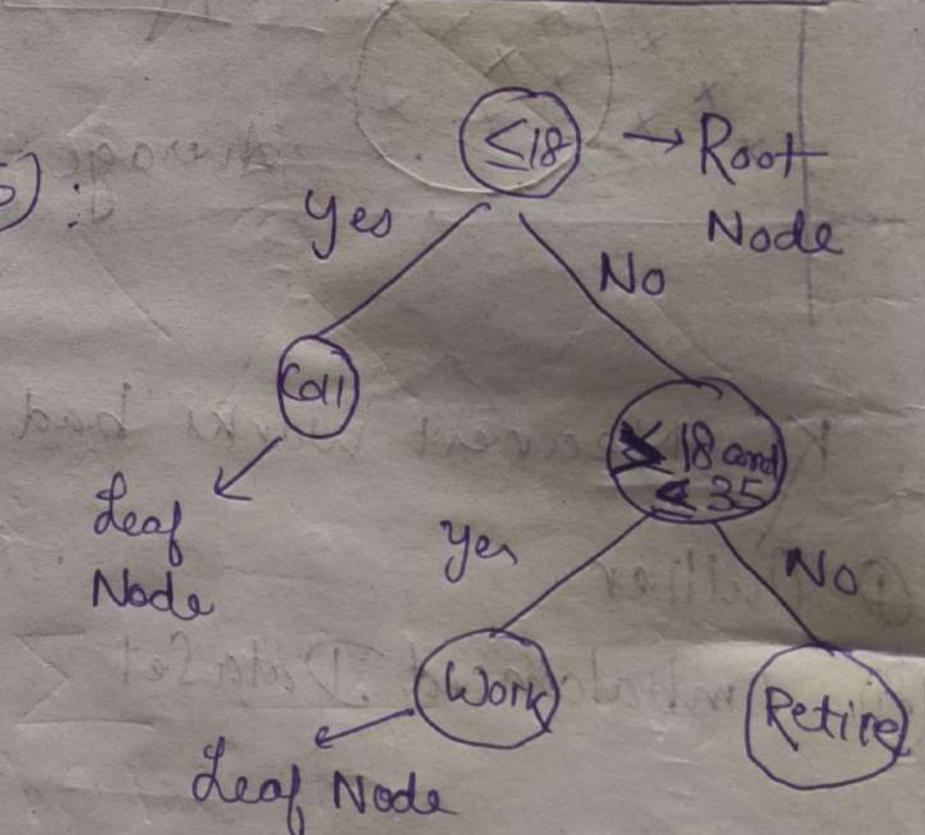
Decision Tree { Solving many usecase }

Regression Classification

```

if (age <= 18):
    print("College")
else if (age > 18 and age <= 35):
    print("Work")
else:
    print("Retire")
  
```

Decision Tree

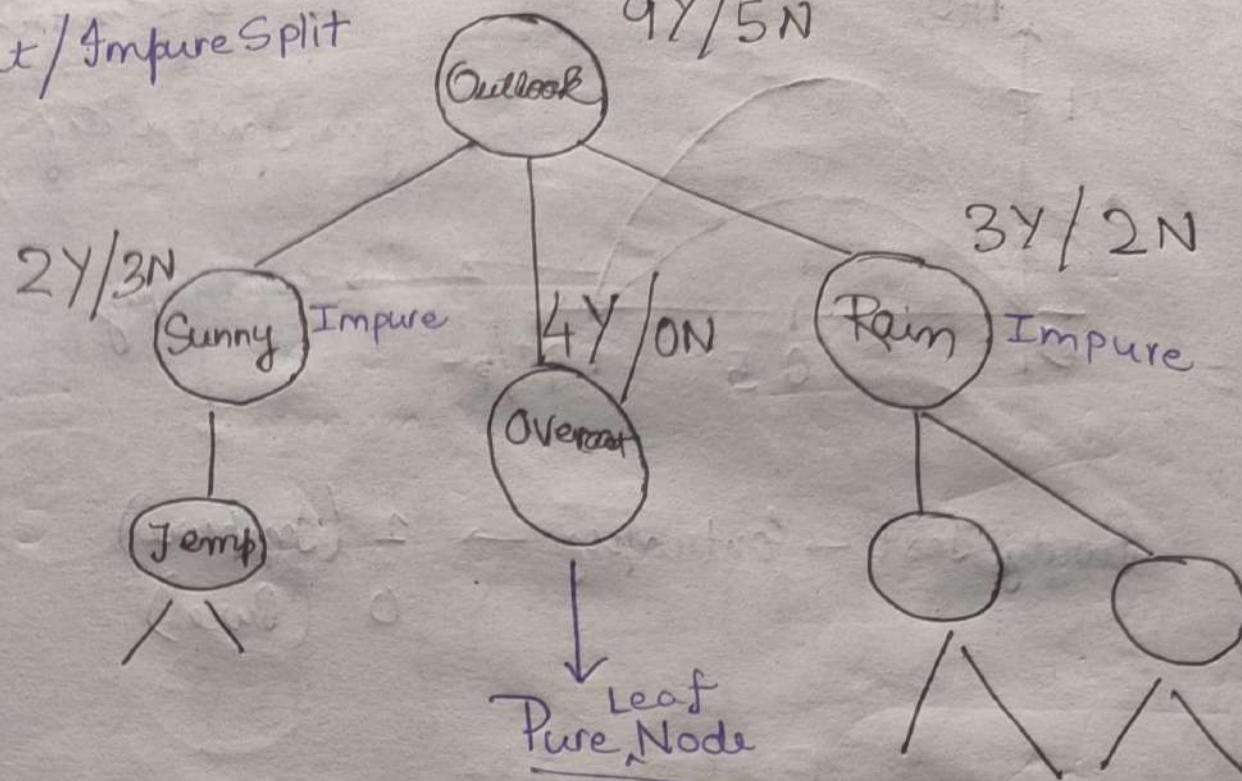


Decision Tree

→ Nested if else \Rightarrow Decision Tree

Refer to Dataset - 1

Pure split / Impure Split



1) How to know purity of a given node? (Mathematically) *Leaf node Pure.*

} Entropy
} Gini Coefficient
} Impurity

2) How the features are selected?

} Information gain }

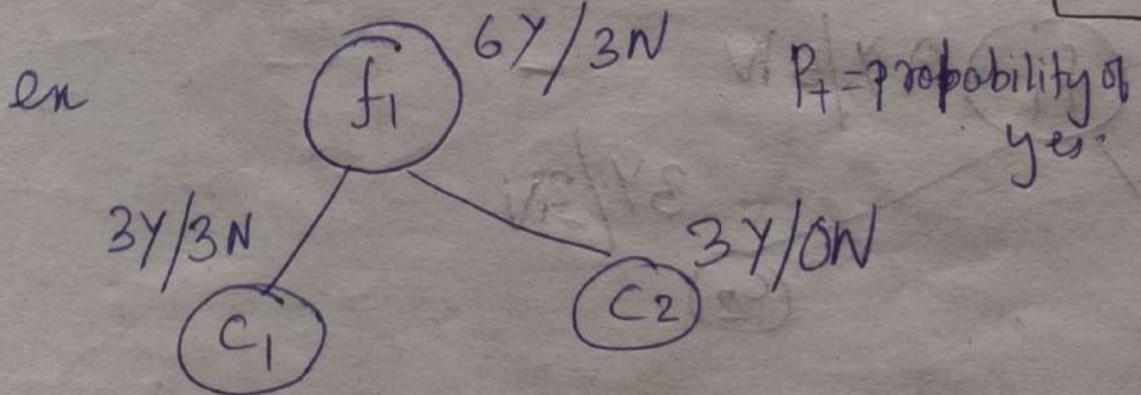
Entropy

$$H(S) = -(P_+) \log_2 P_+ - (P_-) \log_2 P_-$$

Decision tree classifier
by default uses Gini Impurity

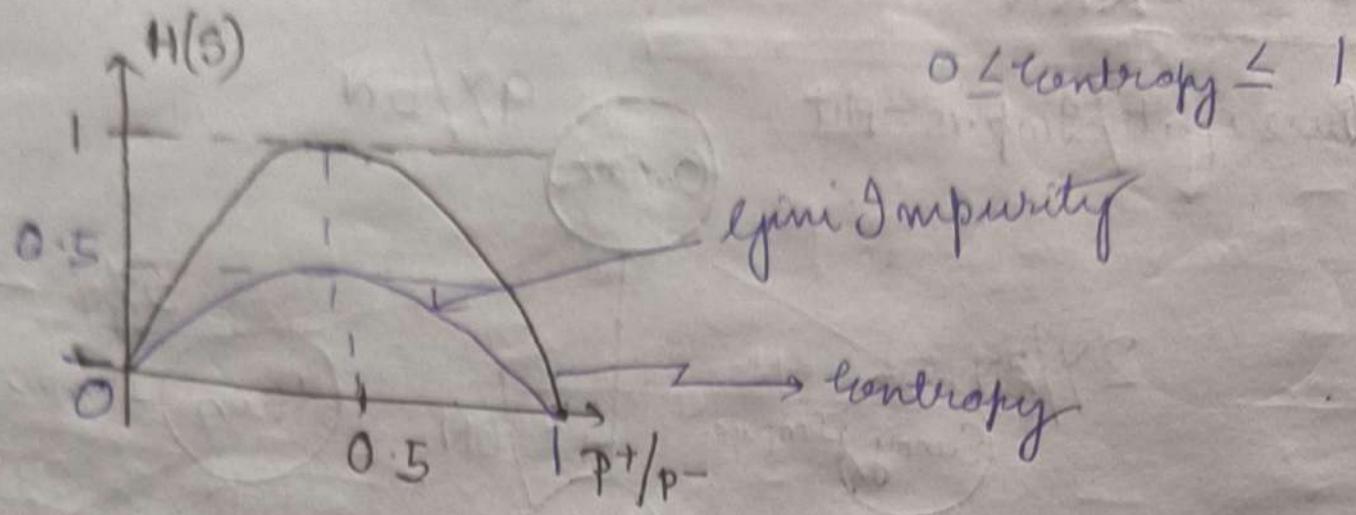
Gini Impurity

$$G.I = 1 - \sum_{i=1}^n (P_i)^2$$



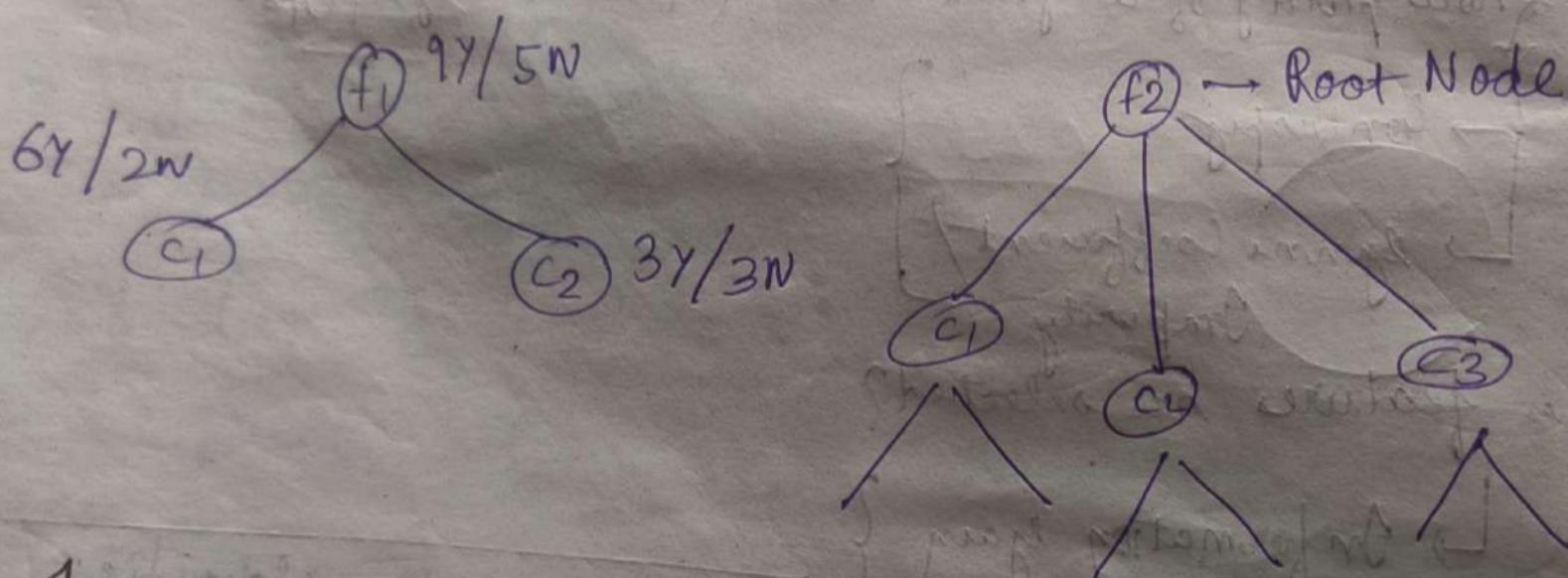
Entropy $H(S_{C2}) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$

$$H(S_{C1}) = 1 - -1 \log 1 = 0$$



Purity Test \rightarrow Entropy \rightarrow 1 (Impure)
0 (Pure)

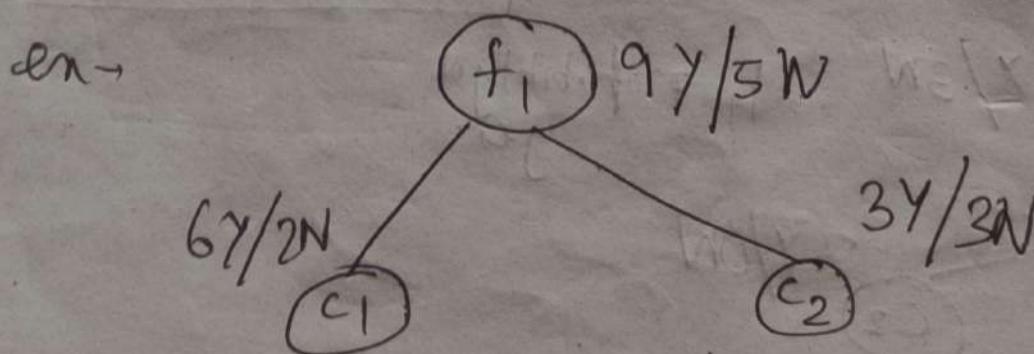
② Which feature to split?



Information gain

$$\text{Gain}(S, f_1) = H(S) - \sum_{\text{Value}} \frac{|S_V|}{|S|} H(S_V)$$

Entropy of Root Node



$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right) \approx 0.94$$

$$H(S_{VC_1}) = -\frac{6}{8} \log_2 6/8 - \frac{2}{8} \log_2 2/8 = 0.81$$

$$H(S_{VC_2}) = 1$$

$$g_{\text{gain}}(S, f_1) = 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$\begin{aligned} g_{\text{gain}}(S, f_1) &= 0.049 \\ g_{\text{gain}}(S, f_2)_{\text{if}} &= 0.051 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad \begin{array}{l} g_{\text{gain}}(S, f_2) >> g_{\text{gain}}(S, f_1) \\ \text{f}_2 \text{ should be selected} \end{array}$$

* Gini Impurity

$$GI = 1 - \sum_{i=1}^n p_i^2 \quad n = 2 \text{ output} \quad \left. \begin{array}{l} \text{yes} \\ \text{no} \end{array} \right\}$$

$$= 1 - ((p_+)^2 + (p_-)^2)$$

$$= 1 - \left(\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right)$$

$$= 1 - \left(\frac{1}{2}\right)$$

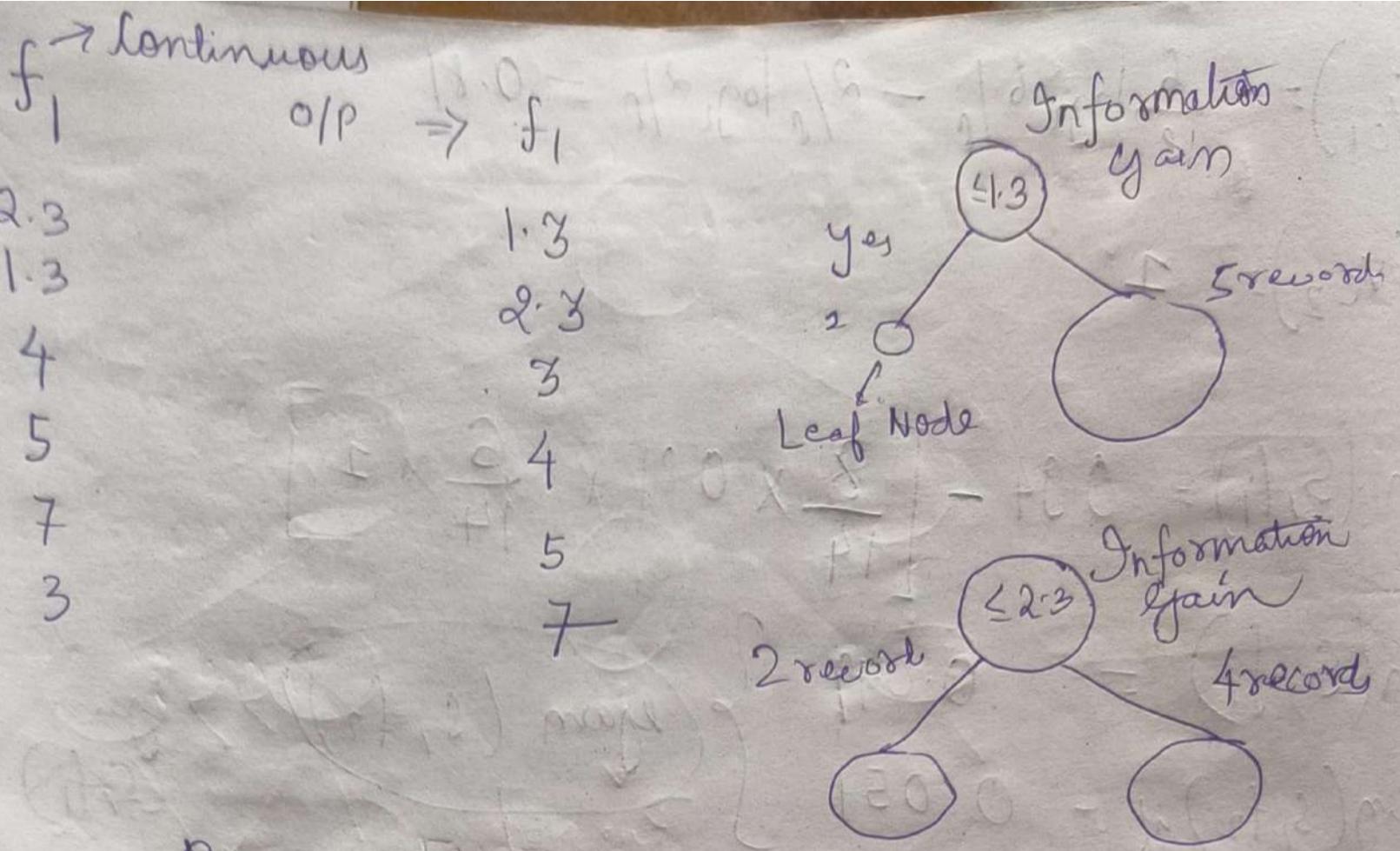
$$= \underline{0.5}$$

* Decision tree has one of the worst time complexities.

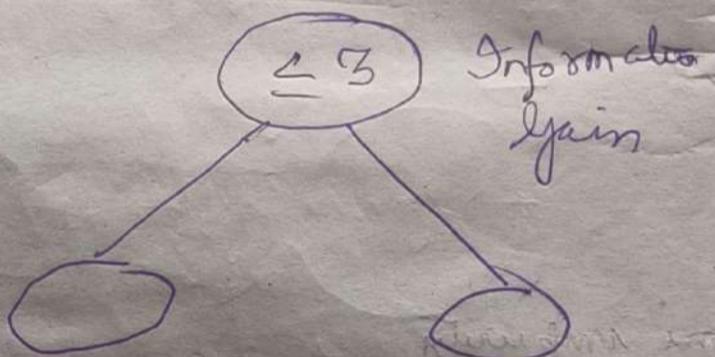
→ More execution time → Entropy

→ Gini Impurity is faster

huge no. of feature apply GI
less no. of feature - Impurity

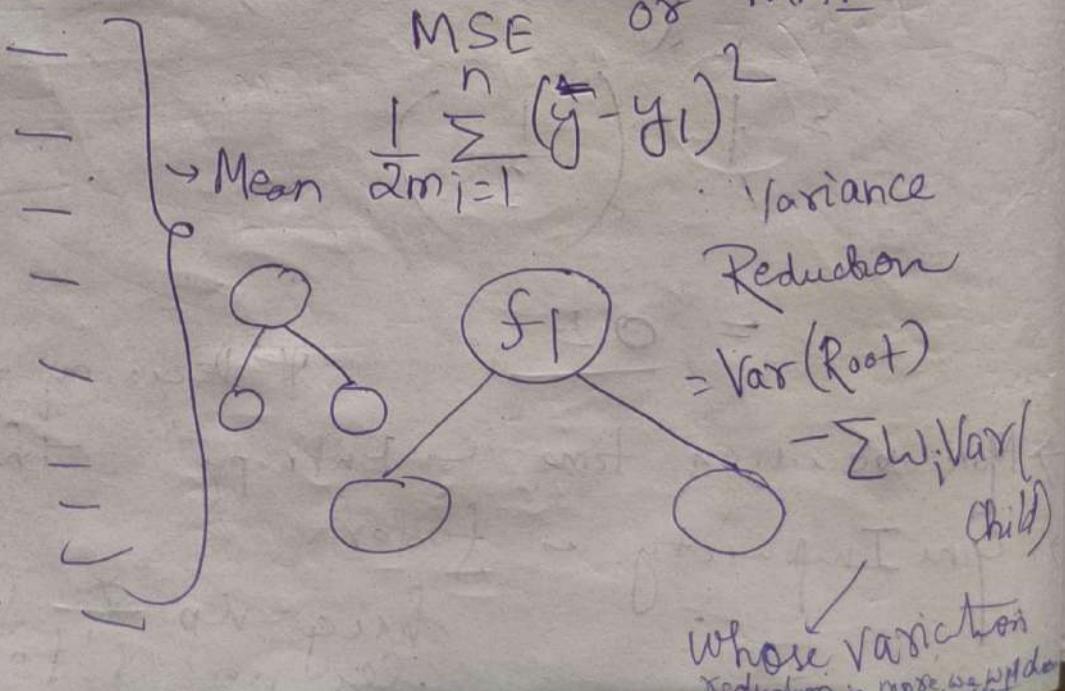


Best information gain will get selected and from the splitting will happen.



Decision Tree Regressor

f_1, f_2 O/P
Continuous

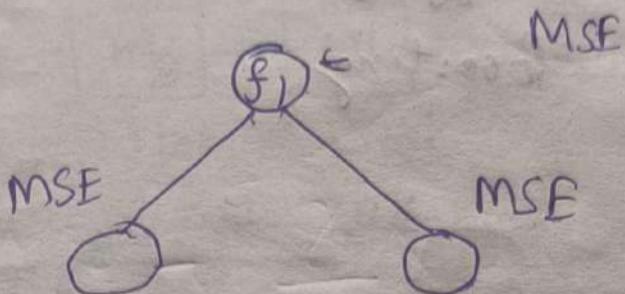


We take mean if two values present at leaf node.

Hyperparam

Decision tree regressor

O/P
f_1
20
24
—
26
—
28
—
30
:



Hyperparameter

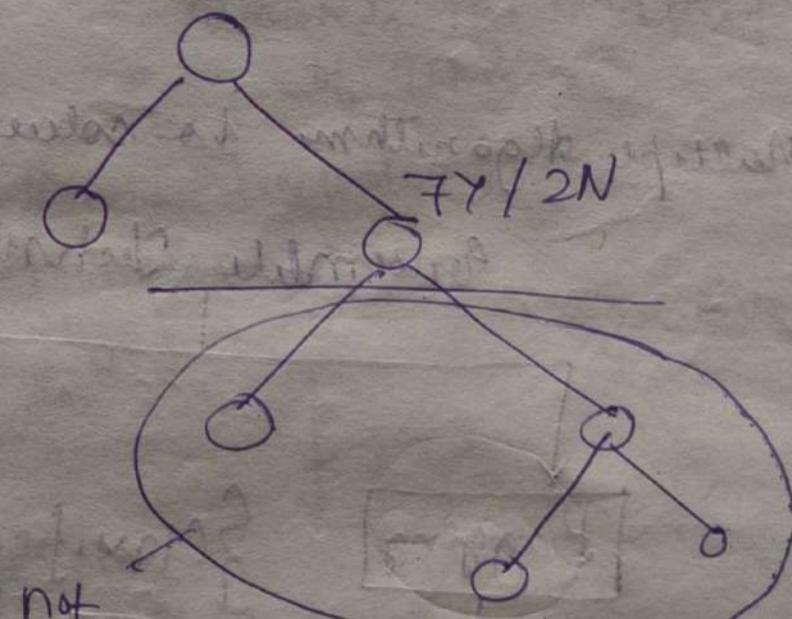
ENCODED - DECODING

Decision → Overfitting

{① Post Prunning }
② Pre Prunning

To avoid overfitting.

Post Prunning



Pre prunning hyperparameters

max_depth, max_leaf

Grid Search CV

Agenda

1) Ensemble Techniques → Bagging

Boosting

2) Random Forest

3) AdaBoost

4) Xg boost

ENSEMBLE TECHNIQUES

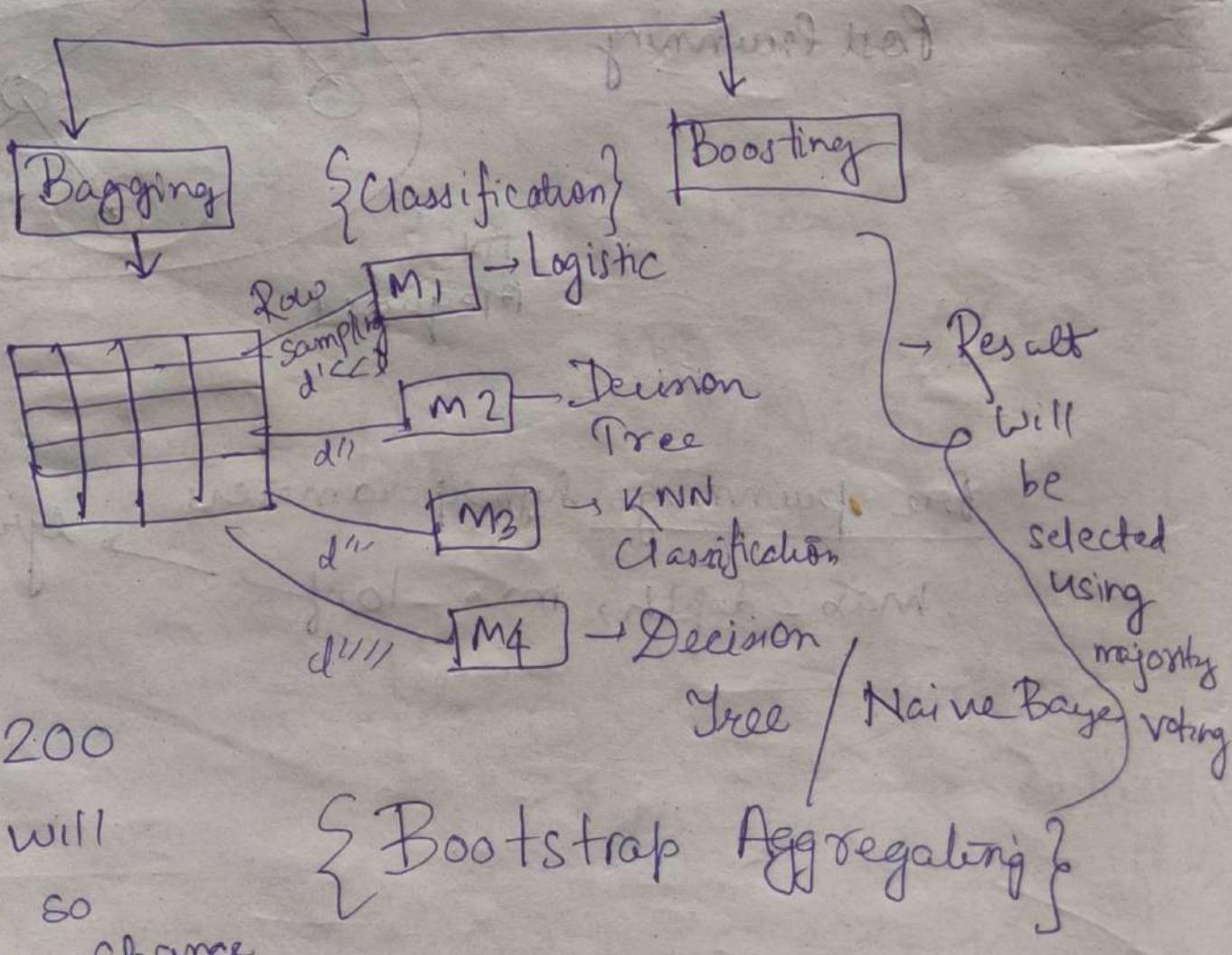
① Classification and Regression

1 Algorithm → Classification
Regression

Multiple algorithms to solve a problem.

Ensemble Techniques

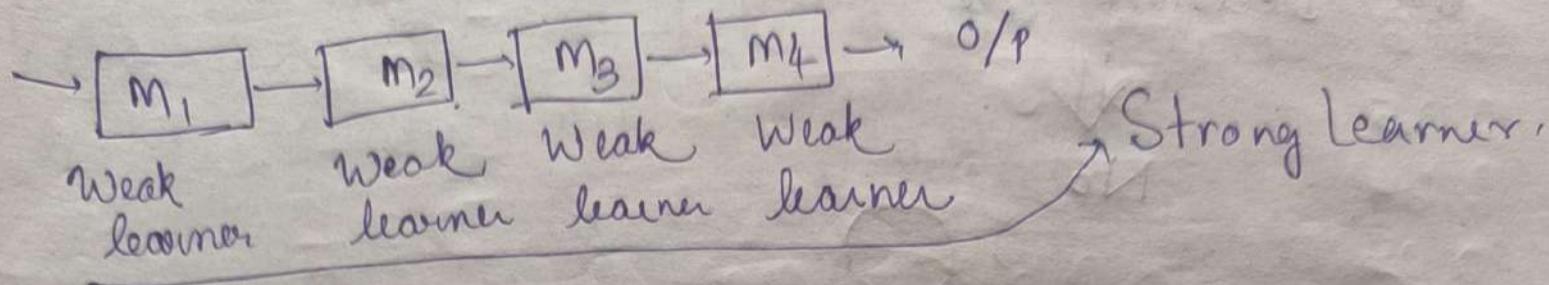
In case
of regression
avg of
all the D
outputs
will
be taken.



100 to 200

models will
be used so
very less chance
of tie.

Boosting (Models will be present sequentially)



Bagging

① Random Forest classifier

② Random Forest Regressor

Boosting

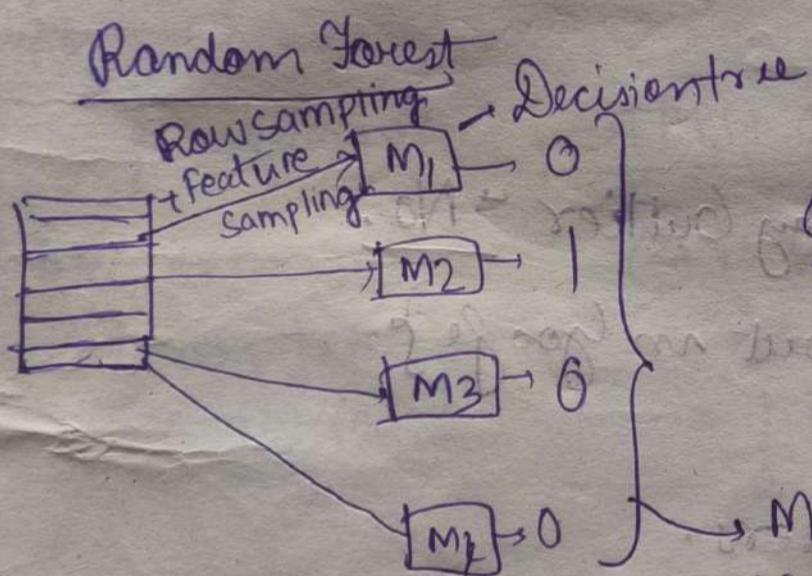
① Adaboost

② Gradient Boost

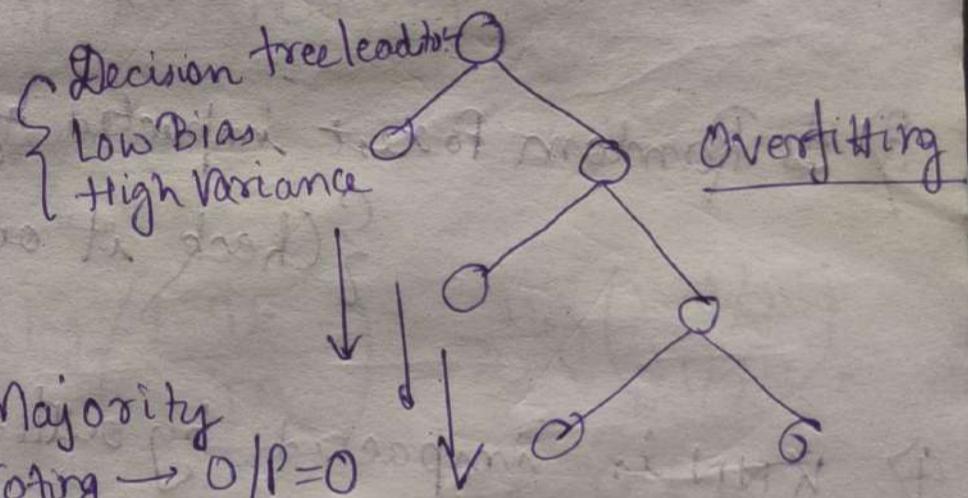
⑤ XG Boost

(Extreme Gradient Boost)

⇒ Random Forest classifier and Regressor



what is the main problem of DT?



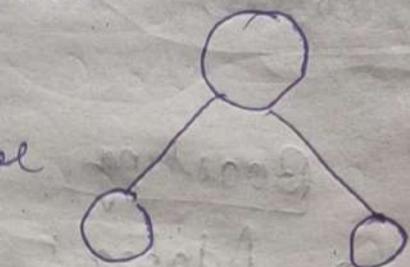
- * In random forest, only decision tree is used.
- * Because of the overfitting condition of the decision tree, we are combining many decision trees so that we can get a generalised model which has low bias and low variance.

- * In random forest regressor, the mean is taken out of the values given by various decision trees.

3) ~~Q~~ Is normalization required in Random Forest?
or Decision Tree

↓
No

Decision tree
does split, so



even if we minimize the data also
standardization required in KNN? → of that
the splits will not be
much importance.
Yes, because it has Euclidean and
Manhattan distance, so that
the computational distance becomes easy.

3) Is Random Forest impacted by outliers → No.
{Check it out in Google}

4) KNN is impacted by outliers.

Bagging → Random Forest is used
mostly

or

Custom Bagging Techniques

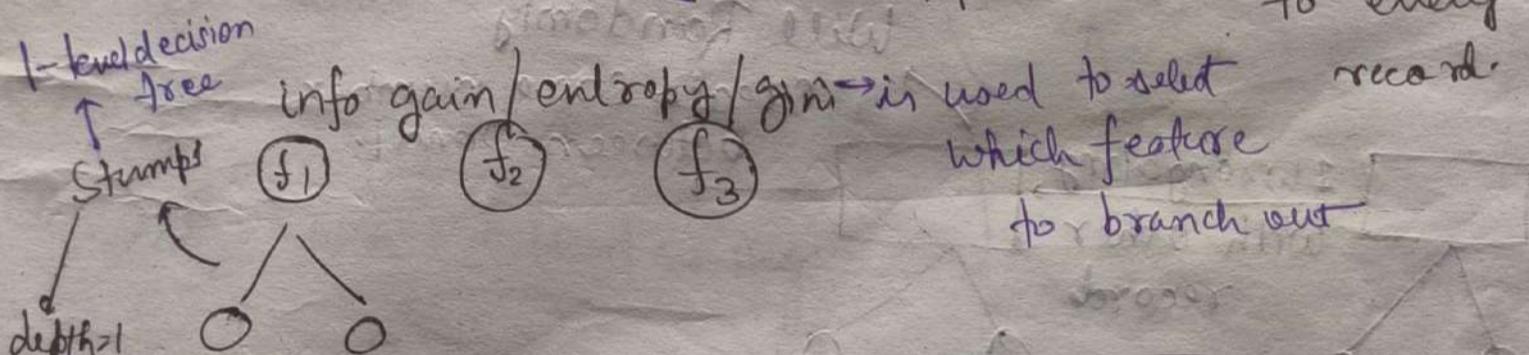
can be used:-

Custom

D → Can
D → be
D → done
D → manually
D → with
the help of hands

BOOSTING

Decision Tree		Adaboost		Decision Tree specifically		Weight		Weight	
f_1	f_2	f_3	f_4	O/P				$\Sigma \text{Weight} = 1$	
-	-	-	-	Yes		$1/7$	Update	0.05	
-	-	-	-	No		$1/7$	$\rightarrow 0.05$	$\frac{1}{7} + \frac{1}{7} + \frac{1}{7} + \frac{1}{7} + \frac{1}{7} + \frac{1}{7} + \frac{1}{7}$	
-	-	-	-	-		$1/7$	$\rightarrow 0.05$	$+ \frac{1}{7} > 1$	
-	-	-	-	-		$1/7$	$\rightarrow 0.349$		
-	-	-	-	-		$1/7$	$\rightarrow 0.05$	* Initially	
-	-	-	-	-		$1/7$	$\rightarrow 0.05$	equal weight	
-	-	-	-	-		$1/7$	$\rightarrow 0.05$	is given =	
-	-	-	-	-		$1/7$	$\rightarrow 0.05$	to every	
1 level decision		stack of tree							



{Weak learners}

STEPS OF ADABOOST

Step ① Total error from
1st decision tree

$$= 1/7 \text{ (assume)}$$

S- ② Performance of Sun

$$= \frac{1}{2} \log_e \left(\frac{1 - TE}{TE} \right) = \frac{1}{2} \log \left(\frac{1 - 1/7}{1/7} \right) \\ = 0.895$$

S-3 Updation of Weight

New Sample Weight =

Correct Record

~~orrect below~~ Weight * e⁻ (Perf. of stump)

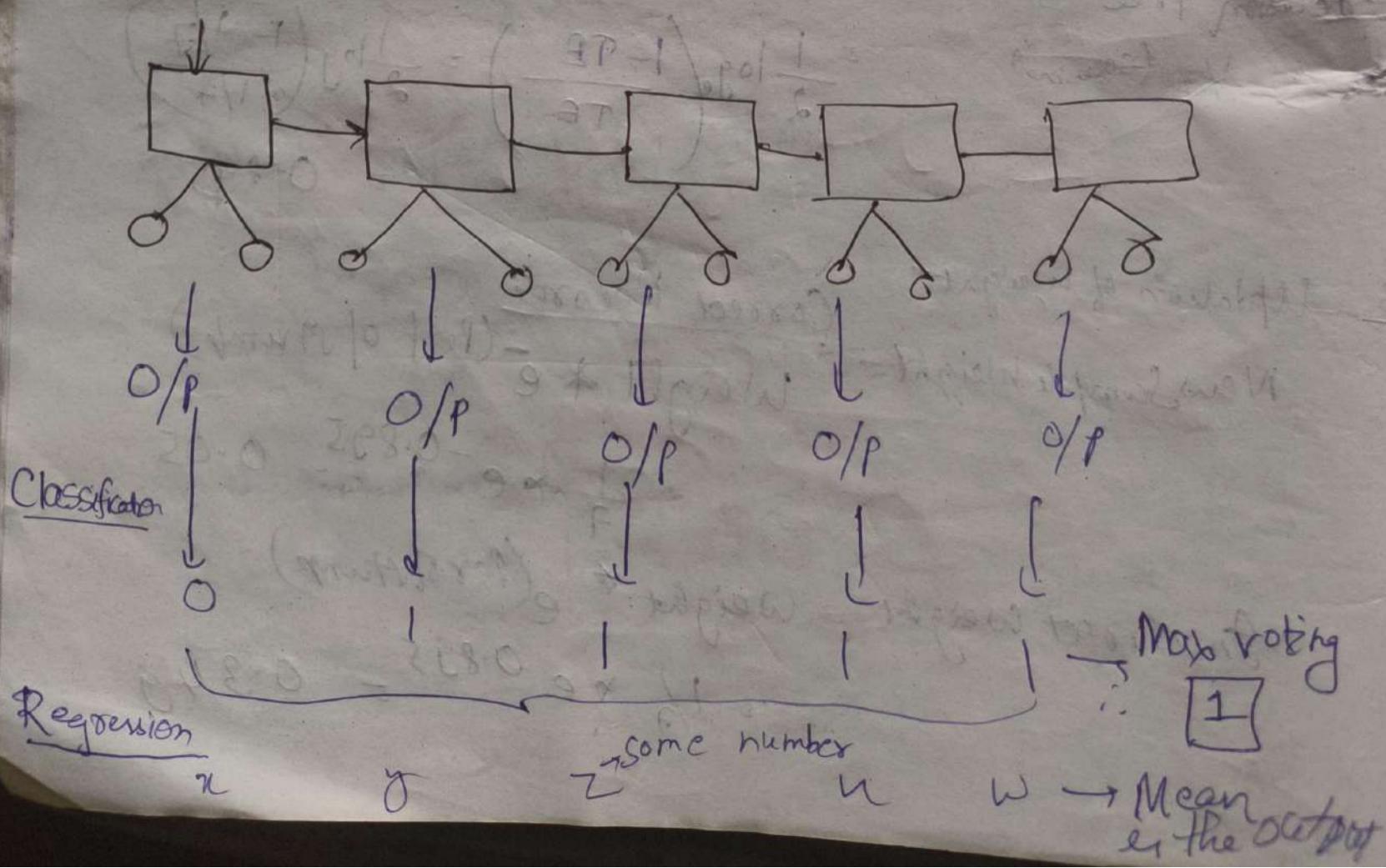
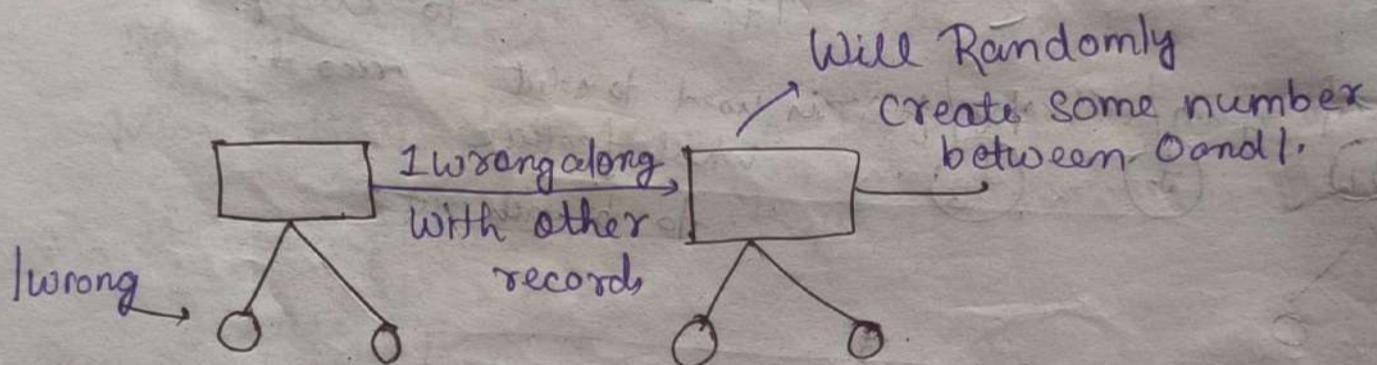
$$\frac{1}{7} \times e^{-0.895} = 0.05$$

$$\text{Incorrect weight} = \text{Weight} \times e^{(\text{Per cent})}$$

$$1/7 \times e^{0.895} = 0.349$$

→ New Weight

	Normalized Weight	Buckets
$0.05 =$	$\frac{0.05}{0.649} \rightarrow 0.07$	$[0 - 0.07]$
$0.05 =$	$\frac{0.05}{0.649} \rightarrow 0.07$	$[0.07 - 0.14]$
0.05	$\frac{0.05}{0.649} \rightarrow 0.07$	$[0.14 - 0.21]$
0.349	$\frac{0.349}{0.649} \rightarrow 0.537$	$[0.21 - 0.747]$
0.05	$\frac{0.05}{0.649} \rightarrow 0.07$	$[0.747 - 0.85]$
0.05	$\frac{0.05}{0.649} \rightarrow 0.07$	$[-]$
0.05	$\frac{0.05}{0.649} \rightarrow 0.07$	$[- 2]$
has the $buckets$ + 0.05		
$\Sigma \neq 1$		
	$= 0.649$	



Black box model vs White Box Model

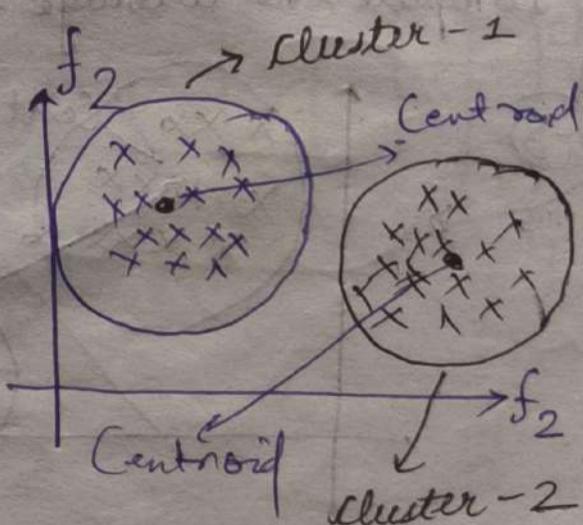
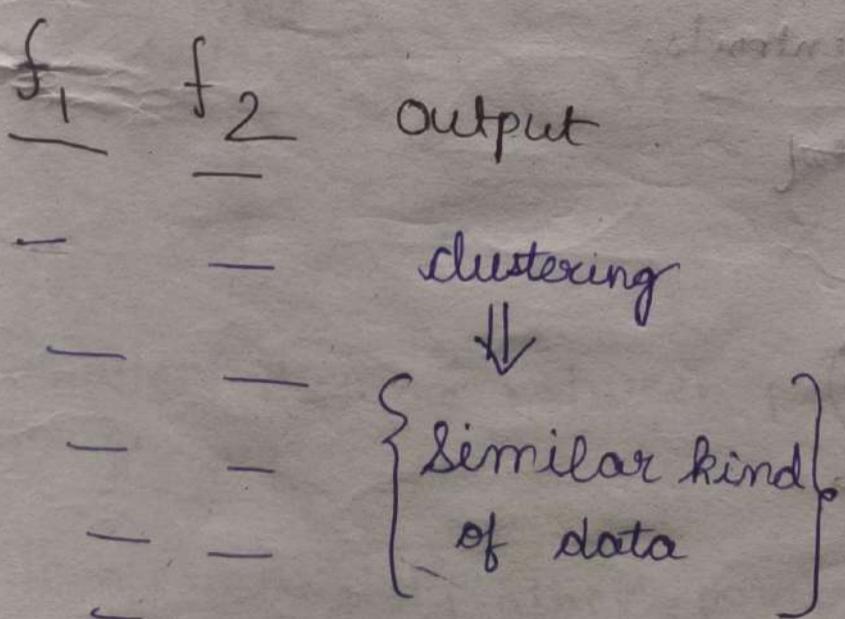
- Linear Regression → White box
- Random Forest → Black box
- Decision Tree → White box
- ANN → Black box

UNSUPERVISED ML MACHINE LEARNING

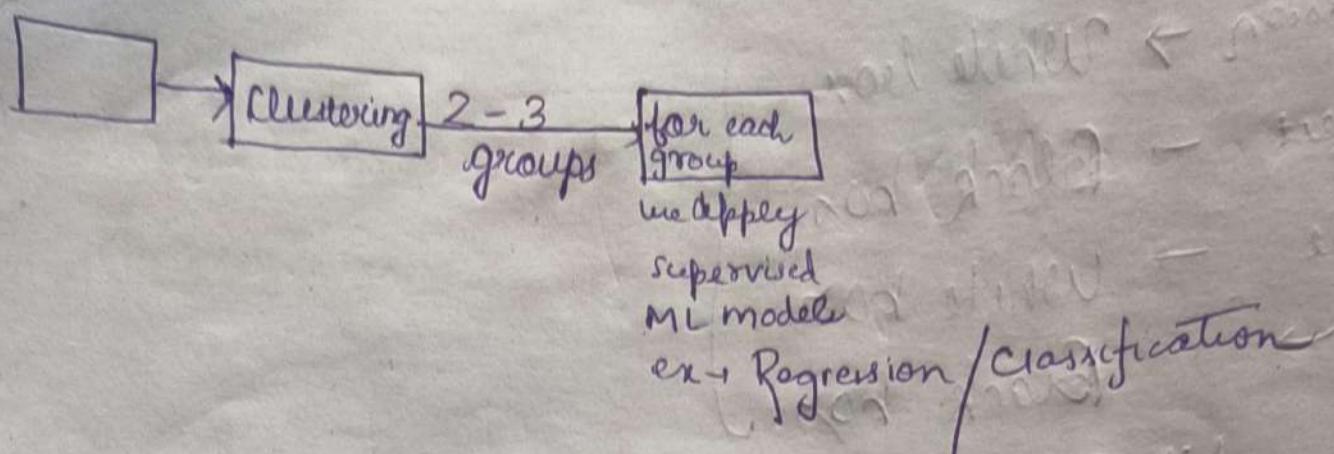
- ① K-Means Clustering
- ② Hierarchical Clustering
- ③ Silhouette Score
- ④ DB-Scan Clustering

Agenda

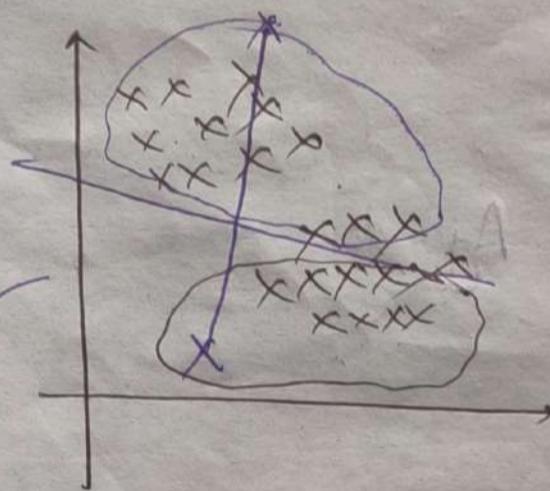
K Means Clustering



Custom Ensemble Techniques



K - Means → $k = \text{Centroids}$

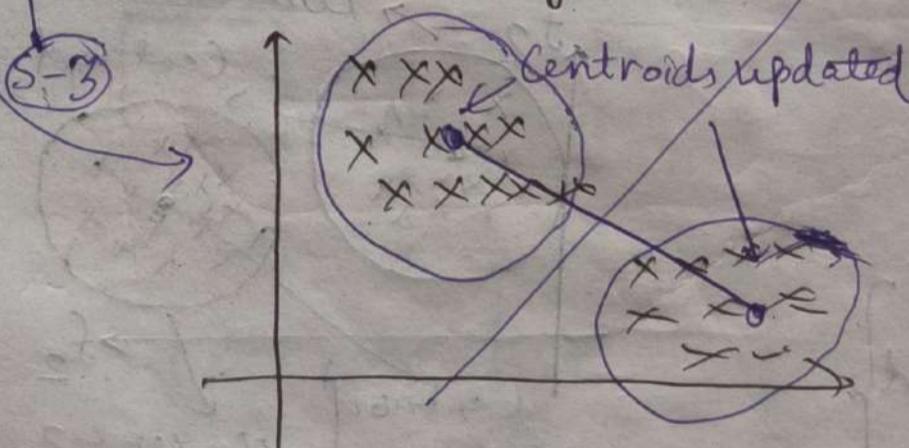


$k = 2$

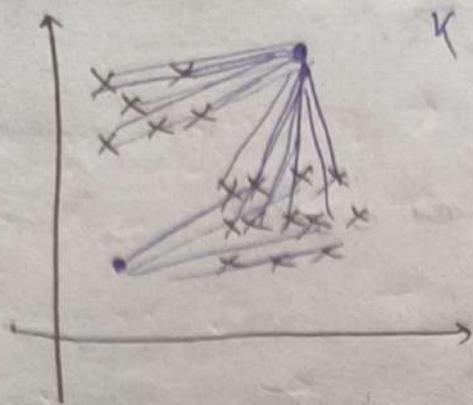
Euclidean Distance

Steps :-

- ① We try k -values \Rightarrow suitable $k = 2$
- ② Initialize k -number of centroid
- ③ Compute the average to update centroids.

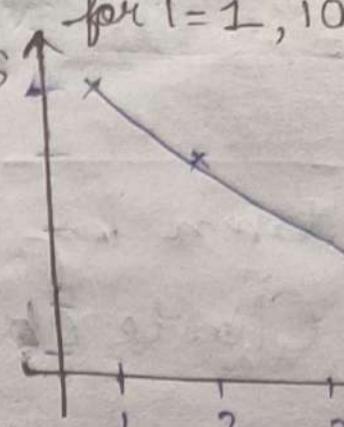


Elbow method (helps to initialize k -value
or optimise k -value
or what k -value
to choose)



for $i=1, 10$

WCSS



elbow curve

this K -value is selected.

WCSS \rightarrow within cluster sum of squares

for finding K -value : we use elbow curve

abrupt change from here and we take this K -Value

for validating our K NN model : we use Silhouette score



for $K=4 \Rightarrow 4$ groups

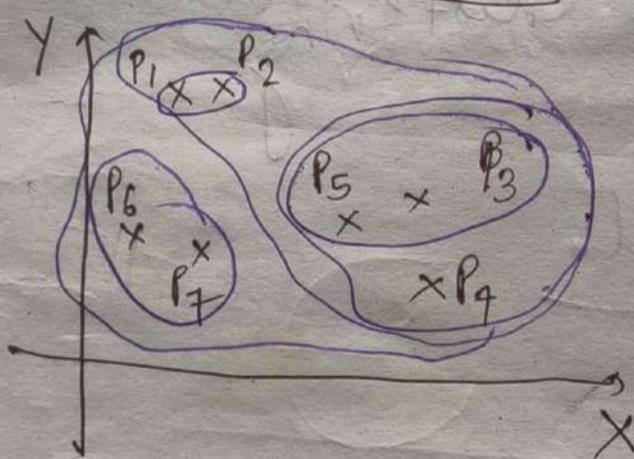
$K=n \Rightarrow n$ groups

$\{K \text{ means} + 1\}$

should be far away from the cluster.
you need to find the longest vertical line that has no horizontal

Dendrogram

- ① * The centroid initially taken
- ② Hierarchical clustering



S-1

* Smallest distance between pairs are grouped.

S-2

* Then, consequently other points are grouped

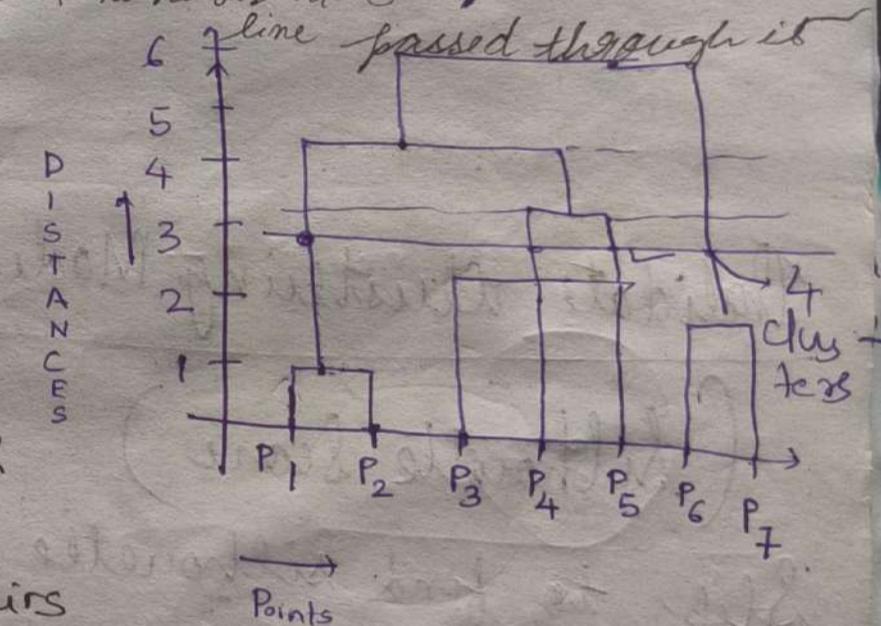
S-3

* All the points have been grouped into 1-group.

S-4

* Find the longest vertical line, which does not have any horizontal line passing through it

S-5 Draw a horizontal line passing through it and find the



Number of intersecting points passing through the horizontal line
The number of intersections equals the 'number of groups'.

* Q → Mar-time is taken by K-Means or
Hierarchical Clustering?

↳ Hierarchical

Clustering

(due to Dendrogram)

* Dataset is small → Hierarchical
Clustering

* Dataset is large → K-Means
Clustering

Validate Clustering Model:

(Silhouette Score)

Steps to find silhouette score →

① a_i find out

average of



all the distances from the centroid in one of the clusters.

② b_i is the average of all the distances from
all the points of one cluster to all
the points of another cluster.

③ for a good model
 $b(i) \gg a(i)$

for a worse model
 $a(i) \gg b(i)$

④ Silhouette score varies from -1 to 1.

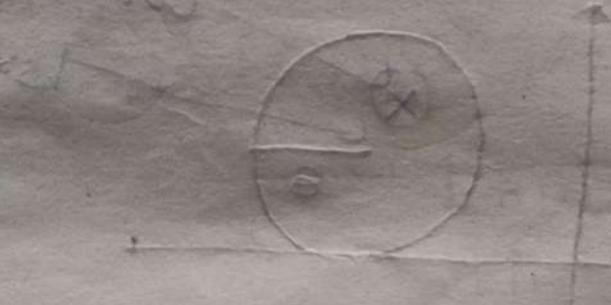
-1 being the worst
to 1 being the best

⑤ Formula of silhouette score (S)

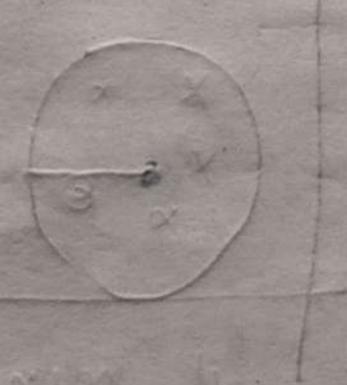
$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

, if $|C_i| > 1$
 \uparrow
 $i^{\text{th}} - \text{cluster}$

$$S(i) = 0, \text{ if } |C_i| = 1$$



+ = if non - s.t.

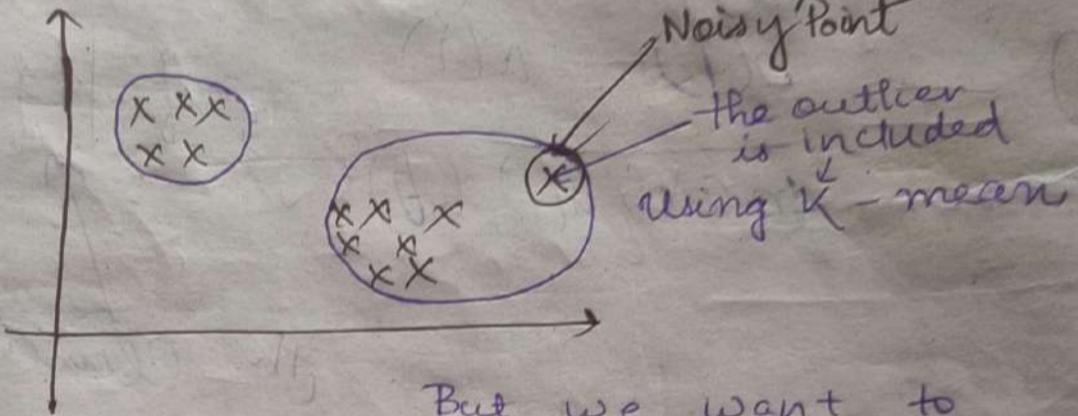


DB-Scan

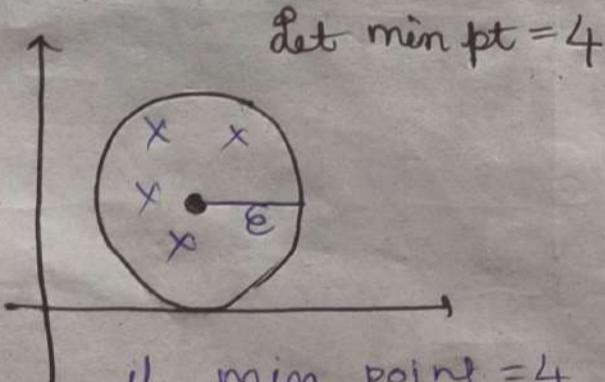
Clustering

(Density Based Spatial Clustering
of Applications with Noise)
(DBSCAN)

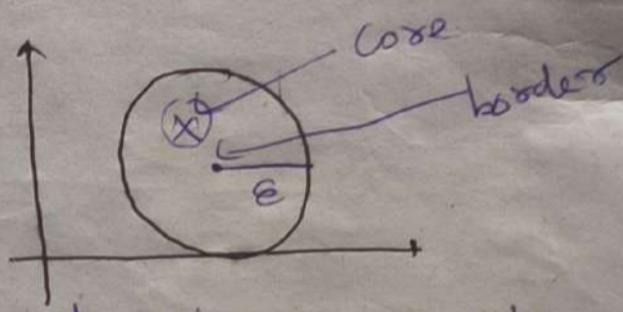
- ① Min pt (one hyperparameter) (e.g. $\text{min} = 3, 5 \text{ etc}$)
- ② Core points
- ③ Border points
- ④ Noise points
- ⑤ Epsilon (Radius of Circle)



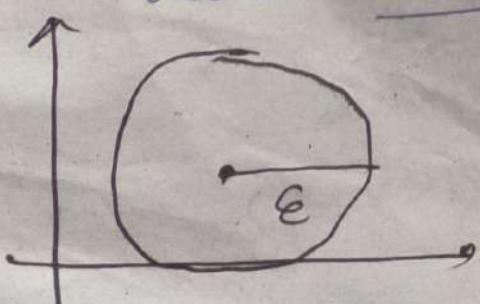
But we want to leave out the outliers, in some cases therefore we use DB SCAN.



If $\text{min point} = 4$, and no. of points inside circle = 4, then the centre becomes core point.



If $\text{min point} = 4$, and no. of point ~~is~~ is greater than 0 but less than 4, then the center becomes border point.



If no point inside the circle, the center becomes Noisy point.

* We neglect all the noise points, and combine
and core point.

DB Scan is much better
than K-Means

* Practical implementation

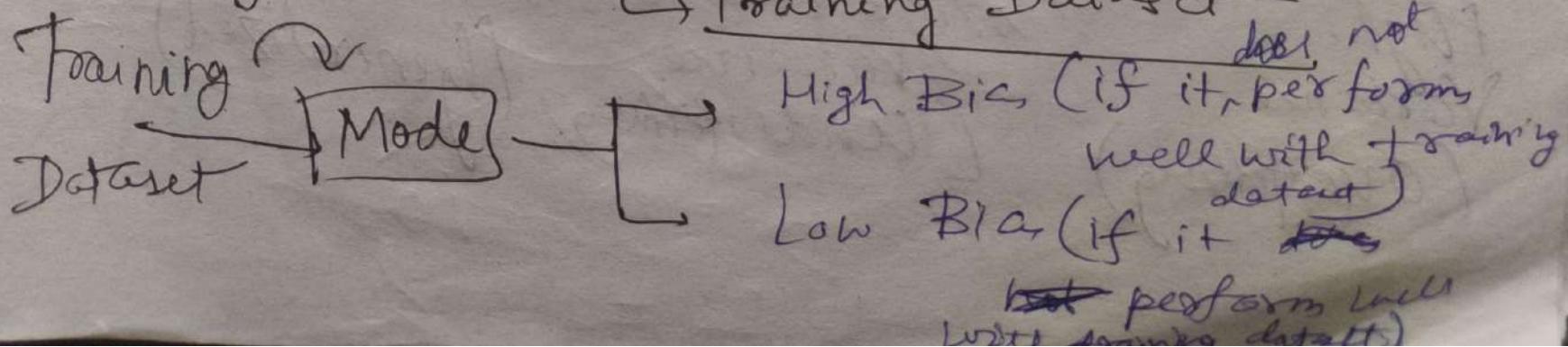
→ we use Silhouette score for validation & clustering
if $\text{sum_cluster} \rightarrow \text{high}$ we get
choice, ex for $n=2$ or $n=4$,
choose the larger model as it
helps in generalisation

Defⁿ of BIAS AND VARIANCE

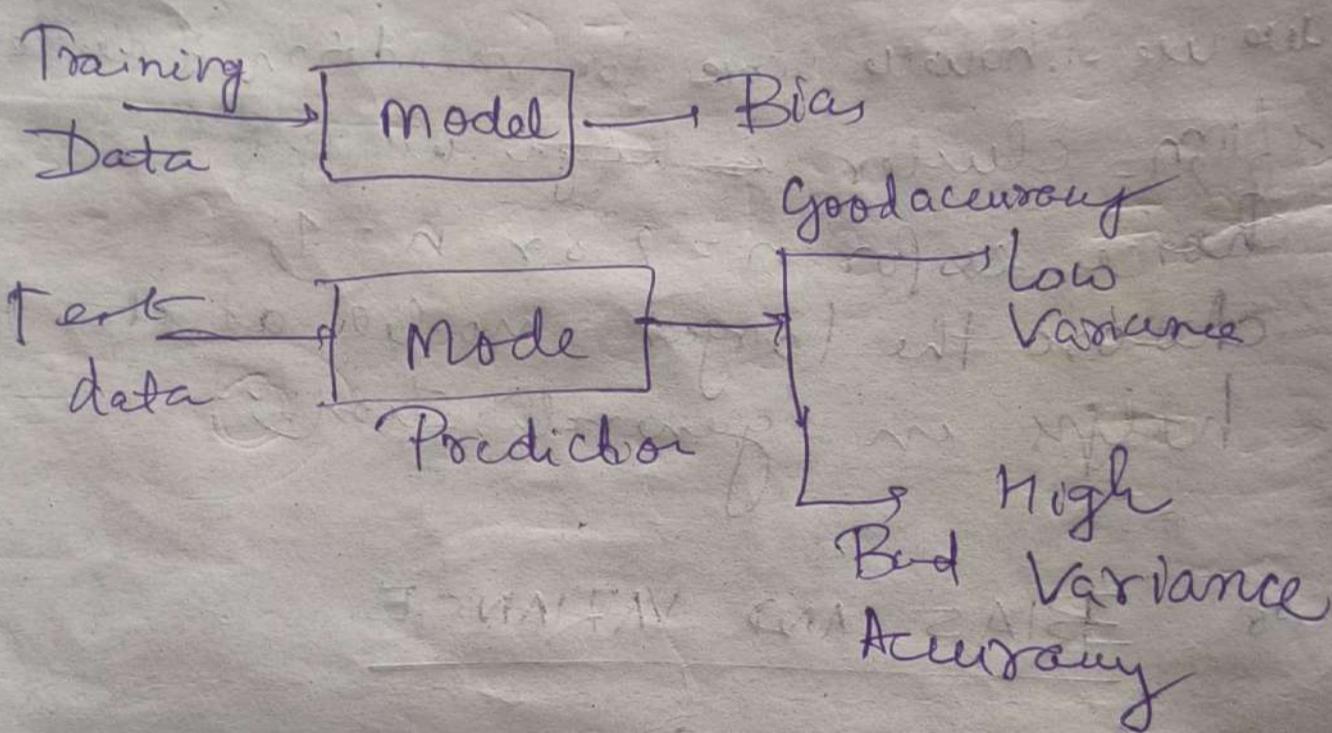
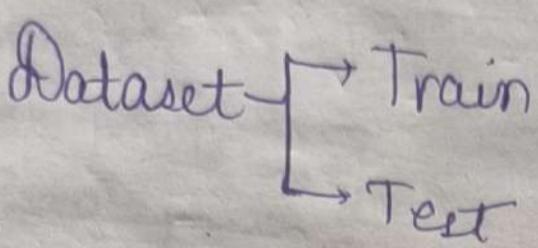
Training Dataset = 80%
Test Dataset = 20% } \Rightarrow Overfitting

Low Bias ✓
High Variance ✓

① Bias - It is a phenomenon that skews the
result of an algorithm in favour of
against an idea.



Variance: Variance refers to the changes in the model when using different portions of training or test data.



Model 1

Train: 90%
Acc: 90%
Test: 75%
Acc: 75%

Low Bias

High Variance
(Overfitting)

Model 2

60%

55%

High Bias

High Variance

(Underfitting) (Generalized)

Model 3

90%

92%

Low Bias

Low Variance

DAY 7: XGBOOST CLASSIFIER AND REGRESSOR

Agenda:

- ① Xgboost Classifier
- ② Xgboost Regressor
- ③ SVM
- ④ SVR

① XGBOOST CLASSIFIER

Extreme Gradient Boosting

{Dataset}

(internally uses ~~dot~~ decision tree)

Salary	Credit	Approval	Residual	
≤ 50	B	0	-0.5	Dummy model
< 50	G	1	0.5	→ Base model (Weak learner) $P_r = 0.5$ in case of classif cut
≤ 50	G	1	0.5	$P_r = 0.5$
> 50	B	0	-0.05	Default
> 80	G	1	0.05	
> 50	N	1	0.5	
$\leq 50k$	N	0	-0.5	

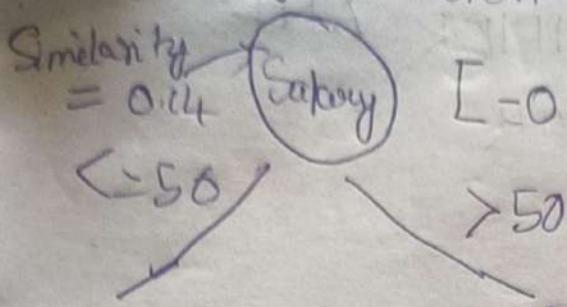
Step① Create a Binary decision tree using features

Step② Calculate Similarity Weight

$$= \frac{(\sum \text{Residual})^2}{\sum P_r(1-P_r) + \gamma}$$

Step③ Calculate Information gain

Use Binary decision tree in Xgboost always



$$[-0.5, 0.5, 0.5, -0.5] \quad [-0.05, 0.5, 0.5]$$

↓ Considering

= Similarity = 0

$$= (-0.5)^2 + 0$$

$$= \frac{(-0.5 + 0.5 + 0.5 - 0.5)}{2}$$

$$0.5(1-0.5) + 0.5(1-0.5)$$

$$+ 0.5(1-0.5) + 0.5(1-0.5)$$

$$= \frac{(-0.5 + 0.5 + 0.5)}{0.5(1 - 0.5)}$$

$$+0.5(1-0.5)$$

$$+0.5(1-0.5)$$

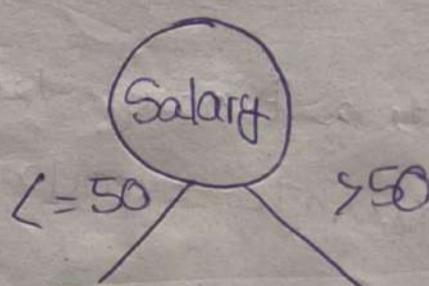
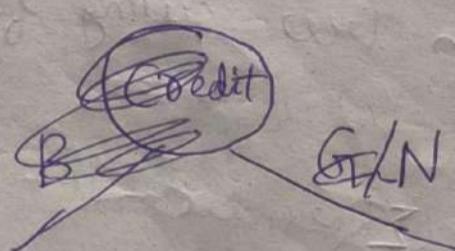
$$= \frac{0.25}{0.75}$$

0.33

③ Information Gain

$$= 0 + 0.33 - 0.14$$

$$\approx 0.19$$



$$\frac{0.25}{0.5(1-0.5)} =$$

[-0.5]

$$\text{Info gain} = 1 + 0.3$$

$$\frac{0.2}{0.2}$$

Sim weight = 1

$$= 0 \\ = 1.33$$

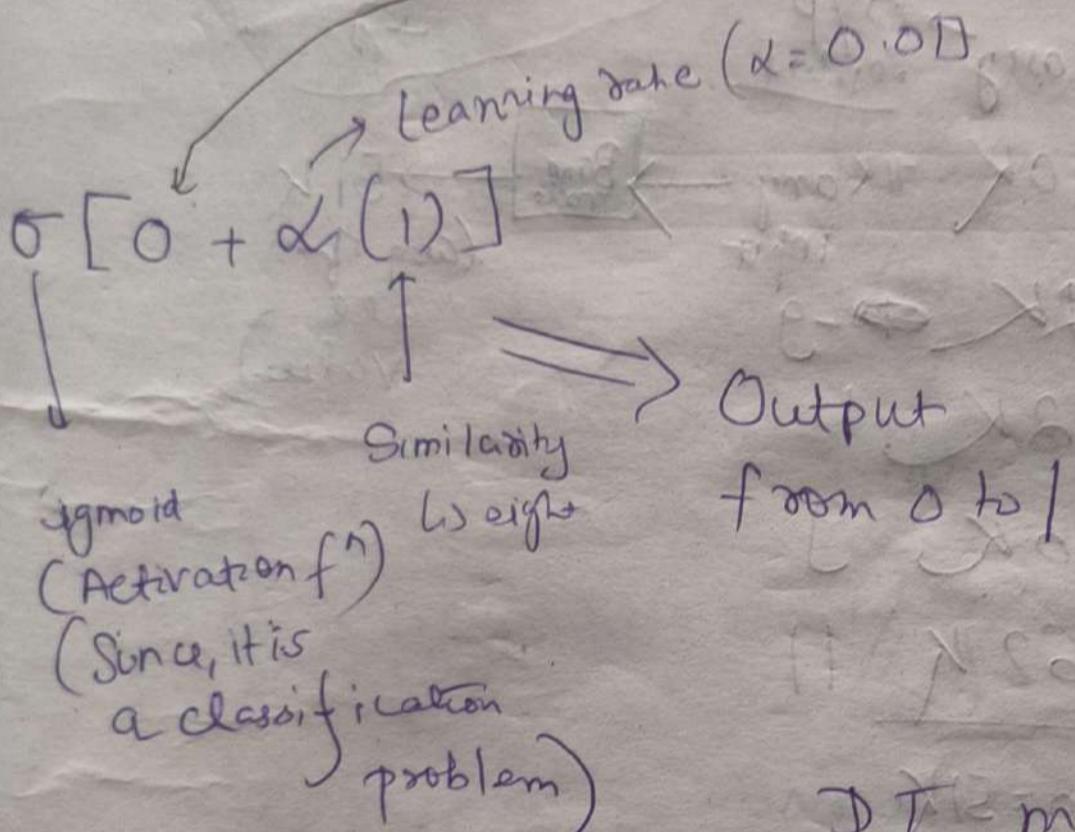
$$\begin{aligned} \text{Sim weight} \\ = 0.23 \end{aligned}$$

for base model only, ex -> for first record as in dataset

$$\text{Real probability} = \log\left(\frac{P}{1-p}\right)$$

$$= \log\left(\frac{0.5}{0.5}\right) = 0$$

Binary decision tree is a specialist in decision



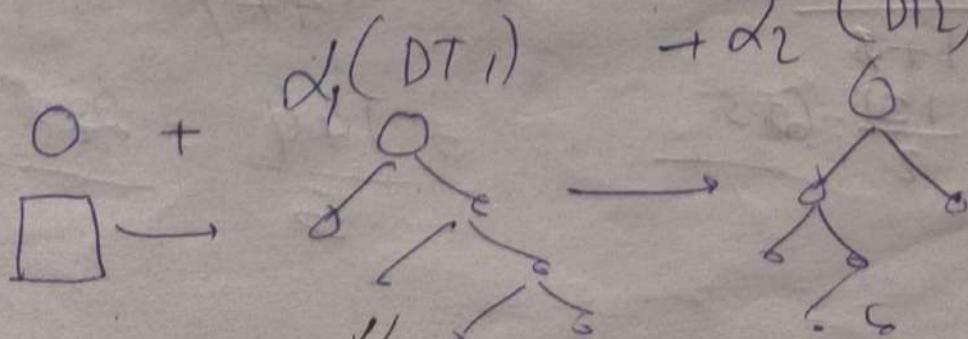
$$\sigma[0 + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) +$$

$$\dots + \alpha_4(DT_4) + \dots + \alpha_n(DT_n)]$$

New inference output
Record output

- * Xg boost is a black box model
- * Pre Pruning done (hyperparameter tuning)
- * Combining many decision tree

$$+ \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n)$$



based on independent features

* λ is a hyperparameter

* λ is set up by cross-validation

X-G Boost Regression

Similarity
Weight = $(\sum \text{Residual})^2$
No. of Residual + λ

Exp	gap	Salary	O/P
2	yes	40K	71K any
2.5	yp	42K	input
3	NO	52K	1100
4	NO	60K	9
4.5	yes	62K	11

$$\text{Avg} = 51K$$

$$\boxed{\text{Exp}} \rightarrow S_{\text{WS}} = \frac{(-11+9+1-9+11)^2}{6} = \frac{1}{6}$$

$$S^2 \quad \lambda = g(\text{let}) \quad [-1] \quad [-9, 1, 9, 11]$$

Since Similarity Weight

$$= \frac{121}{1+1}$$

$$= \frac{121}{2} = 60.5$$

$$\frac{(-9+1+9+11)^2}{4+1}$$

$$\frac{144}{5} = 28.8$$

$$\text{Information Gain} = 60 \div 5 + 28.8 - \frac{1}{16} \approx 89.13$$

Exp

$$SW = \frac{1}{6}$$

$$\begin{aligned} &\leq 2.5 \\ &> 2.5 \end{aligned}$$

$$\begin{aligned} SW = \frac{(-11-9)^2}{2+1} &= \frac{400}{3} = 133.33 \\ SW = \frac{(1+9+11)^2}{4} &= \frac{(21)^2}{4} = 110.25 \end{aligned}$$

$$IG = \overline{\left[133.3 + 110.25 - \frac{1}{16} \right]}$$

better than previous split, we will use this split

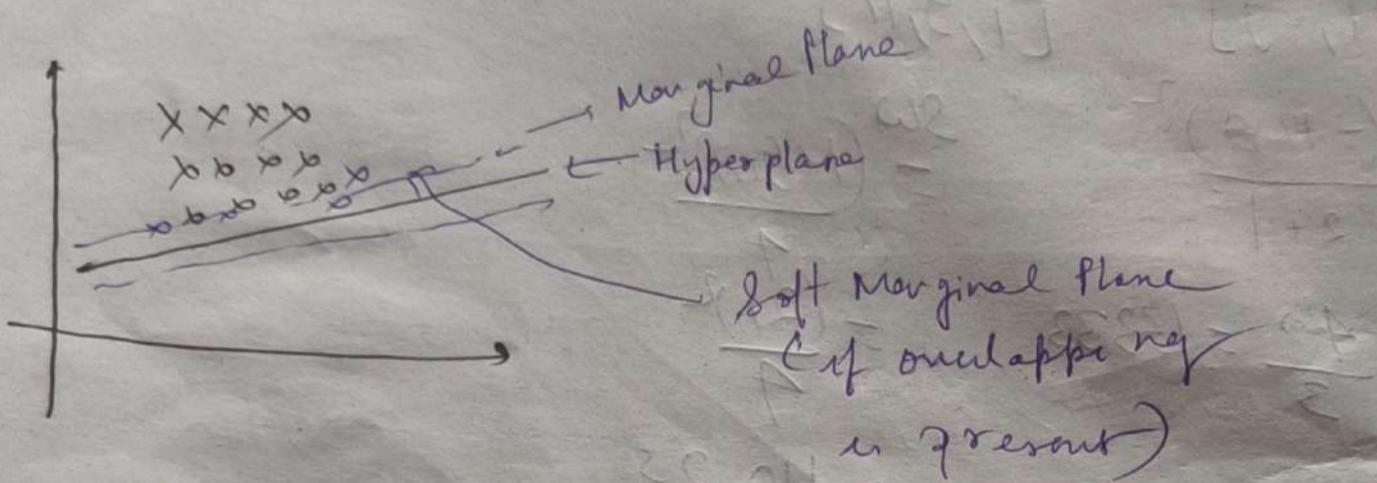
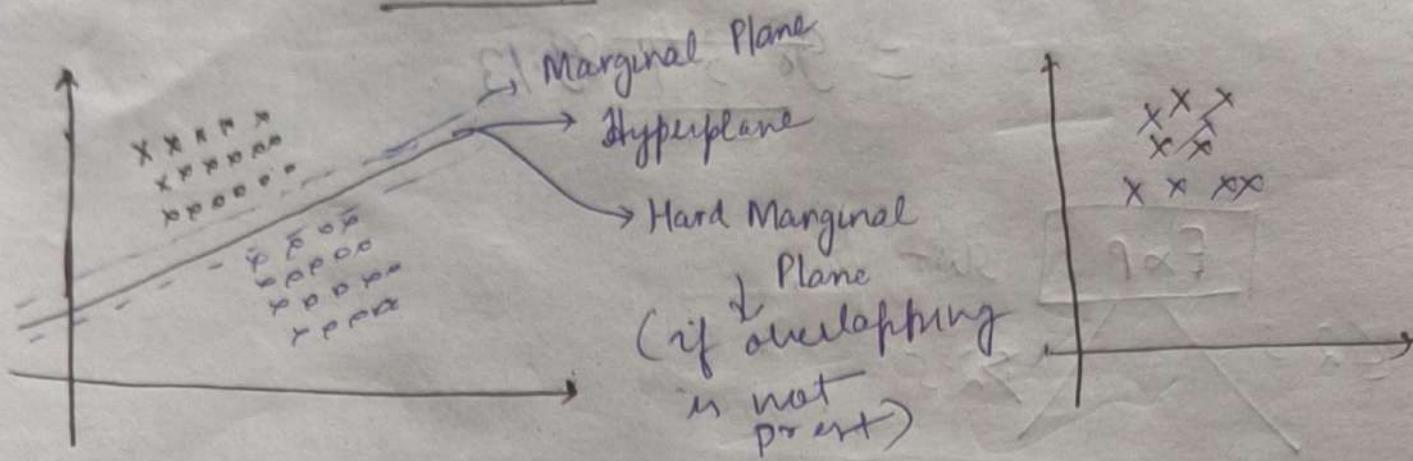
$$\underline{Op} = 51 + \alpha_1 \left(\frac{-11-9}{2} \right) + \alpha_2 (DT_2) + \alpha_3 (DT_3)$$

or

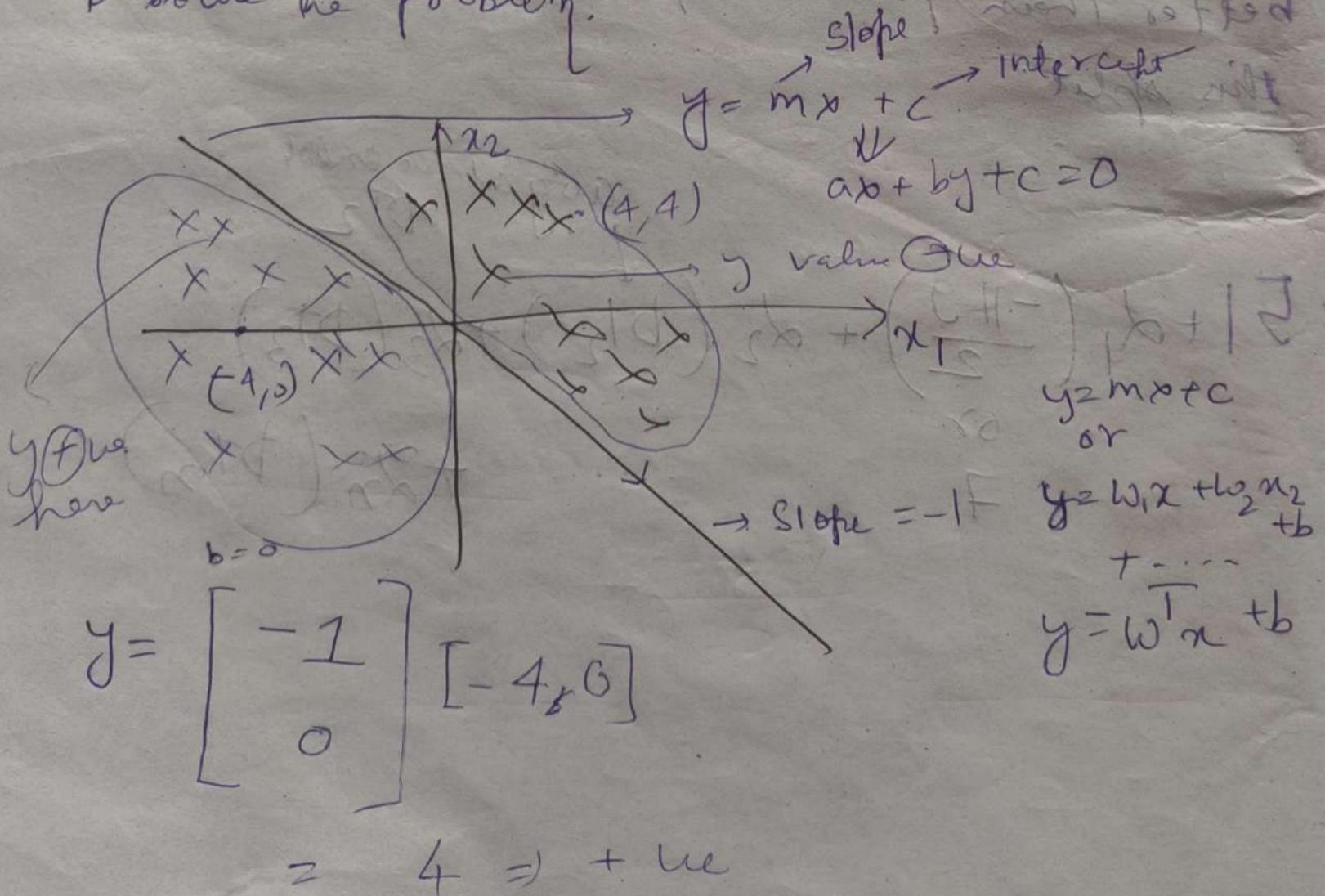
$$\underline{Op} = 51 + \dots + \alpha_n (DT_n)$$

avg at the end of DT₂

SVM



- * We focus on creating marginal plane with larger distance, even though there are errors we use hyperparameters to solve the problem.

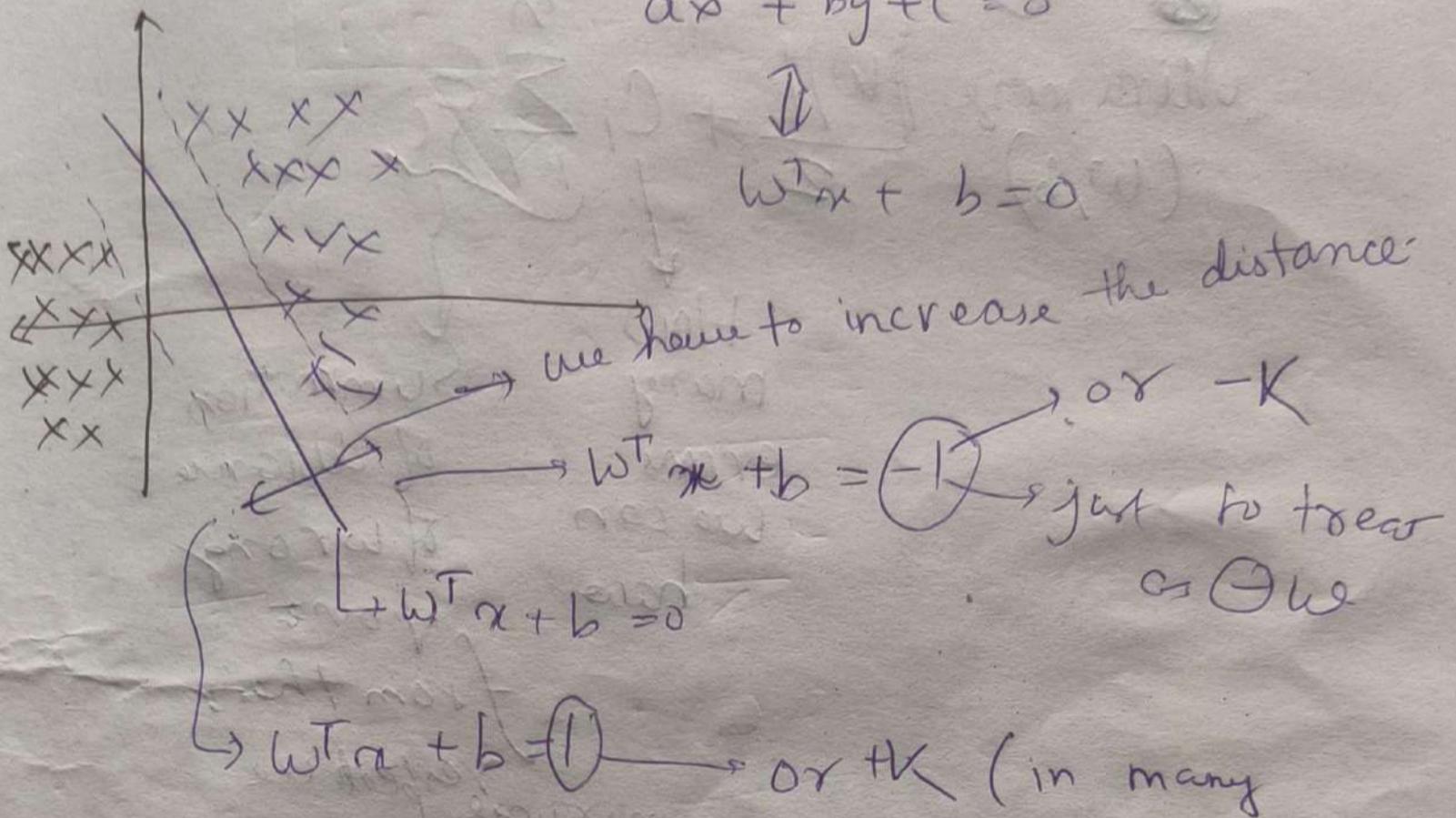


$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} 4, 4 \end{bmatrix} = -4$$

$$ax + by + c = 0$$



$$w^T x + b = 0$$



$$w^T x_1 + b = 1$$

$$w^T x_2 + b = -1$$

— — —

$$\underline{w^T(x_1 - x_2) = 2}$$

Any vector has magnitude

$$\frac{w^T(x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|}$$

Maximize $\frac{2}{\|w\|}$ such that $y_i = \begin{cases} +1 & w^T x_i + b \geq 1 \\ -1 & w^T x_i + b \leq -1 \end{cases}$

Map x to $y_i * (w^T x_i + b) \geq 1$ for correct points

$$\underset{(w,b)}{\text{Maximize}} \frac{2}{\|w\|} \Leftrightarrow \underset{(w,b)}{\text{Minimize}} \frac{\|w\|}{2}$$

~~but~~ \downarrow

$$\underset{(w,b)}{\text{Minimize}} \frac{\|w\|}{2} + C_i \sum_{j=1}^n \xi_j$$

How many errors we can have

Summation of distance of wrong point from their

marginal line of their group

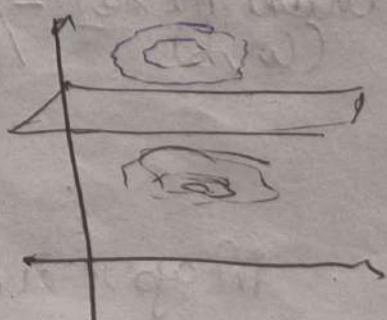
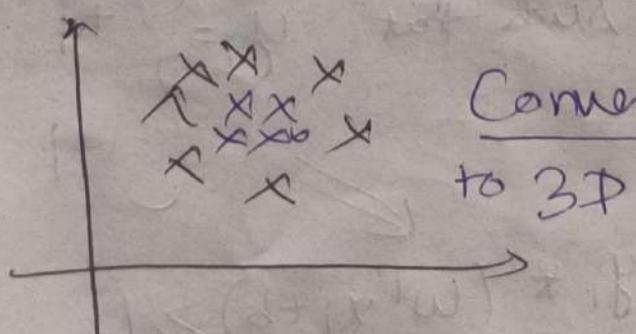
SVR

$$\underset{(w,b)}{\text{Minimize}} \frac{\|w\|}{2} +$$

Rest all are same as SVM.

Homework

SVM Kernel



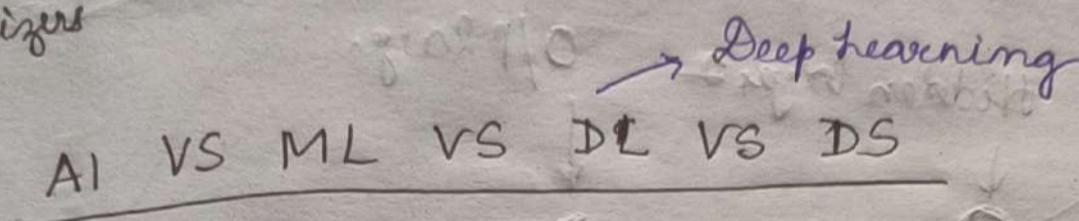
$$f(x_1, x_2) \cdot f^T(x_1, x_2) = K(x_1, x_2)$$

Kernel Function

DEEP LEARNING

Agenda

- ① Deep learning → Perception {AI vs ML vs DL vs DS}
- ② Forward propagation
- ③ Backward propagation
- ④ Loss function
- ⑤ Activation functions
- ⑥ Optimizers



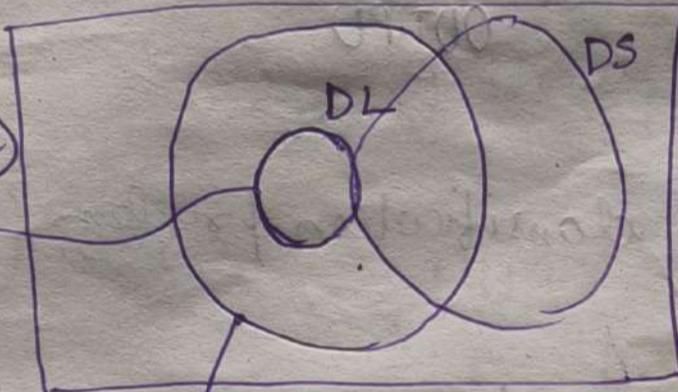
→ Deep learning

Researchers (1958)

Multi
layered
neural
network

{Perceptron}

{Mimics the
human
brain}



→ AI application which can do its own task without any human intervention

- ① Netflix app
- ② Self Driving Car
- ③ Amazon Application
- ④ Sophia

ML → Stats tool to analyze the data, visualize the data, predictions, forecasting, clustering (ML is subset of AI)

Why deep learning is becoming popular?

- ① 2005 → ORRUT, FACEBOOK, INSTA, WHATSAPP, LINKEDIN, TWITTER

DATA → exponentially↑

2008 → {Big Data} → Google

2013 → Company had huge amount of data

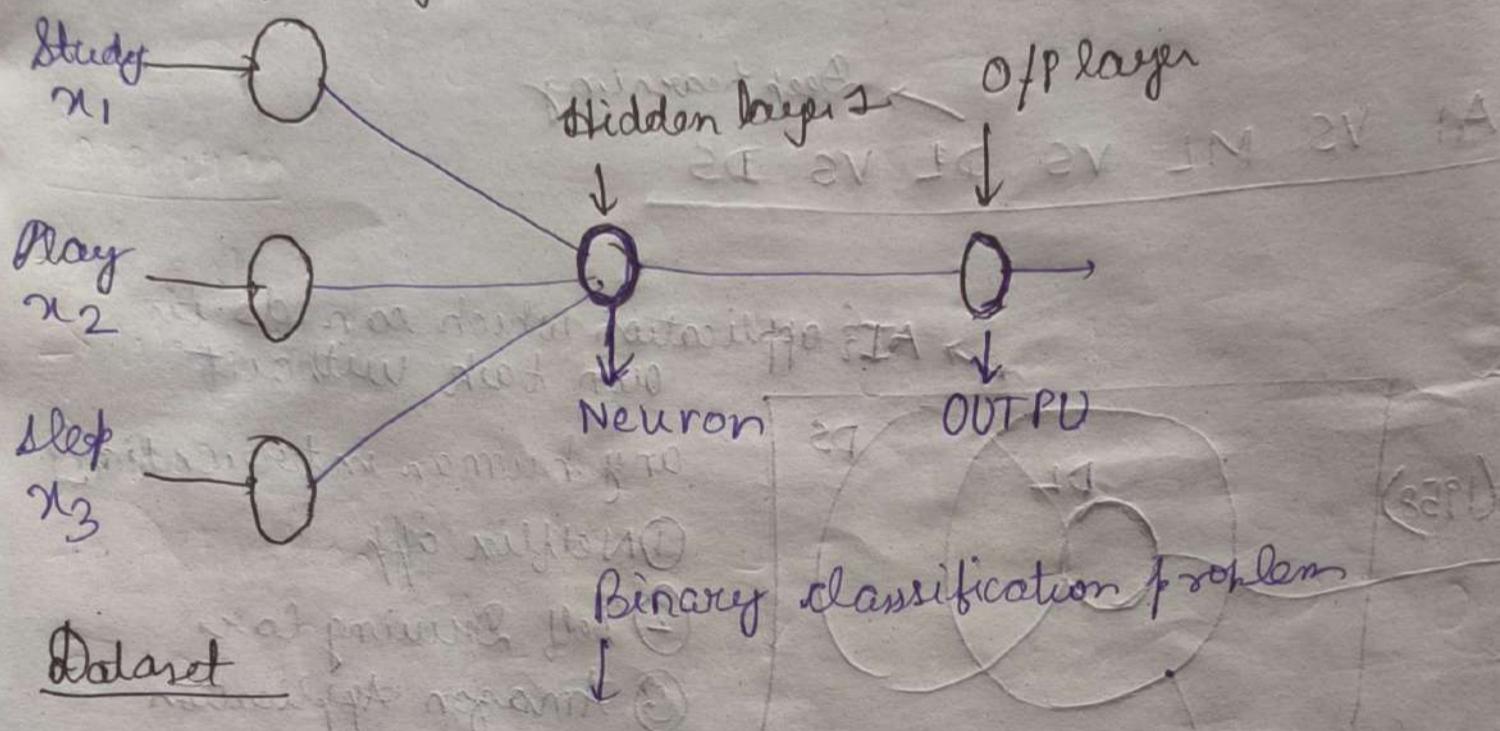
↓
Seamless products {AI became popular}

② Hardware advancements (NVIDIA) → GPU's
 ↓
 Training the model

GPU's → cost ↓↓

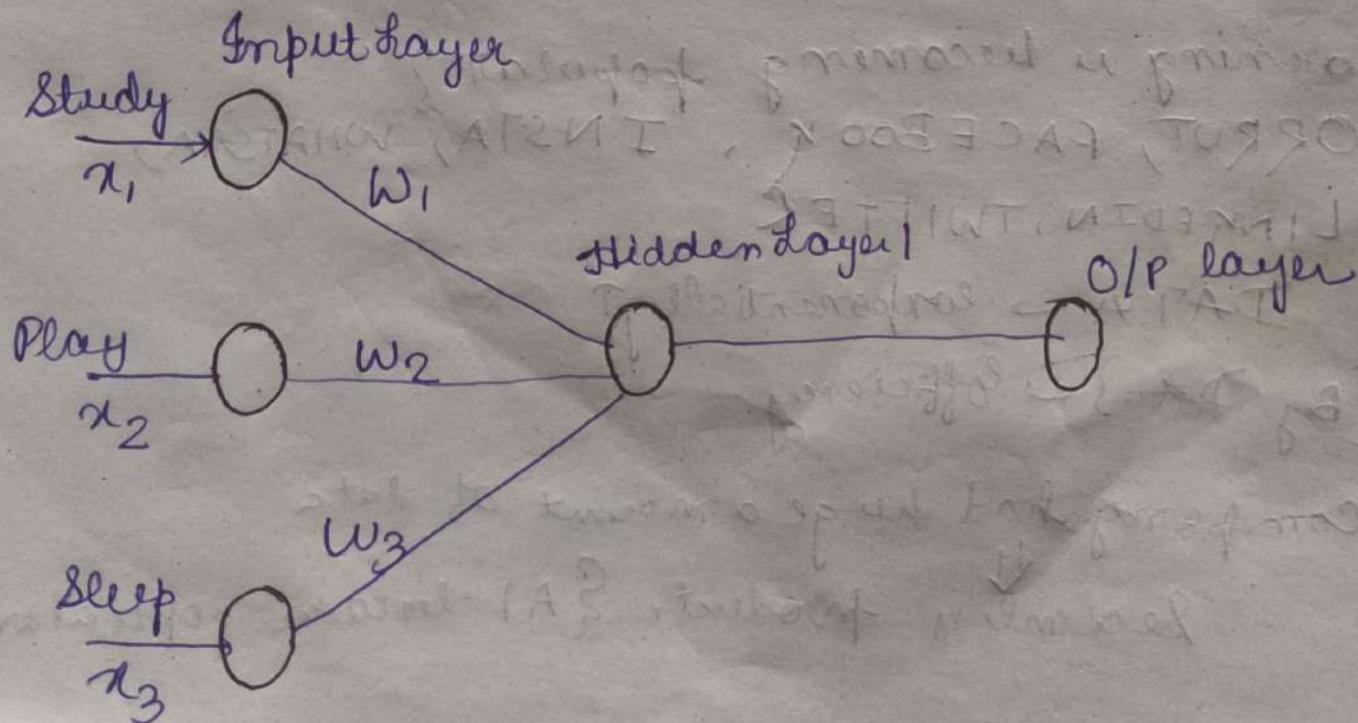
Perception { Singled Layer Neural Network}

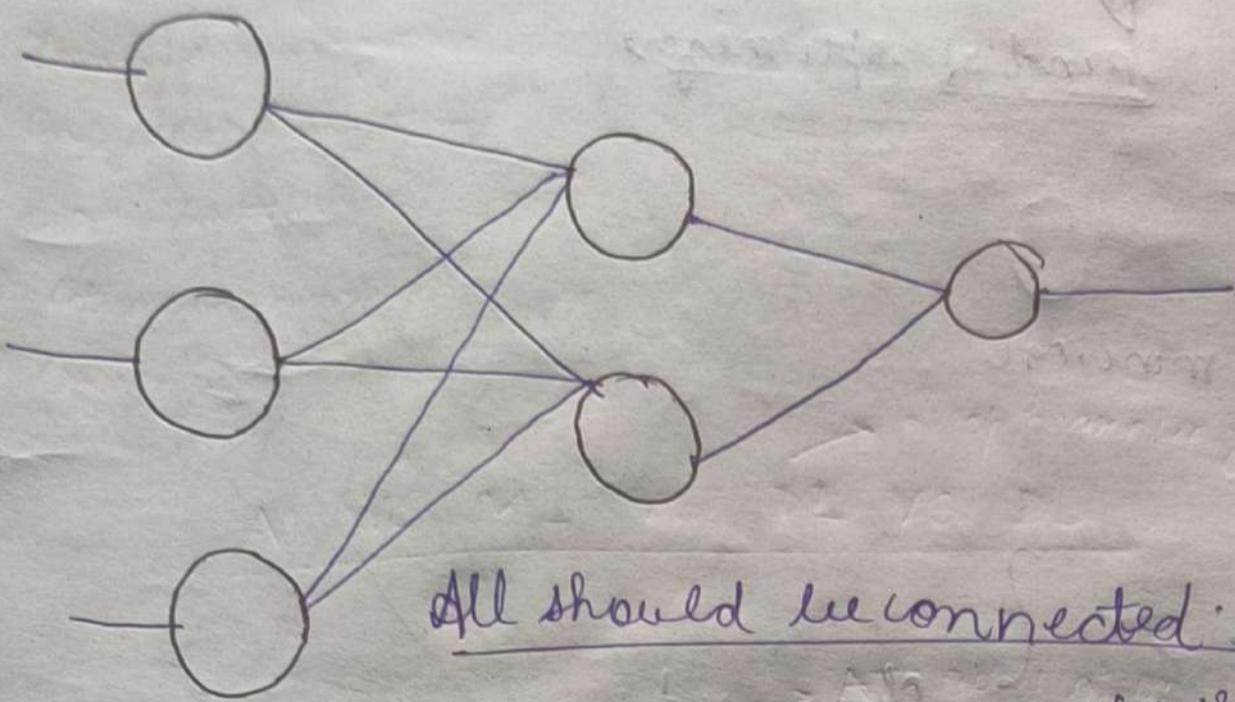
Input layer



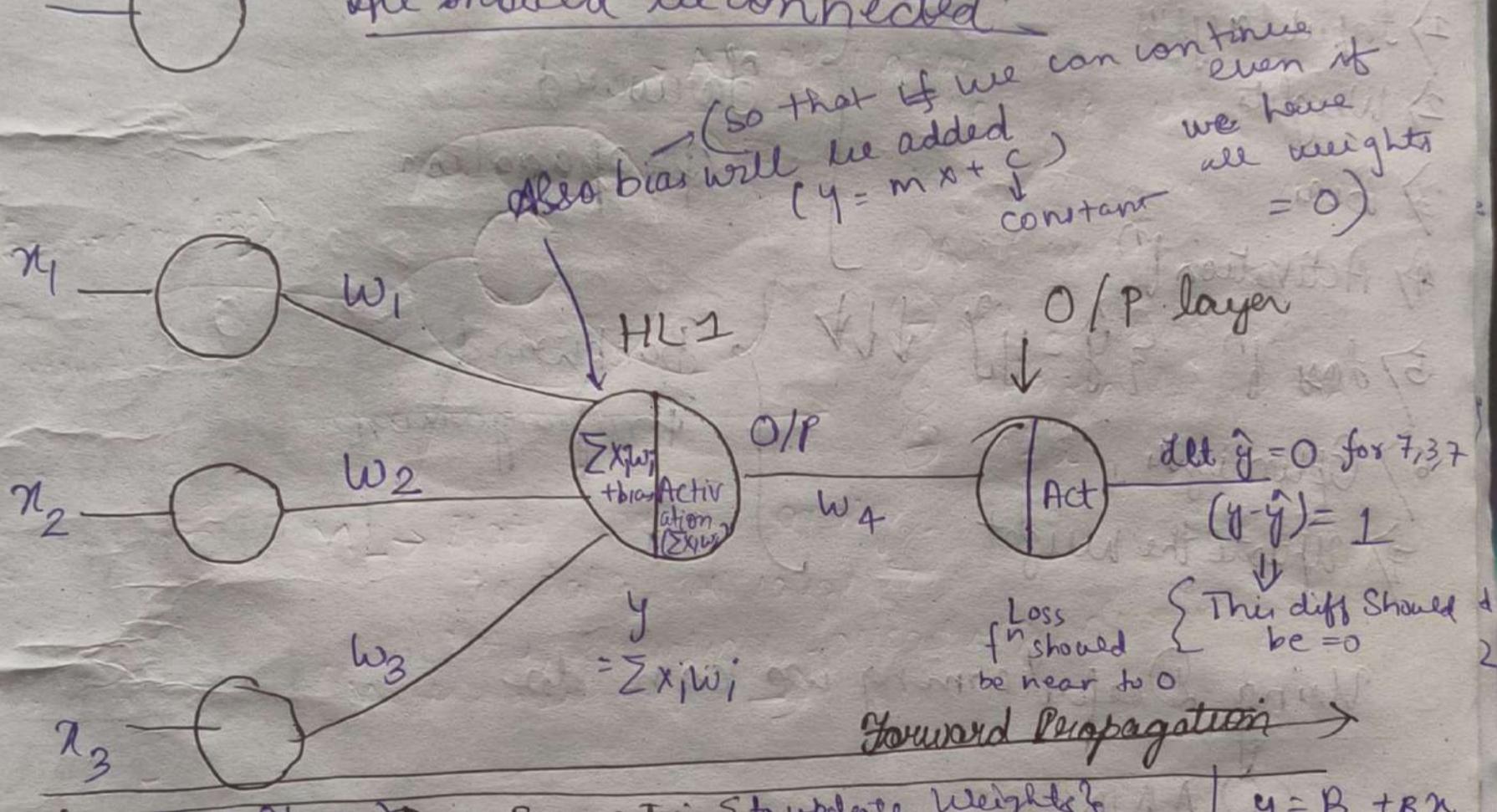
Dataset

	Study	play	sleep	Pass/Fail
7	3	7	7	Fail
2	5	8	8	Pass
4	3	7	1	Pass





All should be connected.



Forward Propagation →

← (Using Optimizers) Back Propagation { to update weights }

$$\sum x_i w_i = x_1 w_1 + x_2 w_2 + x_3 w_3 + b$$

Just like update in gradient descent in linear Regression

$$= w^T x + b$$

$$y = \beta_0 + \beta x$$

$$y = \beta^T x$$

$$y = mx + c$$

for Binary Classification

Sigmoid Activation = $\frac{1}{1 + e^{-y}} = \frac{1}{1 + e^{-(\sum w_i x_i + b)}} = [0.02]$

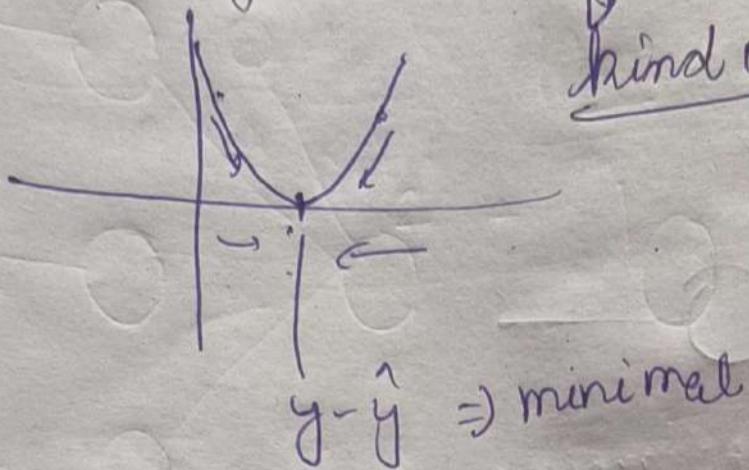
whether the neuron is activated or not

$$\begin{cases} > 0.5 = 1 \\ < 0.5 = 0 \end{cases}$$

In Simple Linear Regression,

Gradient Descent

A kind of optimizer



Conclusion

1) Input layers

2) Weights

3) Bias

4) Activation f^n

5) Loss $f^n \{ \hat{y} - y \}$

6) Optimizer

7) Update the weight

between forward & backward pass

forward

propagation

backward

propagation

Using Deep Learning we can do

① ANN

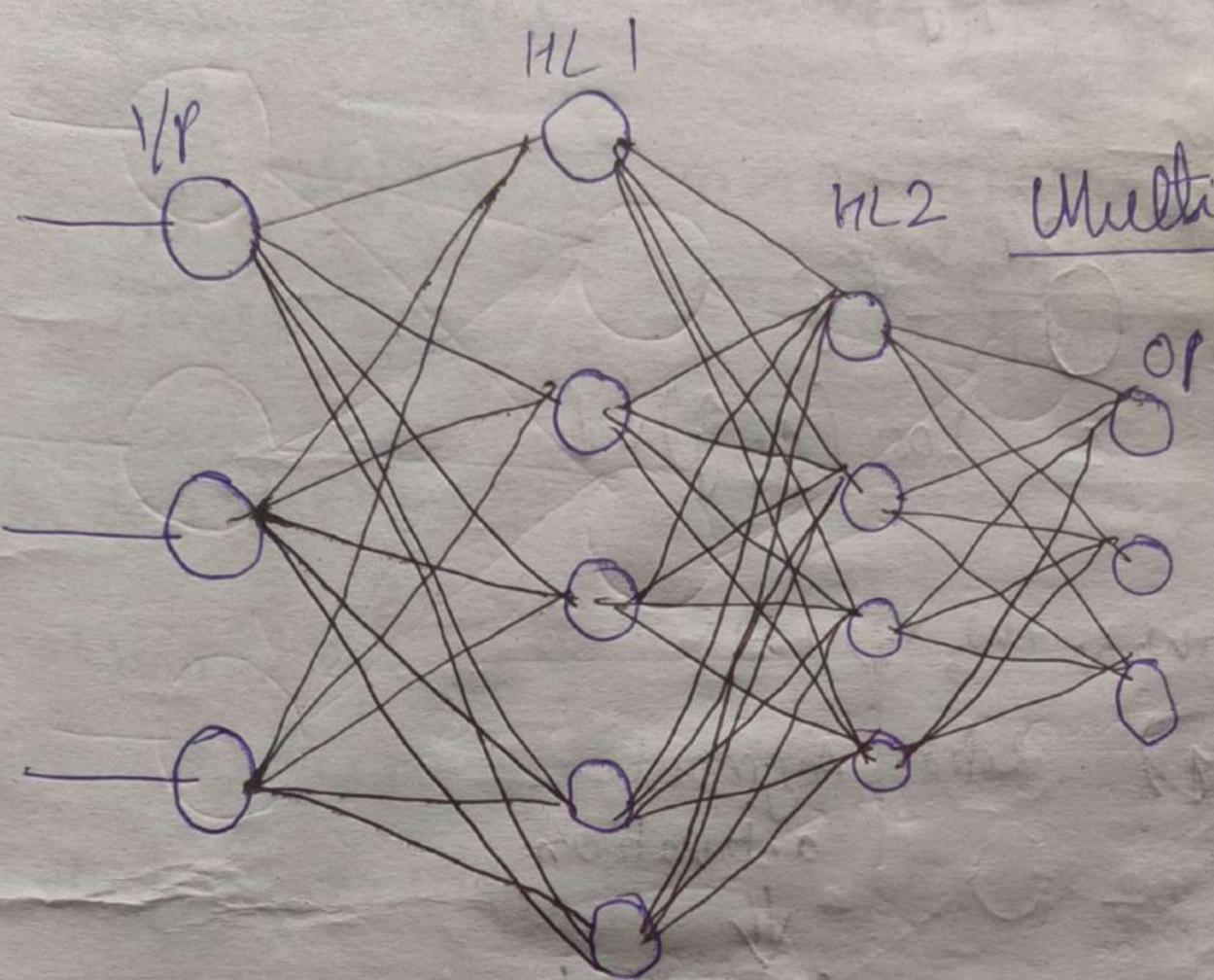
② CNN

③ RNN

④ Object Detection

Multi Layered Network

In DL, we can also have more than 1 neuron in output layer (medical classification)



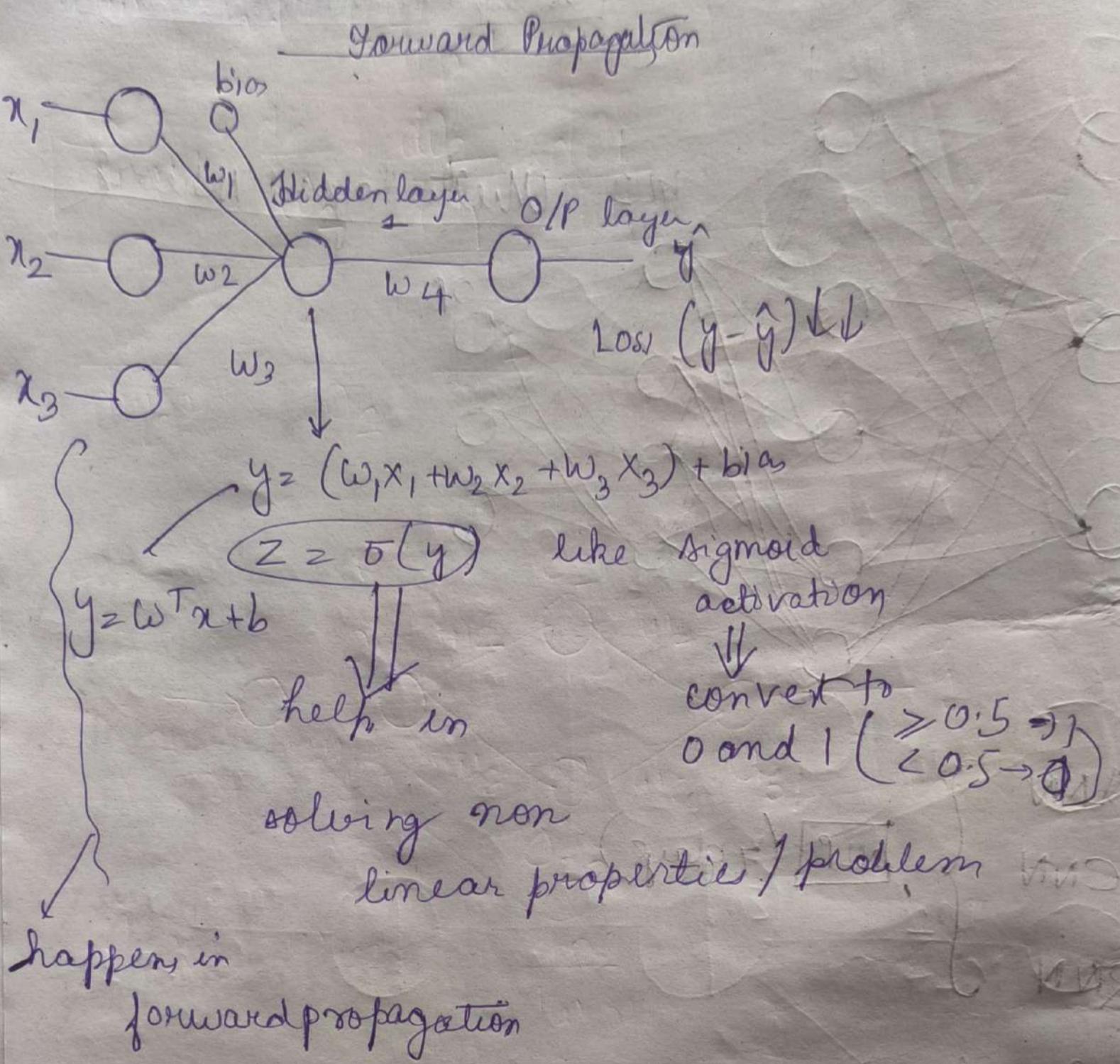
Multi Layered Network

- 1) ANN
 - 2) CNN
 - 3) RNN
- } NLP (7 days)

Day - 2 - Deep Learning

Agenda

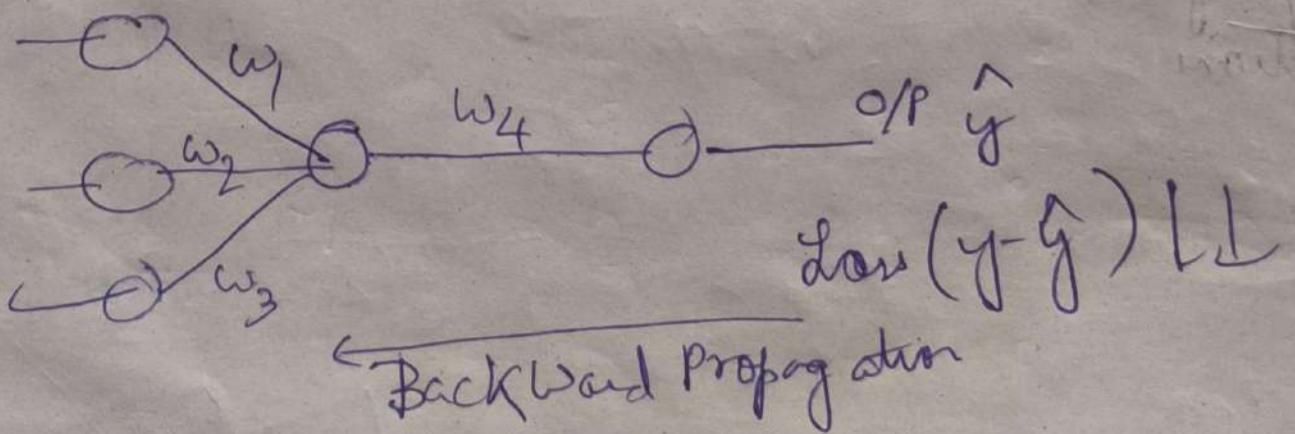
- ① Forward Propagation
- ② Chain rule of derivative
- ③ Vanishing gradient Problem
- ④ Loss functions



② Back Propagation

① Weight updation formula

② Chain Rule of Differentiation



Weight updation formula

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dh}{dw_{\text{old}}} \rightarrow \text{loss}$$

$$\frac{\partial h}{\partial w_{\text{old}}} \rightarrow \text{slope}$$

n is small number $\Theta 0.01$ or 0.001
If n is large, there will be jumps

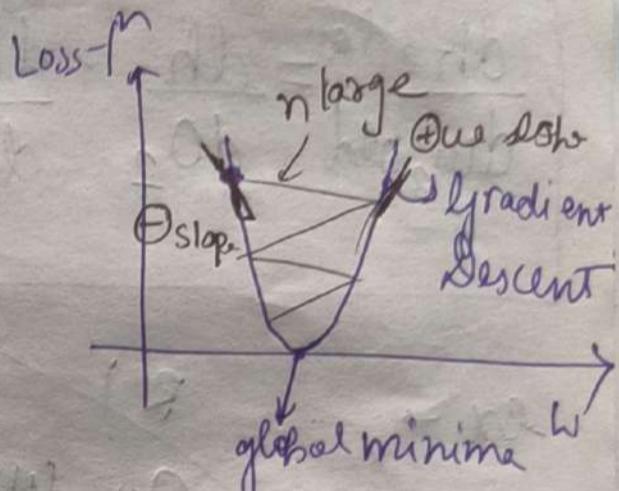
For simple linear regression

In case of One slope,

$$w_{\text{new}} = w_{\text{old}} - \eta (\Theta w)$$

$$= w_{\text{old}} + \eta (\text{some } \Theta w)$$

the $w_{\text{new}} \gg w_{\text{old}}$



In case of Θ ue slope,

$$w_{\text{new}} = w_{\text{old}} - n (\Theta w)$$

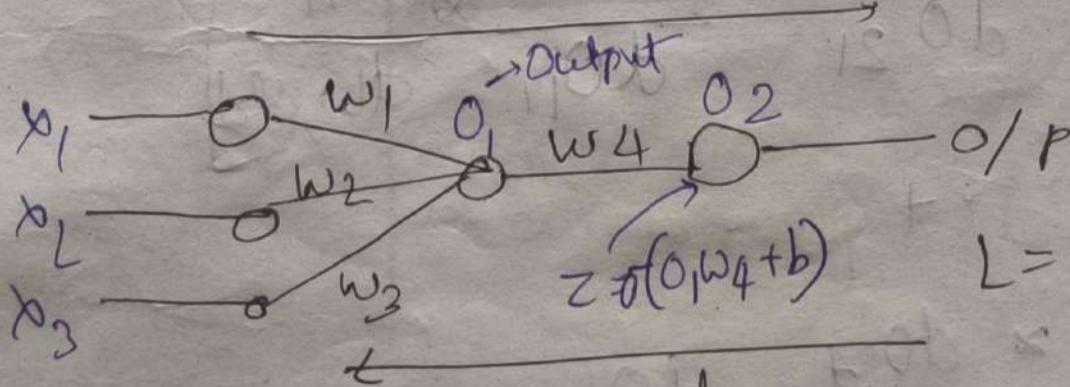
$$= w_{\text{old}} - n (\text{some } \Theta w)$$

the ~~b_{old}~~ $w_{\text{new}} \ll w_{\text{old}}$

$$b_{\text{new}} = b_{\text{old}} - n \frac{dh}{db_{\text{old}}}$$

bias updation formula

② Chain rule of Differentiation



$$w_{\text{new}} = w_{\text{old}} - n \frac{dh}{dw_{\text{old}}}$$

$$- n \frac{dh}{dw_{\text{old}}}$$

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

Backward Propagation

$$\bullet w_{4\text{ new}} = w_{4\text{ old}} - n \frac{dh}{dw_{4\text{ old}}}$$

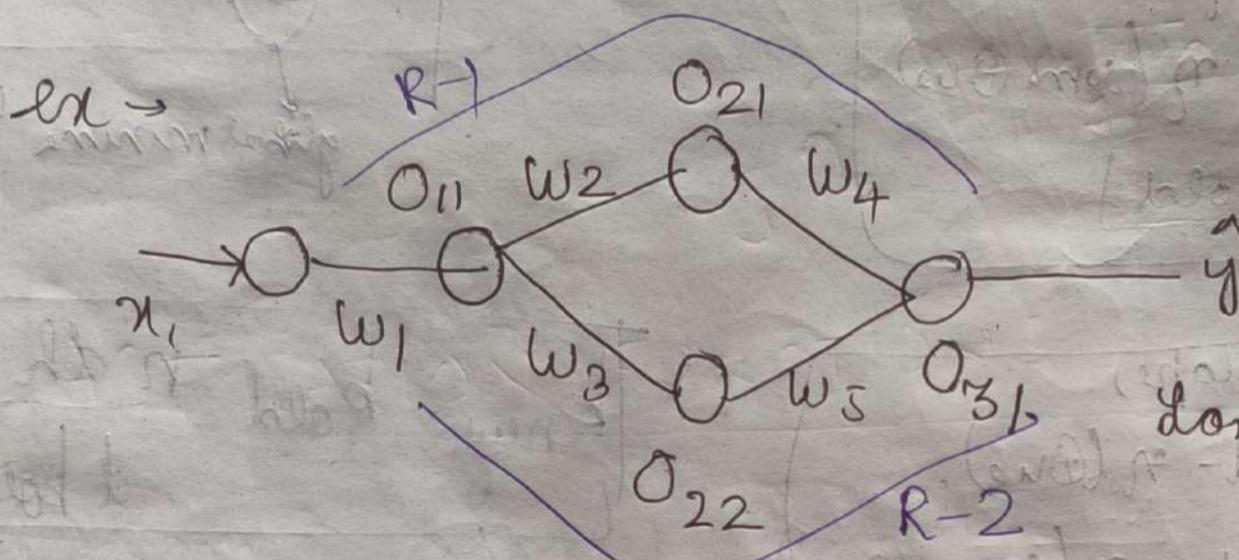
$$\left[\frac{dh}{dw_{\text{old}}} = \frac{dh}{dO_2} \cdot \frac{dO_2}{dw_{4\text{ old}}} \right]$$

$$w_{1,\text{new}} = w_{1,\text{old}} - n \frac{dh}{dw_{1,\text{old}}}$$

$$\frac{dh}{dw_{1,\text{old}}} = \frac{dh}{dO_2} \cdot \frac{dO_2}{dO_1} \cdot \frac{dO_1}{dw_{1,\text{old}}}$$

$$\frac{dO_2}{dO_1} = \frac{dO_2}{dw_4} \cdot \frac{dw_4}{dO_1}$$

$$\frac{dh}{dw_{2,\text{old}}} = \frac{dh}{dO_2} \cdot \frac{dO_2}{dO_1} \cdot \frac{dO_1}{dw_{2,\text{old}}}$$



$$w_{1,\text{new}} = w_{1,\text{old}} - n \frac{dh}{dw_{1,\text{old}}}$$

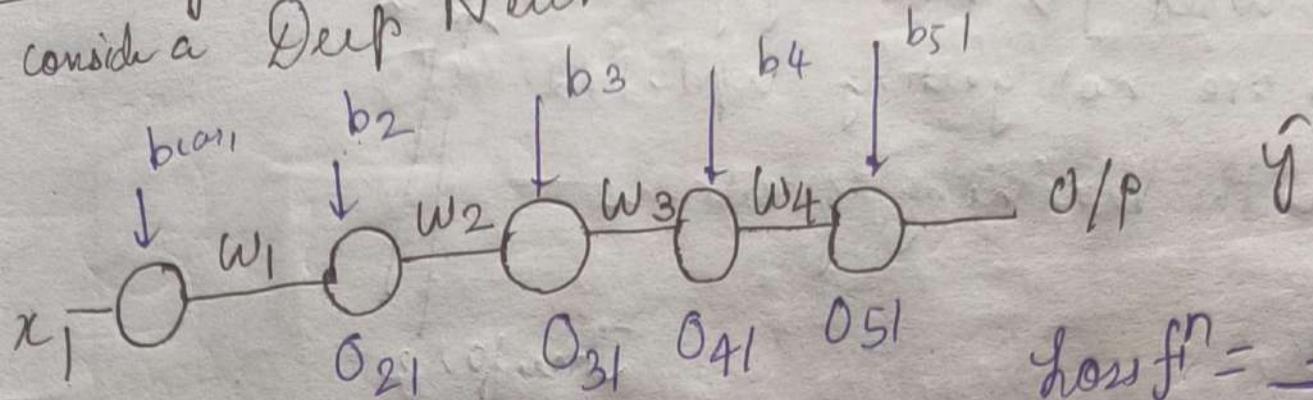
$$\frac{dh}{dw_{1,\text{old}}} = \left[\frac{dh}{dO_{31}} \times \frac{dO_{31}}{dO_{21}} \times \frac{dO_{21}}{dO_{11}} \times \frac{dO_{11}}{dw_{1,\text{old}}} \right]$$

Chain rule of derivatives

$$\frac{dh}{dO_{31}} = \left[\frac{dh}{dO_{22}} \times \frac{dO_{22}}{dO_{21}} \times \frac{dO_{21}}{dO_{11}} \times \frac{dO_{11}}{dw_{1,\text{old}}} \right]$$

Vanishing Gradient Problem

We consider a Deep Neural Network,



$$\text{Loss fn} = \frac{1}{2} (y - \hat{y})^2 \quad (\text{mse})$$

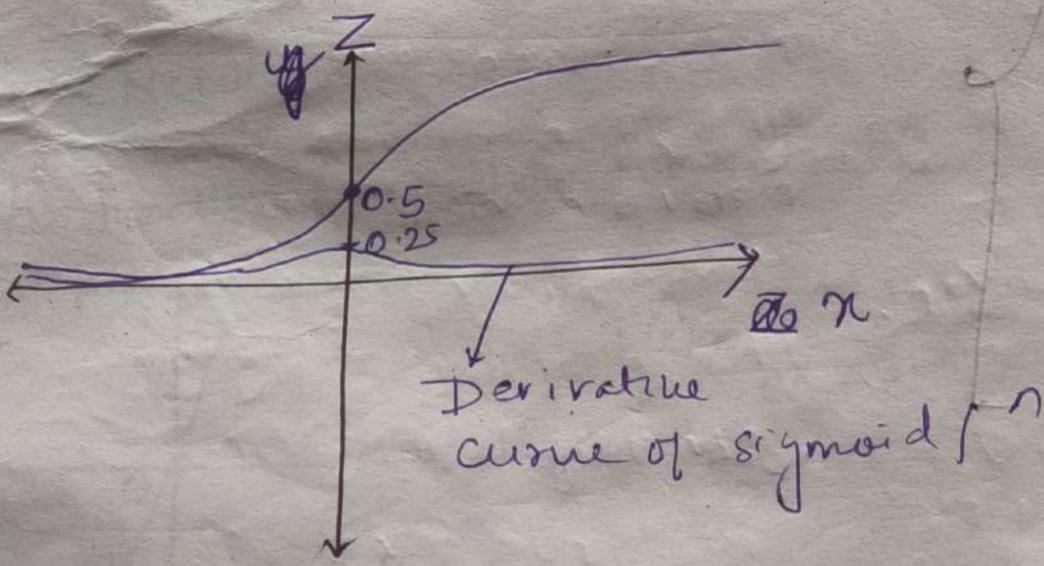
Sigmoid Activation

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dh}{dw_{\text{old}}}$$

$$\frac{dh}{dw_{\text{old}}} = \frac{dh}{dO_{51}} \times \frac{dO_{51}}{dO_{41}} \times \frac{dO_{41}}{dO_{31}} \times \frac{dO_{31}}{dO_{21}} \times \frac{dO_{21}}{dw_{\text{old}}}$$

$$\text{Sigmoid Activation} \quad \left(\frac{y}{z} = \frac{1}{1 + e^{-x}} \right)$$

$\rightarrow 0 \text{ or } 1$
 $> 0.5 \rightarrow 1$
 $< 0.5 \rightarrow 0$



Derivative Cond'n

$$0 \leq \sigma(y) \leq 0.25$$

for ex-

$$\frac{dh}{dw_{\text{old}}} = \frac{dh}{dO_{51}} \times \frac{dO_{51}}{dO_{41}} \times \frac{dO_{41}}{dO_{31}} \times \frac{dO_{31}}{dO_{21}} \times \frac{dO_{21}}{dw_{\text{old}}}$$

Let $= 0.25 \times 0.15 \times 0.10 \times 0.05 \times 0.12$
 $= \text{Small value}$

$$w_{\text{new}} = w_{\text{old}} - \eta / \text{small number}$$

So,
 $O_{51} = \sigma((O_{41} \times w_4) + b)$
 So, it is for all case
 $O_{51}, O_3, \text{ etc } 0.25$

So their derivative range from 0 to 0.25

So, after some time w_{new} will hardly change
No change $\in [w_{new} \approx w_{old}]$ \Rightarrow Vanishing Gradient Problem
in weights

So, weights are not getting updated.

So, this is the Vanishing Gradient

Problem.

Therefore

to solve

Vanishing Gradient

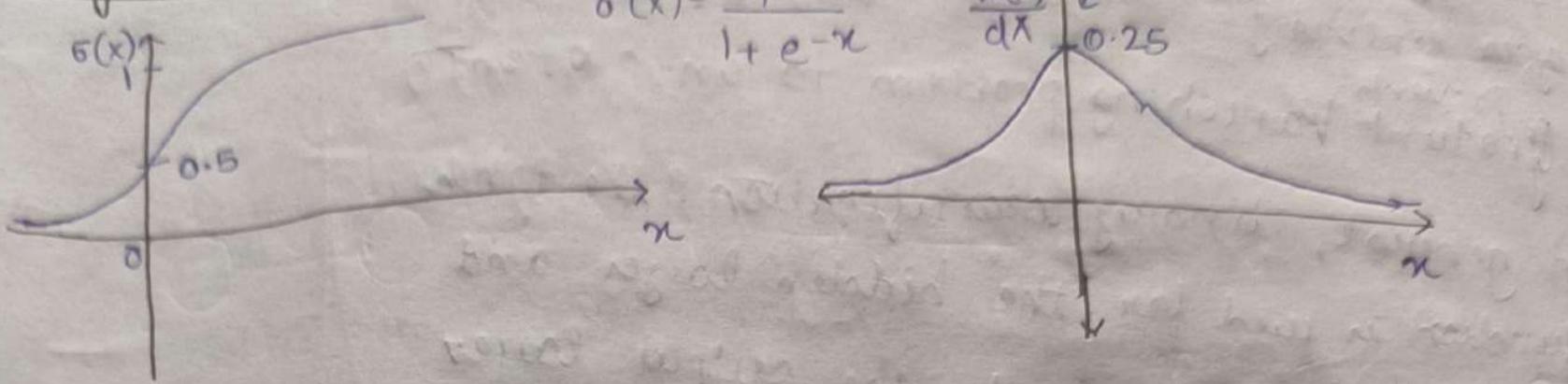
Problem, we use
another Activation

fn.

ACTIVATION FUNCTIONS

- ① Sigmoid
- ② Tanh
- ③ ReLU
- ④ Leaky ReLU
- ⑤ PReLU

① Sigmoid fⁿ



Advantages:

1) Smooth gradient, preventing jumps in output values
(& Easy Convergence of Weights)

(But the fⁿ is not centered 0, so weight update efficiency will get reduce.)

2) Output values bound between 0 and 1, normalizing the output of each neuron

3) Clear Predictions i.e. very close to 1 or 0

zero centered curve

Disadvantages:

1) Prone to Gradient Vanishing

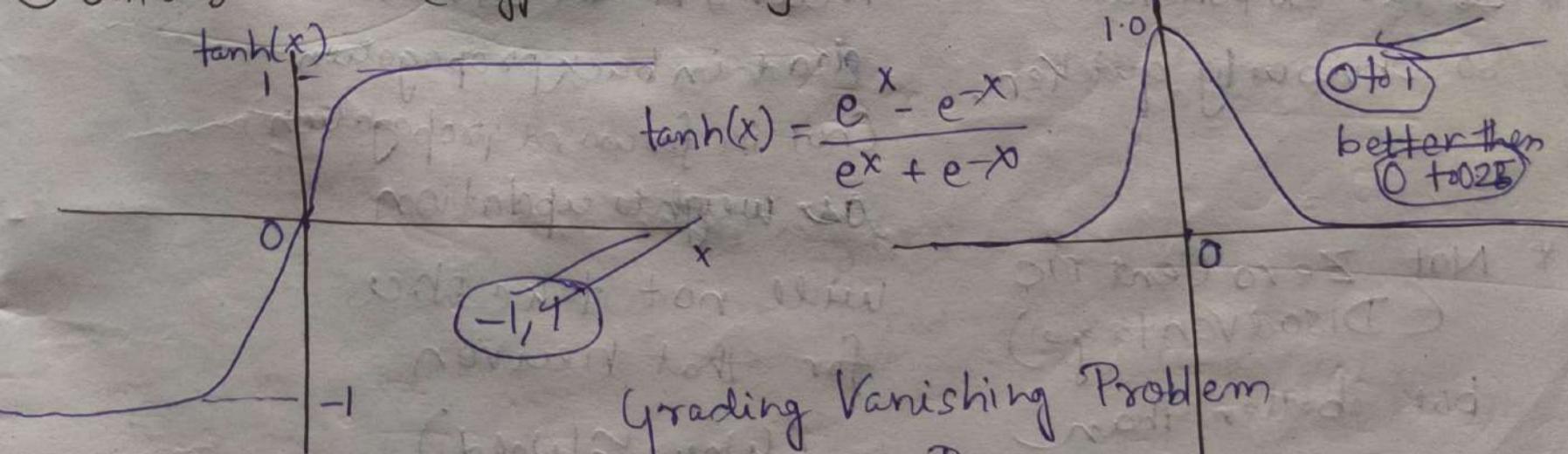
2) Function output is not zero-centered

3) Power ops. are relatively time consuming

So, due to Gradient Vanishing Problem, we cannot create Deep Neural Network.

② tanh function (hyperbolic tangent fⁿ)

$\tanh(x)$



Grading Vanishing Problem
is also Present

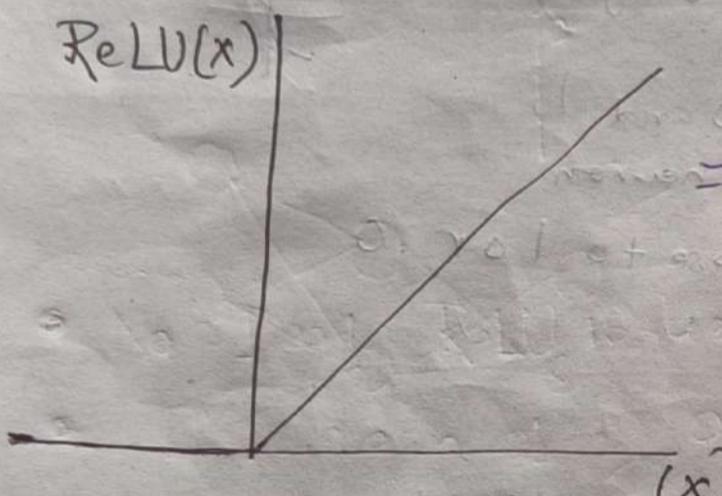
if we create Deep Deep
Neural Network.

Advantages:

- ① Zero centered fn
- ② Prevents Gradient Vanishing problem to an extent

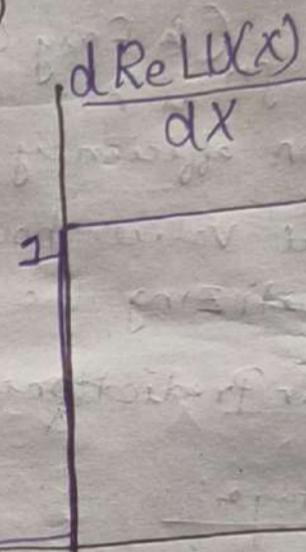
In general, binary classification, the tanh function is used for the hidden layer and the sigmoid fn is used for the output layer

3. ReLU function (Rectified linear unit)



$$\text{ReLU} = \max(0, x)$$

(x)



x

Advantages:

- * Exceptional Performance
- * Solving the Gradient Descent Problem
- * So No exponential so obviously quicker.
- * Not Zero Centric (Disadvantage) but better than Sigmoid and tanh function.

Disadvantage:

- * If one of the derivative becomes 0, as a derivative in ReLU can be either 0 or 1, then the whole neuron becomes dead in back propagation (Not in forward propagation) as weights update will not take place for that neuron

($w_{new} = w_{old}$)

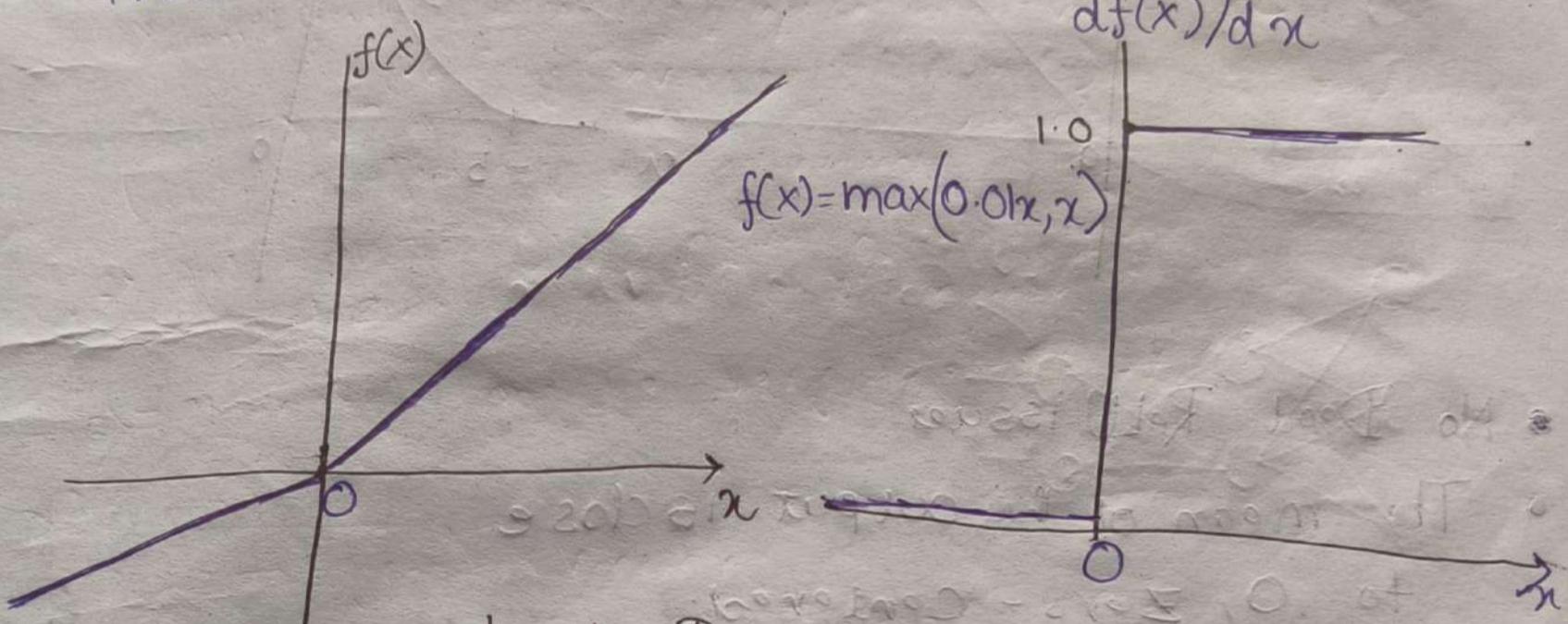
$$\left(\frac{dh}{dw_{old}} = 0 \right)$$

→ Dead

* Once a negative number is entered ReLU is inactive, and ReLU will die (Disadvantage).

4. Leaky ReLU ~~ReLU~~ function

→ To solve the Problem of dead Neuron / ReLU Problem.



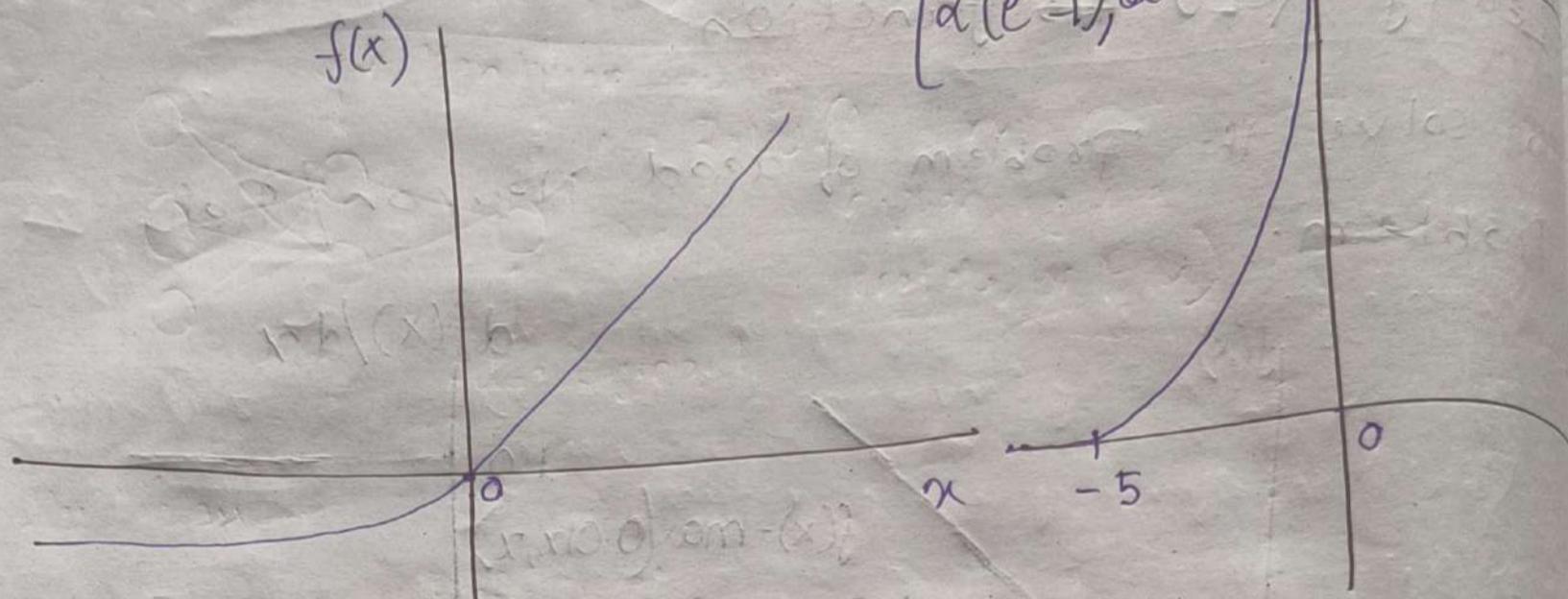
- In theory, Leaky ReLU has all the advantages of ReLU.
- But in actual ops, it has not been fully proved that Leaky ReLU is always better than ReLU.
Plus no proof of avoiding the problem of Dead ReLU

- Another intuitive idea is a Parameter based method.

Parametric ReLU : $f(x) = \max(\alpha x, x)$
which α can be learned from back propagation.

5. Exponential Leaky Units function (ELU)

$$f(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha(e^x - 1), & \text{otherwise} \end{cases}$$



- No Dead ReLU issues
- The mean of the output is close to 0, zero-centered.

One small problem is that it is slightly more computationally intensive. Similar to Leaky ReLU,

although theoretically better than than ReLU, there is currently

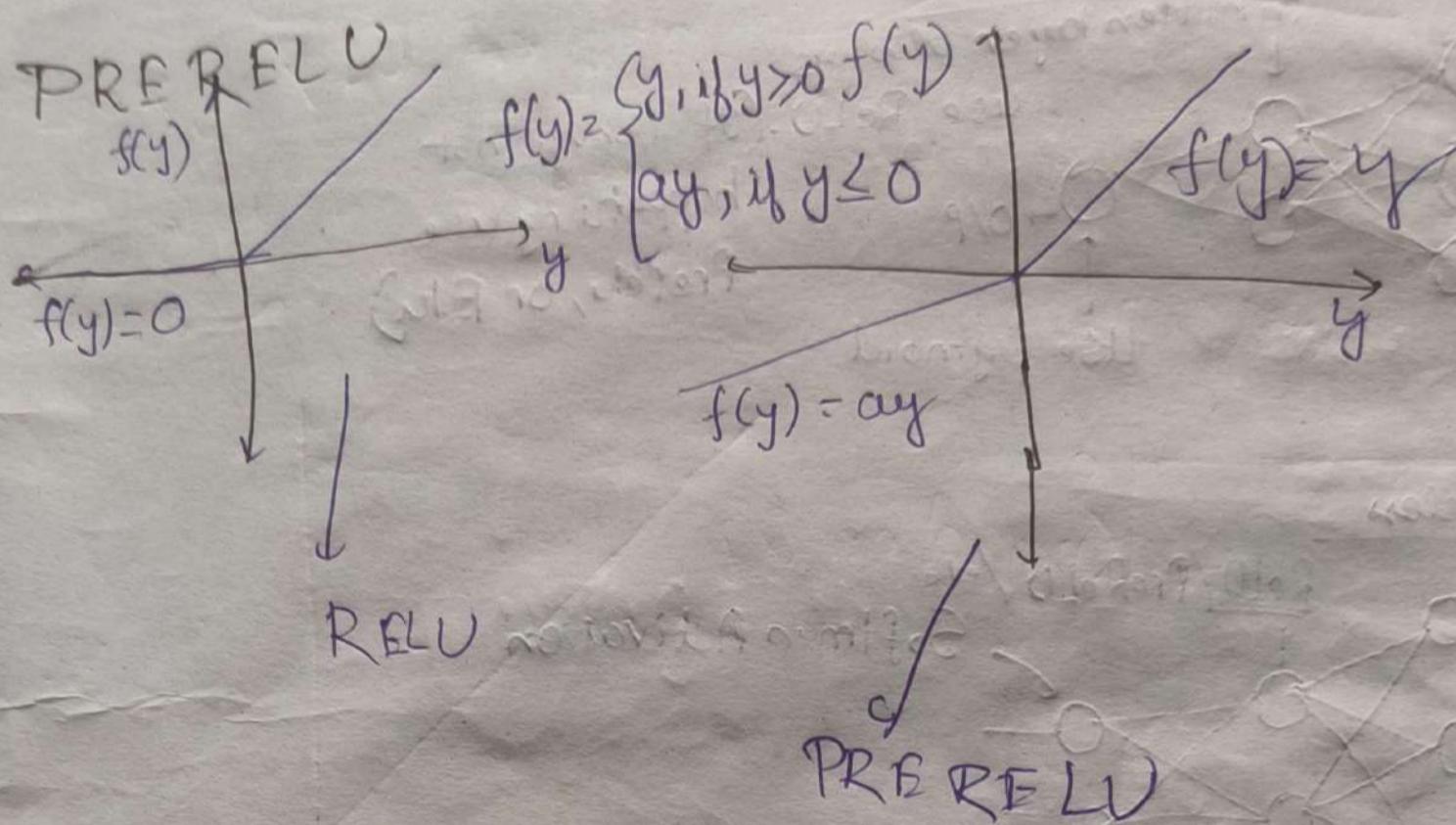
no good evidence in practice

that ELU is always better than ReLU.

6. Softmax

Will Discuss in Activation
fn

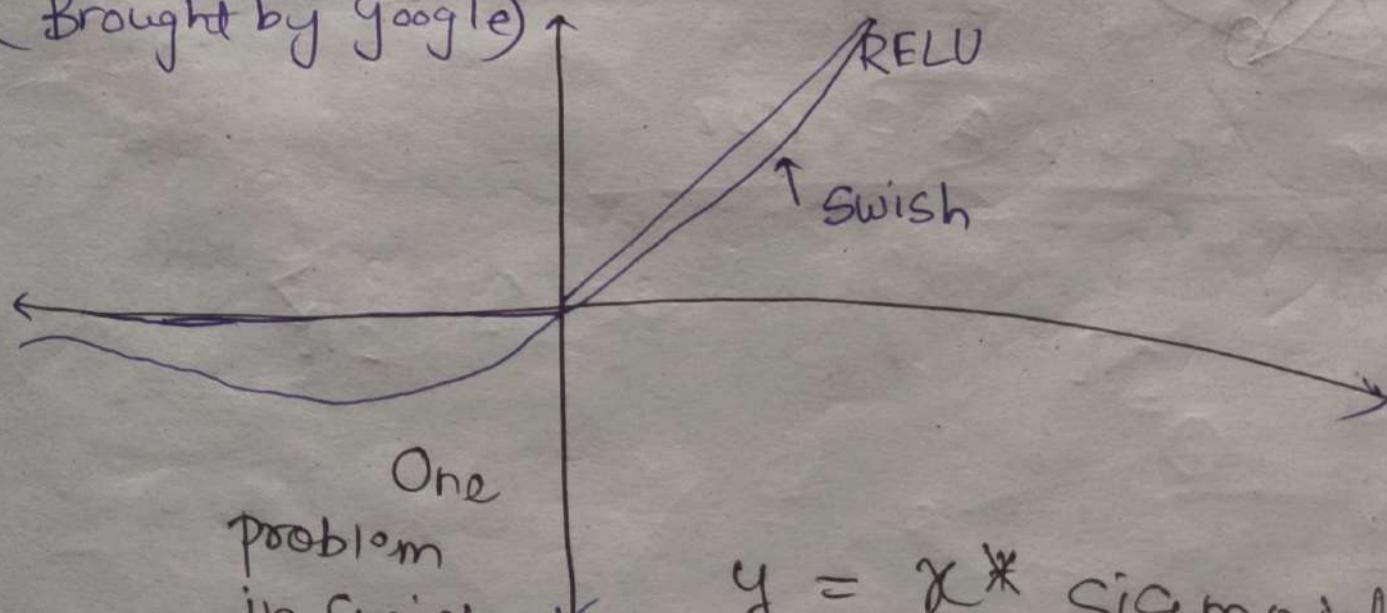
7. PRE RELU



- if $a=0$ it becomes ReLU
- if $a>0$, it becomes leaky ReLU
- if a is a learnable parameter, it becomes

PRelu

8. Swish (A Self-Gated) Function (Brought by Google)



One problem in Swish

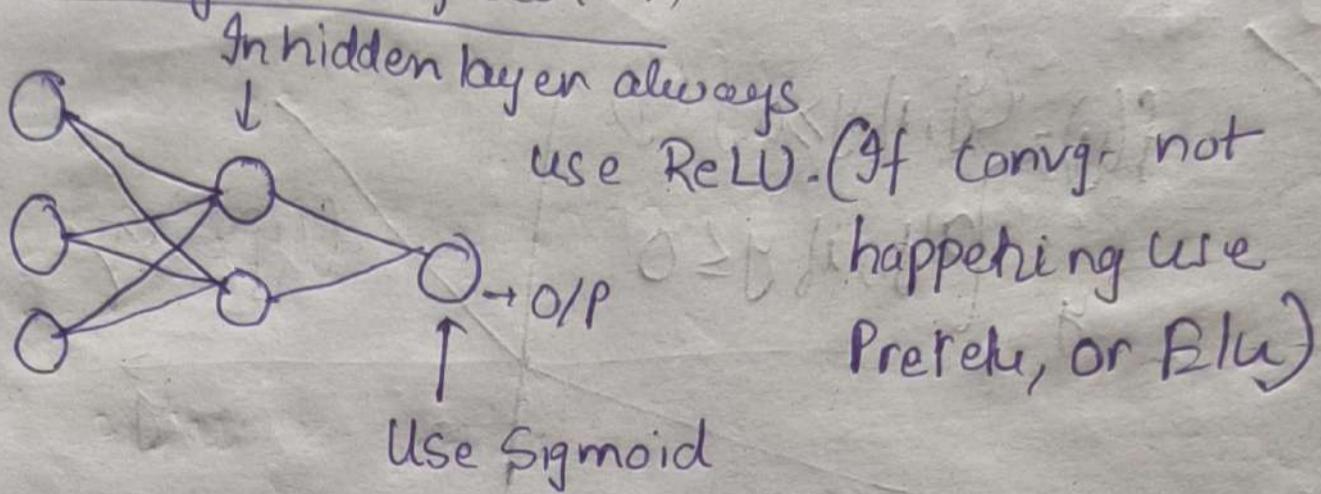
$$y = x * \text{Sigmoid}(x)$$

We cannot find derivative at 0.

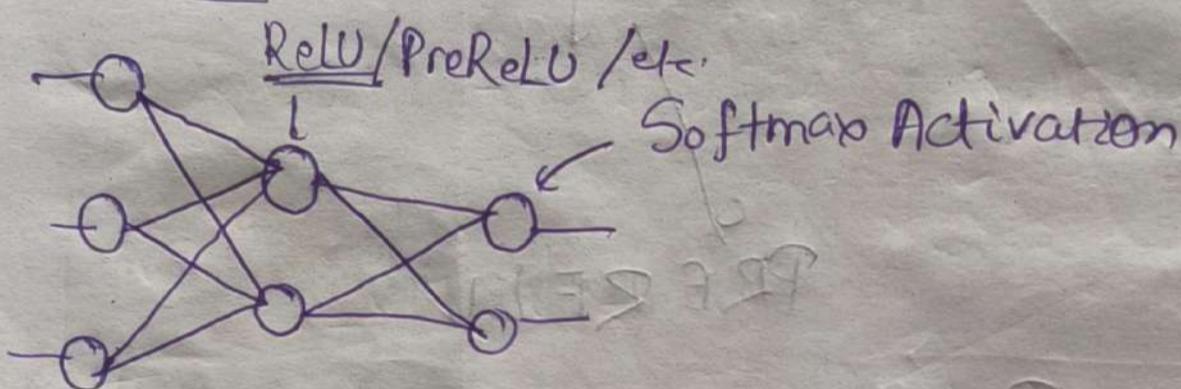
5

Technique which activation f^n should we use ???

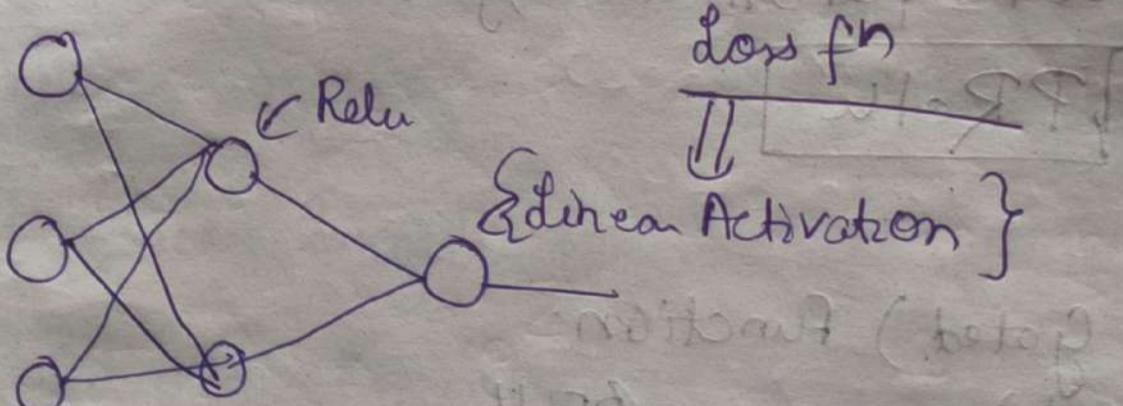
For Binary Classification



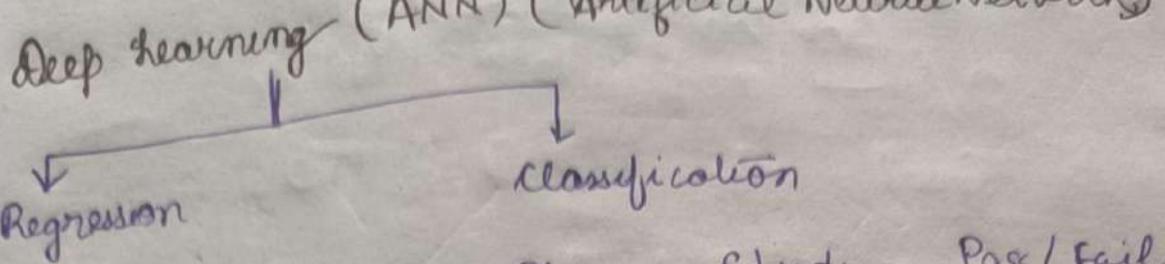
Multiclass



Regression



Loss functions



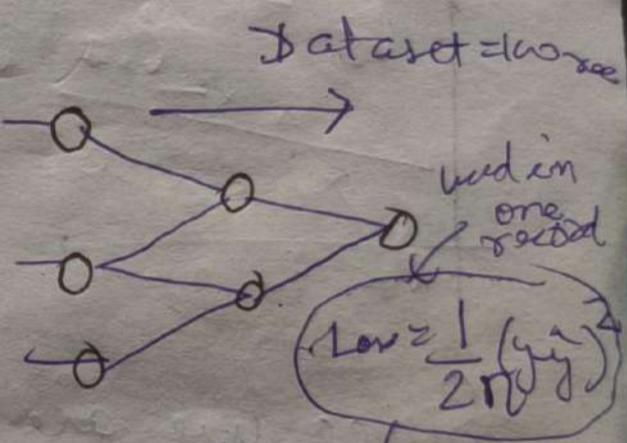
Exp	Deg	Salary	Play	Study	Pass / Fail
10	Pnd	-	10	2	F
-	-	-	4	3	F
-	-	-	5	5	maybe
-	-	-	2	7	Pass

Regression Problem

① Loss f^n in Regression

- 1) MSE (Mean Squared Error)
- 2) MAE (Mean Absolute Error)
- 3) Huber Loss

Loss f^n and cost f^n



① Mean Squared error (MSE).

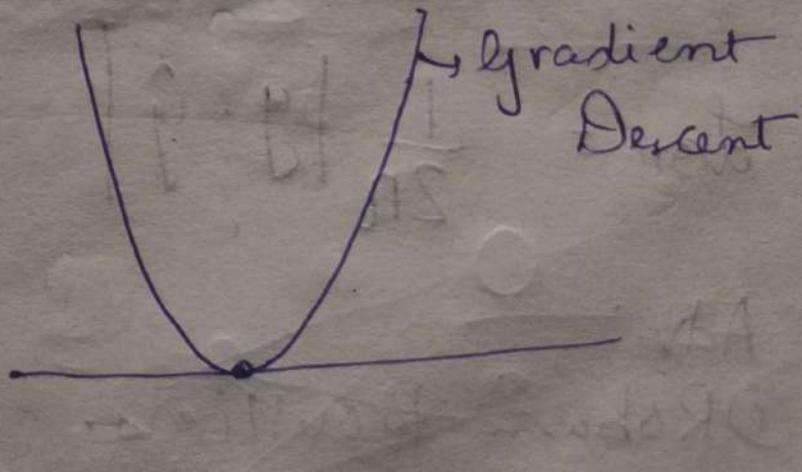
$$\text{Loss } f^n = \frac{1}{2n} (y - \hat{y})^2 \quad \text{cost } f^n = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

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

↓
Quadratic Eq^n

$$(a-b)^2 = a^2 + b^2 - 2ab$$

$$ax^2 + bx + c$$

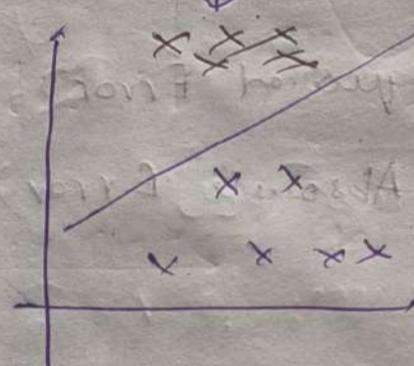
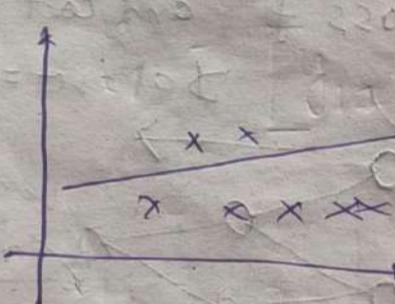


Advantages of using a Quadratic Curve

- ① Differentiable
- ② It has only one local or global minima
- ③ It converges faster.

Disadvantages:

⇒ Not Robust to outliers



Since we are penalizing and squaring the error, the line has a major shift on introduction of outliers.

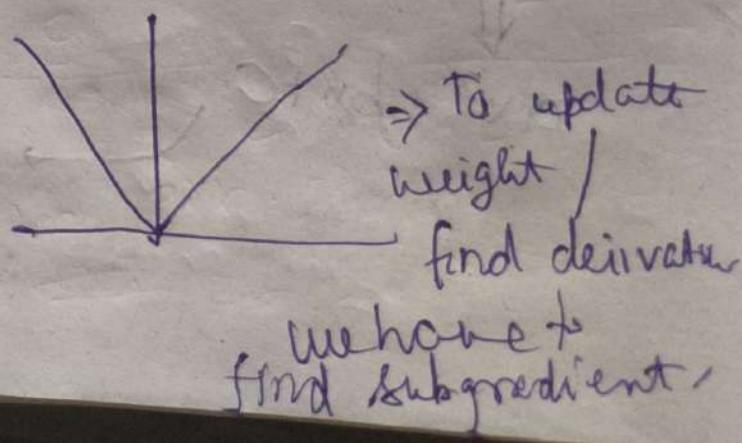
Mean Absolute Error

$$L_{\text{MAE}} = \frac{1}{2n} \sum |y - \hat{y}|$$

$$\text{Cost fn} = \frac{1}{2n} \sum_{i=1}^n |y - \hat{y}|$$

Adv:

- ① Robust to outliers.
(minor shift)



Dis :- Time Consuming ↑

③ Huber Loss:-
Combination of MSE or MAE.

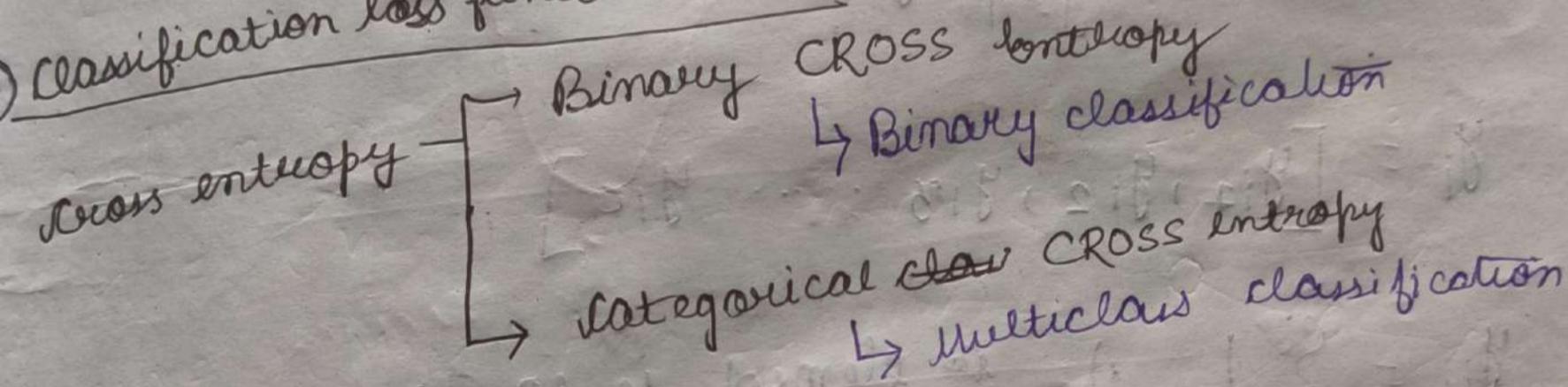
① MSE

② MAE

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

Hyper parameter
i.e. when outliers are not present
(i.e. when outliers are present)

Classification loss functions



Binary CROSS ENTROPY

$$\text{Loss} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \Rightarrow \text{Logistic Regression}$$

$$\text{Loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y = 0 \\ -\log(\hat{y}) & \text{if } y = 1 \end{cases}$$

Binary classification

$$\hat{y} = \frac{1}{1+e^{-x}}$$

i → row j = column

② Categorical class entropy { Multiclass Classification Problem
One hot Encoding
Bad Network

f_1	f_2	f_3	O/P → 1/pod	0	1	0
2	3	4	Good	1	0	0
5	6	7	Bad	0	1	0
8	9	10	Neutral	0	0	1

$c \leftarrow$ no. of categories

$$\text{Loss}(x_i, y_i) = -\sum_{j=1}^c y_{ij} * \ln(\hat{y}_{ij})$$

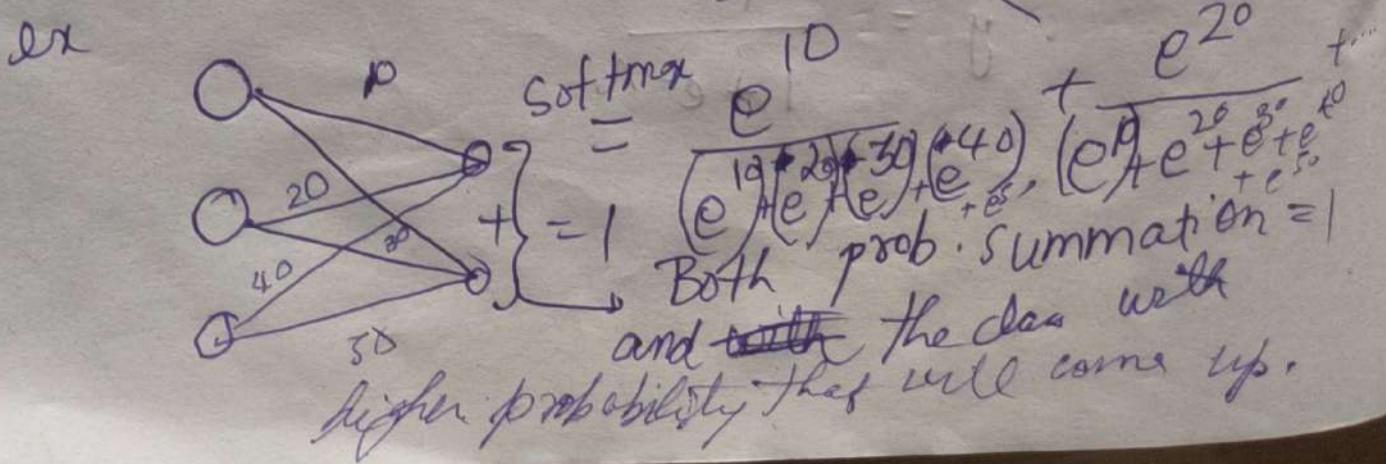
$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}]$$

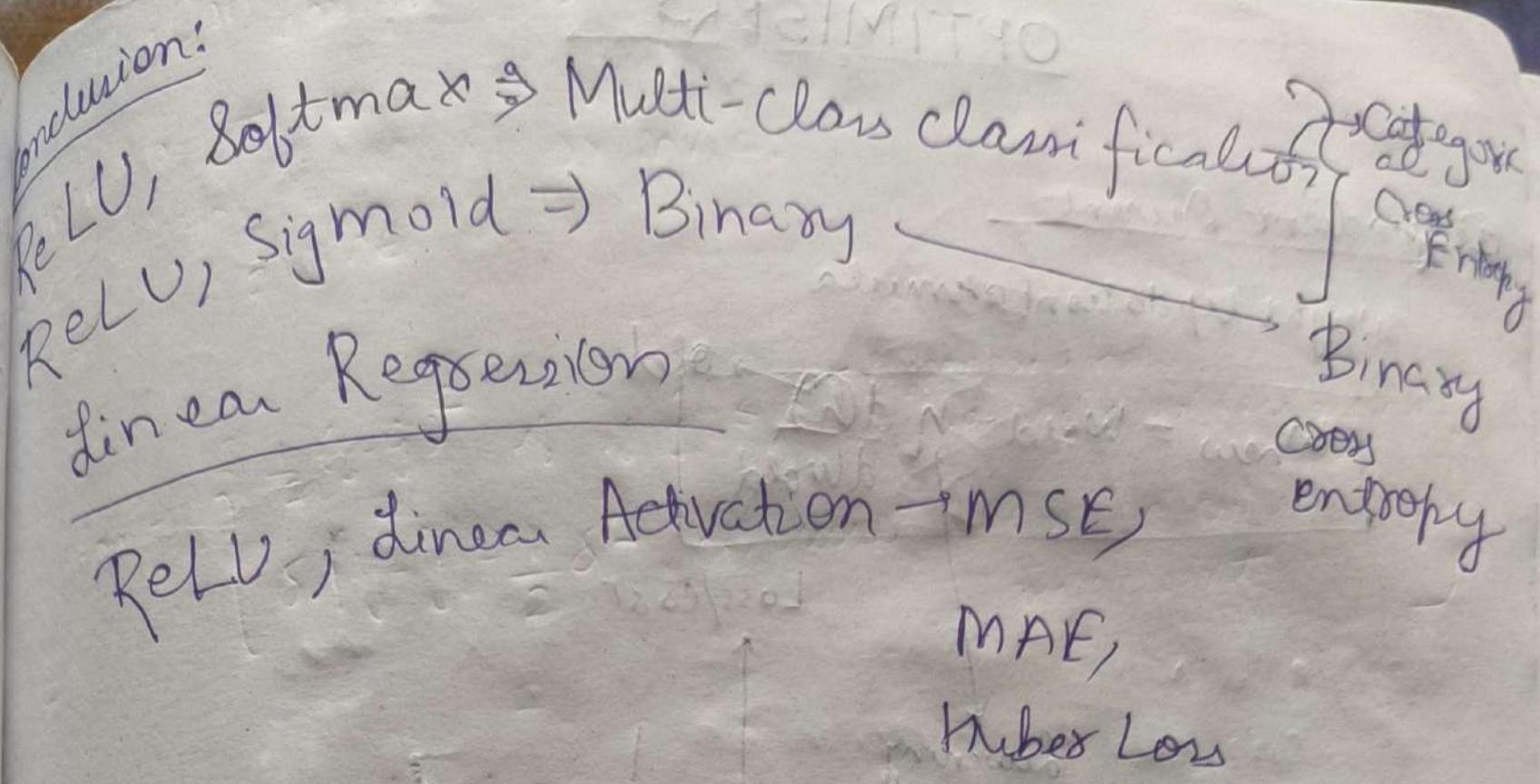
$$y_{ij} = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ element is in class} \\ 0 & \text{otherwise} \end{cases}$$

\hat{y}_{ij} = Softmax Activation layer classification

$$\sigma(z) = \frac{e^{z_i}}{\sum_{j=1}^c e^{z_j}} \rightarrow \text{Softmax activation}$$

probability 0.4, 0.5





OPTIMIZERS AND ANN IMPLEMENTATION

agenda :-

- 1 Optimizers
- 2 Gradient Descent
- 3 SGD (stochastic gradient descent)
- 4 MiniBatch SGD
- 5 SGD with momentum
- 6 Adagrad
- 7 RMSProp
- 8 Adam optimizers

Batch, Epochs, Iterations

ANN

OPTIMISERS

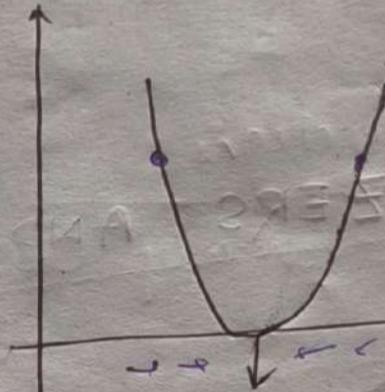
① Gradient Descent

Weight updation formula

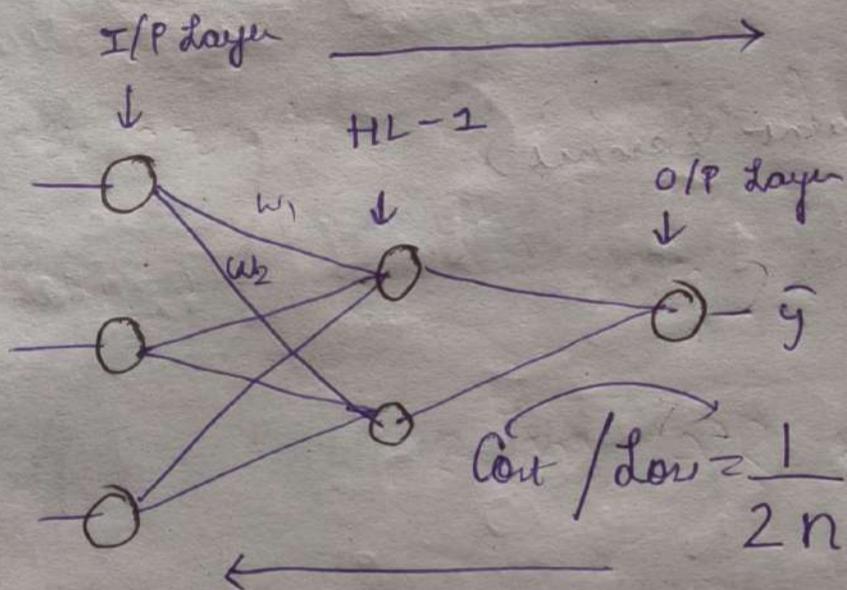
$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dH}{dw_{\text{old}}}$$

Learning rate

Loss/cost



global Minima



$$\text{Cost / Loss} = \frac{1}{2n} \sum_{i=1}^n (y - \tilde{y})^2$$

MSE

Epoch \rightarrow } 1 epoch (1 forward propagation, 1 backward propagation)

Dissadvantages

- 17 Resource intensive { Huge RAM }
{ as million record in 1 epoch }

② Stochastic Gradient Descent

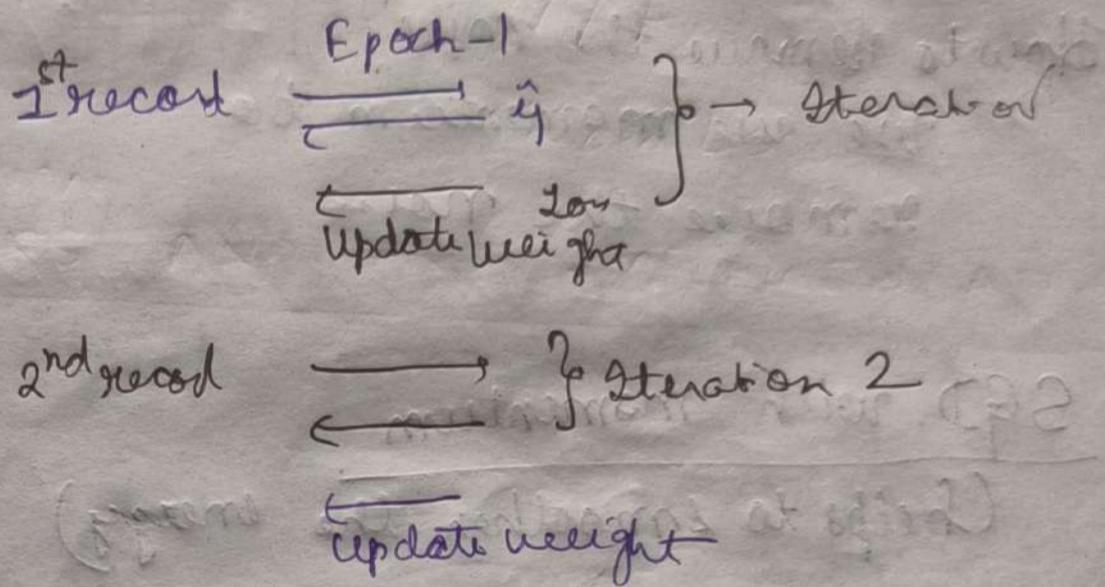
{ we pass 1 record in every epoch }

- ① DRAM ↓↓
but

② Disadvantage

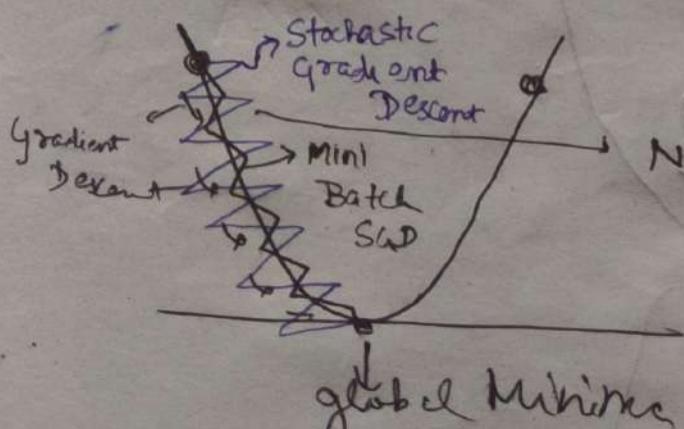
① Convergence will {
be very slow }

② Time Complexity
will also be
high.



③ ~~Batch~~ Mini Batch SGD

- ① Resource intensive
 - ② Convergence will be better
 - ③ Time complexity will improve



Let 100000 data points batch size = 1000
Epoch = 1 } Iteration 1

Epoch = 1
1000 } Iteration 1

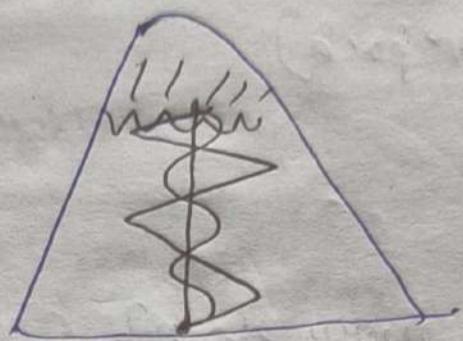
$\left. \begin{array}{l} \text{Iteration} \\ 1600 \end{array} \right\}$

Noise

```

graph TD
    Noise[Noise] --> GradientDescent[Gradient Descent]
    Noise --> MiniBatchSGD[Mini Batch SGD]
    Noise --> FullSGD[Full SGD]
    GradientDescent --> Stochastic[Stochastic]
    GradientDescent --> Batch[Batch]
  
```

if we have noise, it will definitely take more time'



How to remove the noise?

We use momentum to
remove the noise

④ SGD with momentum

(Helps to smoothen the energy)

{Exponential Weighted Average}



Time Series



ARIMA, ARMA,

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dh}{dw_{\text{old}}}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{dh}{db_{\text{old}}}$$

$$w_t = w_{t-1} - \eta \frac{dh}{dw_{t-1}}$$

Exponential Weighted Average

	t_2	t_3	t_4	\dots	t_n <small>← time</small>	a_1	a_2	a_3	a_4	\dots	a_n <small>← value</small>	{Forecasting}

$$V_{t_1} = a_1$$

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

hyperparameters

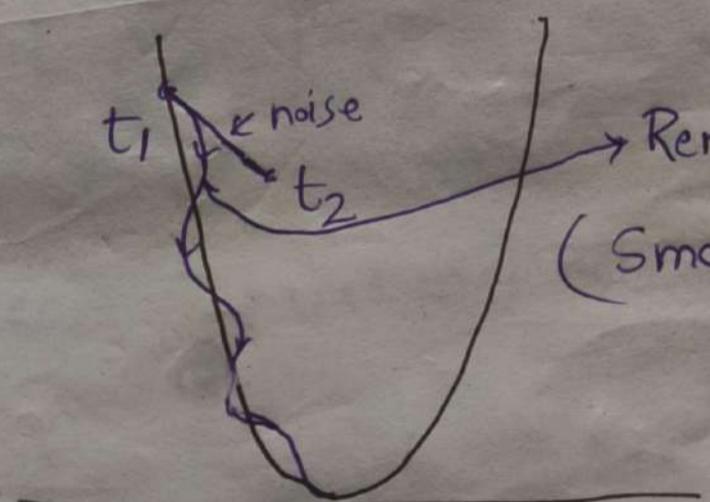
if $\beta = 0 \text{ to } 1$

if $\beta = 0.95$

$$V_{t_2} = (0.95) V_{t_1} + (0.05) a_2$$

Less imp. to a_2

a_2
 a_1 ↓
 Smoothening
 \downarrow will happen
 more control
 to a_1
 (previous
 point)
 than
 current
 point



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

Where to apply this?
We apply this in the derivative of loss to derivative of weight.

exponential Weighted average

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

$$V_{dw_t} = \beta \times V_{dw_{t-1}} + (1-\beta) \times \frac{dL}{dw_t}$$

$V_d \rightarrow$ value of derivative

Smoothing
of curve will happen

- { Reduce the noise
- Minibatch,
- Quicker convergence.

Recap

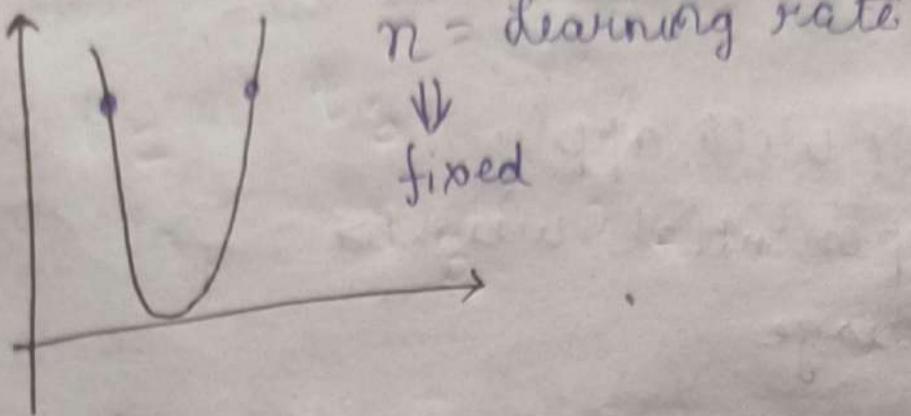
① Gradient Descent

② SGD

③ Minibatch SGD

④ SGD with momentum.

⑤ Adagrad → Adaptive Gradient Descent



$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}}$$

$\eta = \text{fixed} \rightarrow \text{adaptive}$

$$w_t = w_{t-1} - \eta' \frac{dh}{dw_{t-1}}$$

$$\eta' = \frac{\eta \leftarrow \text{Learning rate}}{\sqrt{\alpha_t + \epsilon}}$$

$\alpha_t + \epsilon$ small number
so denominator is $\neq 0$

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

η should decrease as we are reaching global minima.

$$t=1 \quad t=2 \quad t=3$$

$$\eta=0.01 \quad \eta=0.005 \quad \eta=0.002$$

\therefore Learning rate decreases

Since, $d_t + \epsilon$ is large therefore $n' \approx n^{20}$, therefore
very less change in update. So, to avoid this we use:
Adadelta and RMS Prop.

$$n' = \frac{n}{\sqrt{S_{dw} + \epsilon}}$$

Slow decreasing
of n' .

We will apply
exponential weighted
average

$$\begin{aligned} S_{dw} &= 0 \\ S_{dw_t} &= \beta S_{dw_{t-1}} + (1-\beta) \left(\frac{dh}{dw_t} \right)^2 \end{aligned}$$

$$\text{if } \beta = 0.95$$

$$S_{dw_t} = (0.95) S_{dw_{t-1}} + (0.05) \left(\frac{dh}{dw_{t-1}} \right)^2$$

Adam Optimiser (Best Optimizer)

Momentum + RMS Prop

(Adaptive Learning Rate)

Initially

$$V_{dw} = 0 \quad V_{db} = 0$$

$$S_{dw} = 0 \quad S_{db} = 0$$

$$\boxed{\begin{aligned} w_t &= w_{t-1} - n' V_{dw} \\ b_t &= b_{t-1} - n' V_{db} \end{aligned}}$$

$$n' = \frac{n}{\sqrt{S_{dw} + \epsilon}}$$

$$V_{dw} = \beta \times V_{dw_{t-1}} + (1-\beta) \frac{dh}{dw_t}$$

$$V_{db} = \beta \times V_{db_{t-1}} + (1-\beta) \frac{dh}{db_t}$$

- Helps in :-
- ① Smoothening
 - ② Learning Rate becomes adaptive }
③ Adapt

Practical ANN implementation

Day-4 - Deep Learning

① ANN Practical Implementation

② Early Stopping

③ Black Box Model Vs White Box Model

④ CNN Introduction

GOOGLE COLAB

ANN IMPLEMENTATION

#

```
> !pip install tensorflow-gpu
```

```
import tensorflow as tf
```

```
print(tf.__version__)
```

↳ 2.8.0 (should be > 2.0 for Keras integration)

```
## import some basic libraries
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import pandas as pd
```

```
dataset = pd.read_csv('Ourn - Modelling.csv')
```

```
dataset.head()
```

divide the dataset into independent and dependent features

all Row Column
 ↓ ↓

```
X = dataset.iloc[:, :-1]
```

```
y = dataset.iloc[:, -1]
```

Feature Engineering

```
geography = pd.get_dummies(X['Geography'], drop_first = True)
```

↳ One hot encoding

```
gender = pd.get_dummies(X['Gender'], drop_first = True)
```

concatenate these variables with dataframe

```
X = X.drop(['Geography', 'Gender'], axis = 1)
```

```
X = pd.concat([X, geography, gender], axis = 1)
```

Splitting the dataset into Training set and Test Set

```
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split  
    (X, Y, test_size  
     = 0.2)
```

Interview:- ① For which all algorithm feature scaling is required?

ANN? ✓
LR? ✓

k-Means! ✓

Logistic Regression? ✓

Feature scaling is required for:-

Decision? ✗

① Decision Distance based algo ✓

RF? ✗

② Gradient Descent algo ✓

XG? ✗

Do to make values smaller,
do to make calculation faster
or convg faster}

Why fit_transform is used for only train data and not test data?

→ to avoid data leakage

rotate around mean value

↓ ↗ Z-score

feature Scaling

```
from sklearn.preprocessing import StandardScaler
```

SC = StandardScaler()

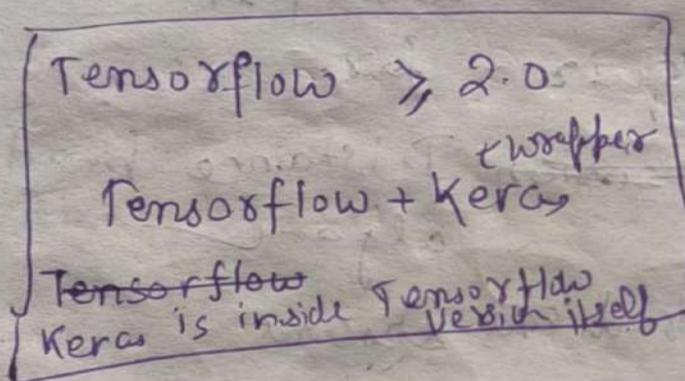
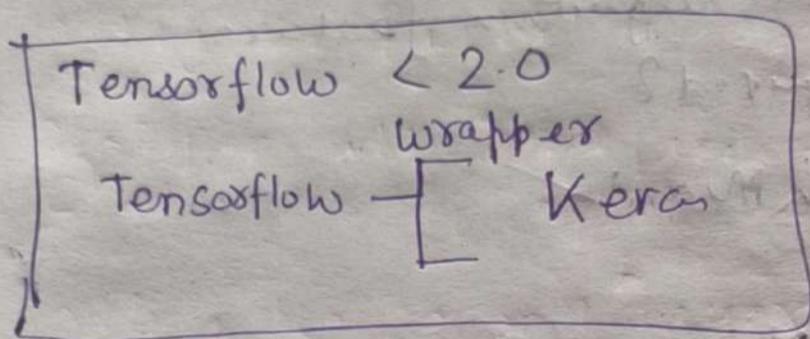
X_train = SC.fit_transform(X_train)

X_test = SC.transform(X_test)

Part 2: Now let's create the ANN

Tensorflow \Rightarrow Google Deepmind

Pytorch \Rightarrow Facebook



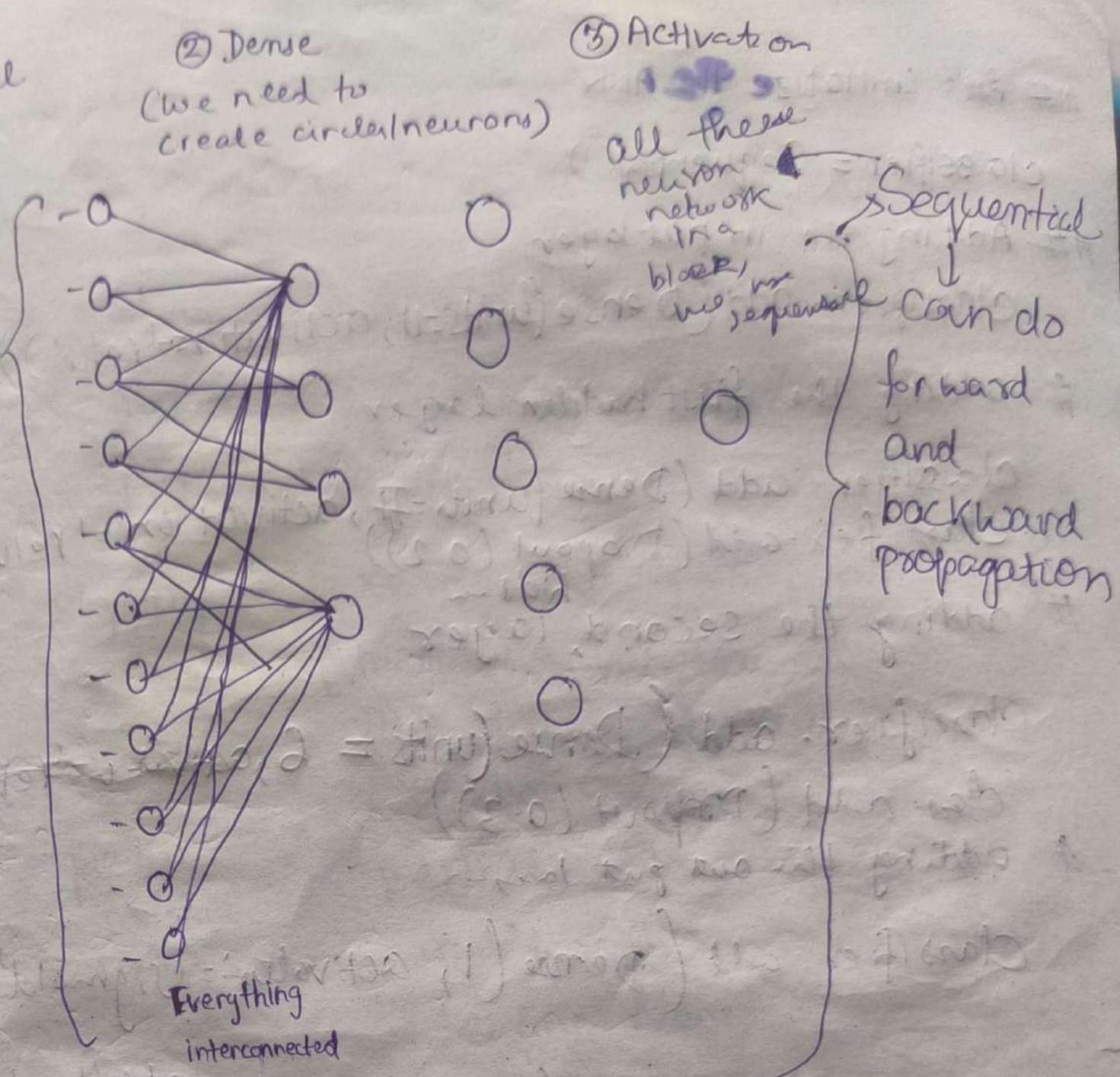
```
from tensorflow.keras.models import Sequential
```

```
from tensorflow.keras.layers import Dense
```

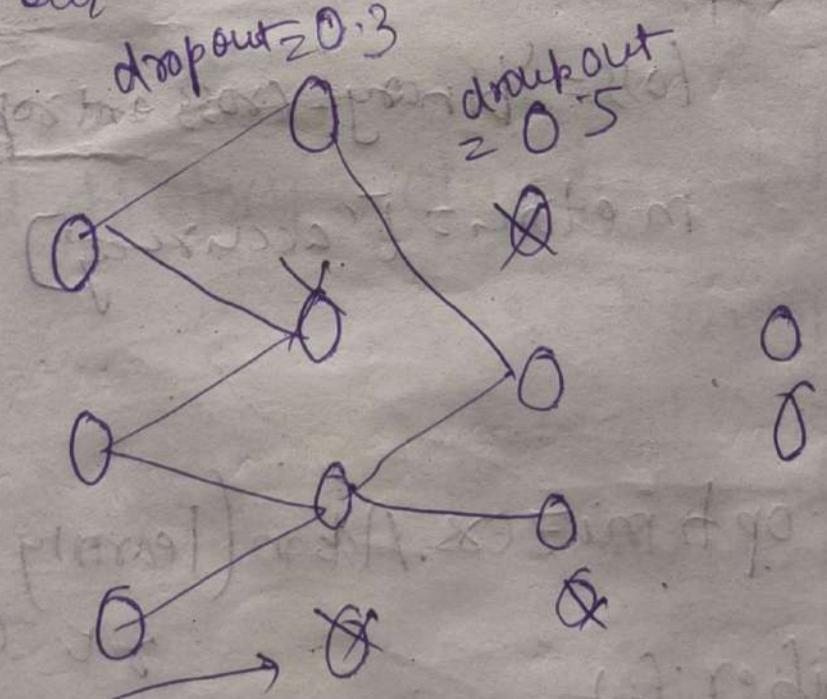
```
from tensorflow.keras.layers import LeakyReLU, PReLU, ELU, ReLU
```

```
from tensorflow.keras.layers import Dropout
```

① Sequential



④ Drop out



This entire neural network sometimes lead to overfitting

To remove overfitting we use, Drop out layer, \rightarrow Regularization

if drop out layer is 0.3, then 30% of neuron will get select rated -

{ L1 Norm }
L2 Norm }

Let's initialize the ANN

classifier = Sequential()

Adding the input layer

classifier.add(Dense(units=11, activation='relu'))

Adding the first hidden layer

classifier.add(Dense(units=7, activation='relu'))
classifier.add(Dropout(0.2))

Adding the second hidden layer

classifier.add(Dense(units=6, activation='relu'))
classifier.add(Dropout(0.3))

Adding the output layer

classifier.add(Dense(1, activation='sigmoid'))

classifier.compile(optimizer='adam',
loss='binary_crossentropy',
metrics=['accuracy'])

import tensorflow

opt = tensorflow.keras.optimizers.Adam(learning_rate=0.01)
model.history = classifier.fit(x_train, y_train)

(x_train, y_train)

Validation_split=0.33,

batch_size=10, epochs=100)

Model

42/1000

Early Stopping

import tensorflow as tf

early_stopping = tf.keras.callbacks.EarlyStopping

(monitor='val-loss')

min_delta=0.0001,

Patience=20,

verbose=1

mode="auto"

baseline=None

restore_best_weights=False

)

model_history = classifier.fit(x_train, y_train)

validation_split=0.33, batch

size=10, epoch=1000,

callbacks=early_stopping)

```
# model.history.history @ keys()
dict_keys(['loss', 'accuracy', 'val_loss', 'val_accuracy'])

# Summarise history for accuracy or val-loss also
# plt.plot(model.history.history['accuracy']) mean plot
plt.plot(model.history.history['val_accuracy'])

plt.title('model-accuracy')
plt.ylabel('accuracy')
plt.xlabel('epoch')
plt.legend(['train', 'test'], loc='upper left')
plt.show()
```

```
y_pred = classifier.predict(X_test)
y_pred = (y_pred >= 0.5)
```

```
# # make confusion matrix
```

```
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
```

```
# # Calculate the accuracy
```

```
from sklearn.metrics import accuracy_score
score = accuracy_score(y_pred, y_test)
```

get the weights
classifier.get_weights()

Black Box Model Vs White Box Model

Random Forest :- Black box

DecisionTree → WhiteBox Model

Linear Regression → White Box Model

CNN/ANN → Black Box Model

xgboost → Black Box Model

Explainable AI { }

Understanding CNN and Implementation

ANN → Theoretical, Practical

CNN → Convolution Neural Network

Image, Video Frame

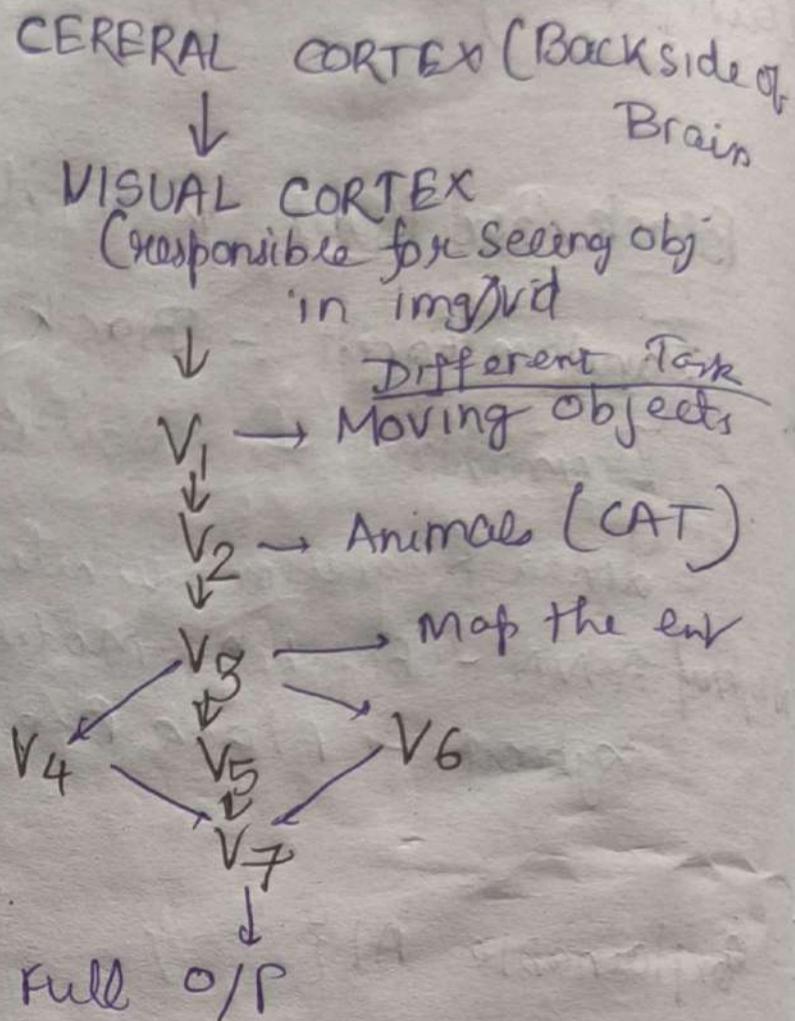
Agenda

- ① CNN vs Human Brain
- ② Convolution operation
 - Convolution
 - Padding
 - stride
 - Filter (Kernel)
- ③ Max Pooling
- ④ Flattening
- ⑤ Practical Implementation CNN

CNN vs Human Brain

CAT img
DOG CAR
BICYCLE
HUMAN

Layers



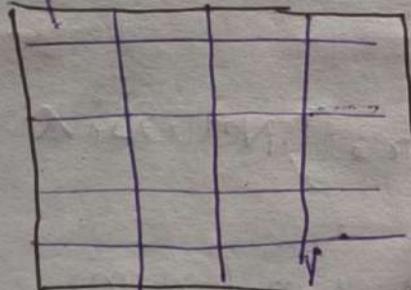
CNN

Convolution

Images =

Black and white images

0 - 255



Black and white OR RGB

1 channel

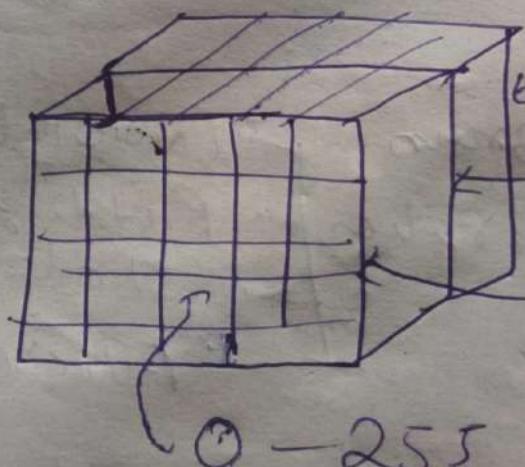
0 → Black

255 → White Color

5x5 pixel img

RGB image

5x5x3



Blue channel

Green channel

Red Channel

0 - 255

wrt each and every channel

1st step of CNN: Convolution

0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255

6x6

3x3 filter or Kernel

1	2	1
0	0	0
-1	-2	-1

exp

{ Horizontal }
{ Edge filter }

4x4 output

Step-1 We bring the values ranging from 0-255 to 0 to 1.

This process is known as min-max scaling.

Step-2 We place the filter on top of input layer.

Stride - How many posn to skip

ie Stride = 1

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

let's consider for first op

0	0	0
0	0	0
0	0	0

1	2	1
0	0	0
-1	-2	-1

$$1(0) + 2(0) + 1(0) + 0(0) + 0(0) + (-1)(0) \\ + (0)(-1) + (0)(-2) + (-1)(0)$$

So output's first cell = 0

0			

Similarly, we continue, and the first row of output becomes:-

0	0	0	0

Since, stride = 1 we move to right by 1.
 And after completion of 1st row in the output layer,
 we move down by 1 unit.

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

$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix} =$$

0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

$$\begin{aligned}
 & 1(0) + 2(0) + (1)(0) + 0(0) + 0(0) + 0(0) \\
 & + (-1)(0) + (-2)(0) + (-1)(0) \\
 & = 0
 \end{aligned}$$

In next step

$$\begin{aligned}
 & 1(0) + 2(0) + 1(1) + 0(0) + 0(0) + 1(0) \\
 & - 1(0) + (-2)(0) - 1(1) = 0
 \end{aligned}$$

H-W → Calculate the above by
 doing convolution process.

These operations are specifically known
 as "Convolution".

Let's say we get some output as

Imagine

0	0	0	0
-4	4	-4	-4
-4	4	-4	4
0	0	0	0

Q-

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

We had done our map scaling, 8×6

Already defined

*
Convolution

1	0	-1
2	0	-2
1	0	-1

→ helps to identify the vertical edge.

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

Vertical edge filter detected or kernel, 3×3

4×4

Output

4×4

0	-4	-4	0
0	-4	4	0
0	-4	-4	0
0	-4	-4	0

Reversing
back from feature

255	0	0	255
255	0	0	255
255	0	0	255
255	0	0	255

scaling

$255 \rightarrow$ white color
 $0 \rightarrow$ black

white black actual image



vertical edge

- * These filters like horizontal edge filter, vertical edge filter, round edge filter etc. helps to identify or extract various info. So, we can have horizontal, vertical or round edge and many more filters.
- * We will get multiple outputs as we use multiple different filters for the same image.

Whenever we pass a (6×6) input image and pass it through (3×3) filter we get 4×4 output. So, how is this happening??

if 'n' is the size of the image.

and

if 'f' is the size of the filter

then we get

$n-f+1$ as output size.

$$\text{if } n=6 \quad f=3$$

$$\therefore n-f+1 = 6-3+1 = \underline{\underline{4}}$$

Here the image size decrease from 6×6 to 4×4 , which should not happen, (we are losing some info therefore we use 'Padding').

Padding is nothing but a layer ^{on top of image} which protects the image.

6×6

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

We apply
padding

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

Here padding = 1
(as padding layer = 1)

Padding
layer

~~Image is 6×6 only~~
~~but (8×8) in total~~

Padding are of different types according to the values p filled in the cells of padding.

① Zero Padding

↳ fill it with 0.

② fill it with the nearest value.

So after padding

$$8 - 3 + 1 = \underline{\underline{6}}$$

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

$$\begin{matrix} * & \begin{matrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{matrix} \end{matrix} =$$

25	0	0	25
25	0	0	25
25	0	0	25
25	0	0	25

So, after padding:

Formula for output image size is $n + 2p - f + 1$ \rightarrow stride = 1

$$n + 2p - f + 1 \quad \text{for } s=1$$

So $n = 6$ but general formula \rightarrow
 $p = 1$
 $f = 3$

$$\frac{n + 2p - f + 1}{s} \quad s \rightarrow \text{stride}$$

$$\therefore 6 + 2(1) - 3 + 1 = 8 - 3 + 1 = 6$$

What is the importance of padding??

- To prevent the info. loss of image we apply different types of padding techniques

* We don't need to hard code the values in filter/kernel. First we randomly choose any values, and then through back propagation, we update these values based on ^{input} image.

So, where how backpropagation will happen, we will have not used any ~~an~~ activation function. So where do we specifically use activation function??

Ans:- After applying convolution ~~ad~~ operation, and getting the output we apply ReLU activation fn on each and ~~as~~ every cell.

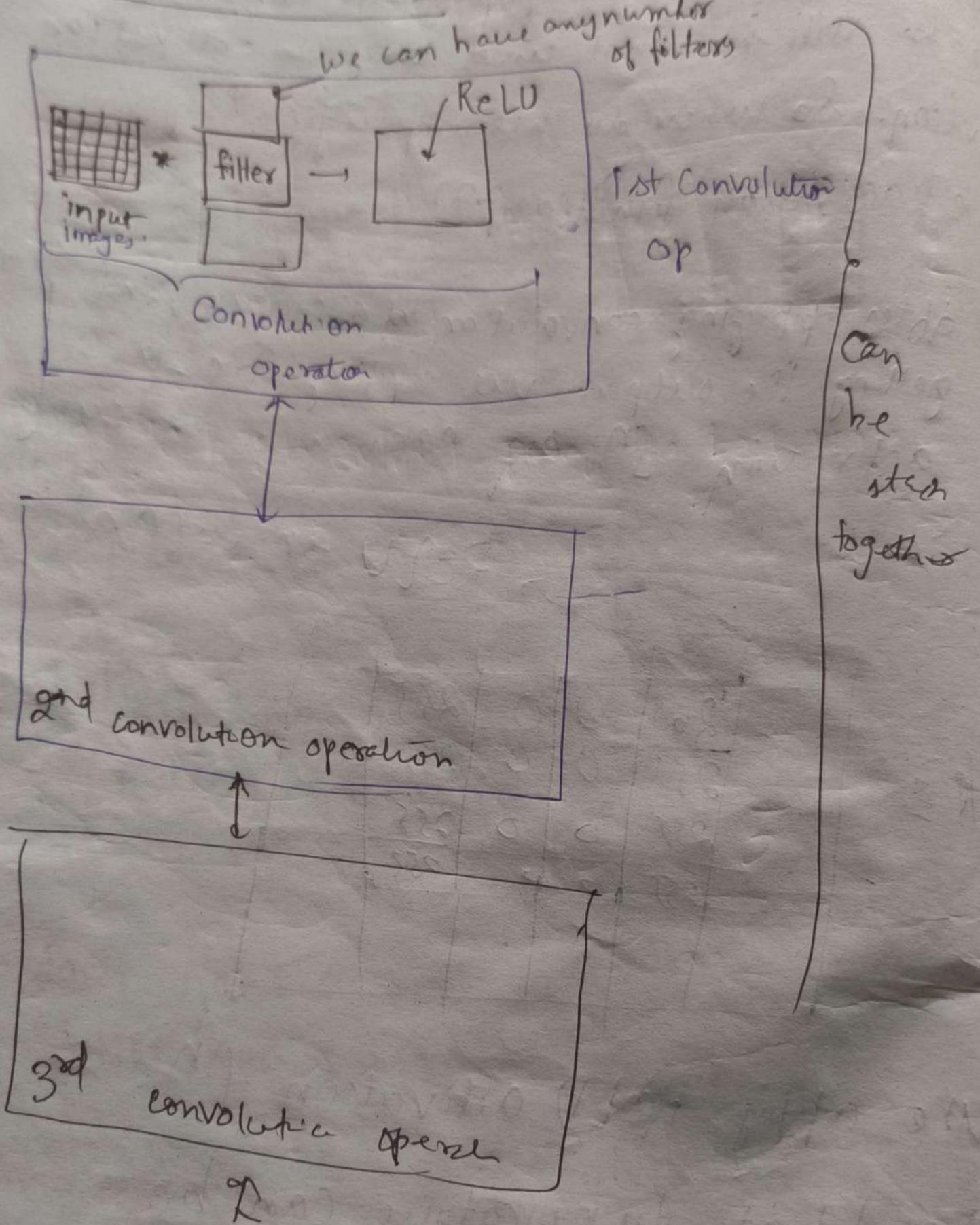
$$\text{ReLU} = \max(0, n)$$

O/p	255	0	0	255
O/p	255	0	0	255
O/p	255	0	0	255
O/p	255	0	0	255

We apply ReLU activation fn

so that derivative can be found out during back propagation which can help us to update filter values.

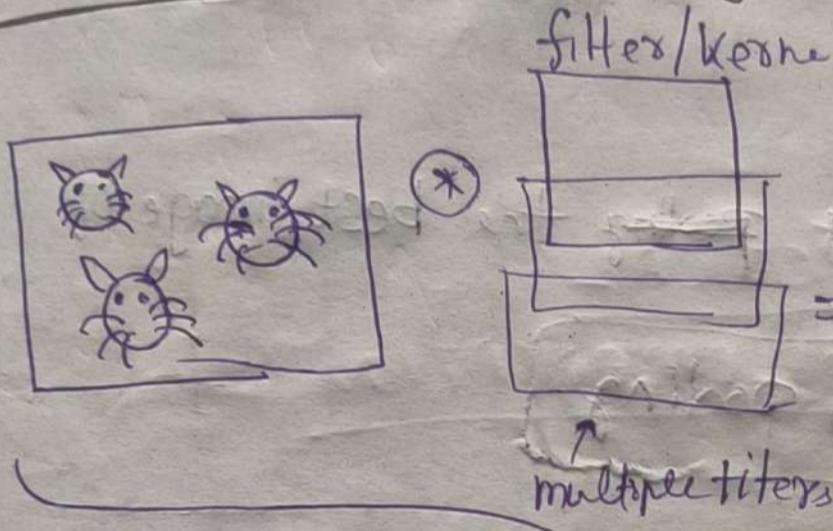
So, Steps of CNN



- * We need to learn from these filters
- * We need to learn and update ~~an~~
filters based on the
input images.

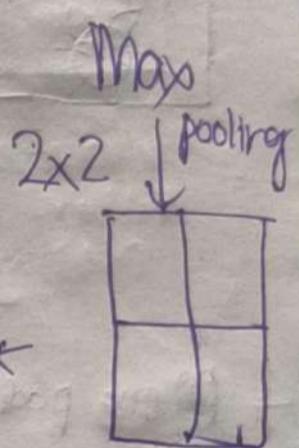
Train
Aim

MAX-POOLING (Max Pooling layer)



{Location invariant}

ReLU



CONVOLUTION

~~EX~~ OI ERATION

Location invariant says that my CNN till it goes to the next layer should be able to extract more clear and clearer information, *location of thing can be anywhere*. Therefore, to extract much clearer information, we use max pooling.

Max Pooling layer

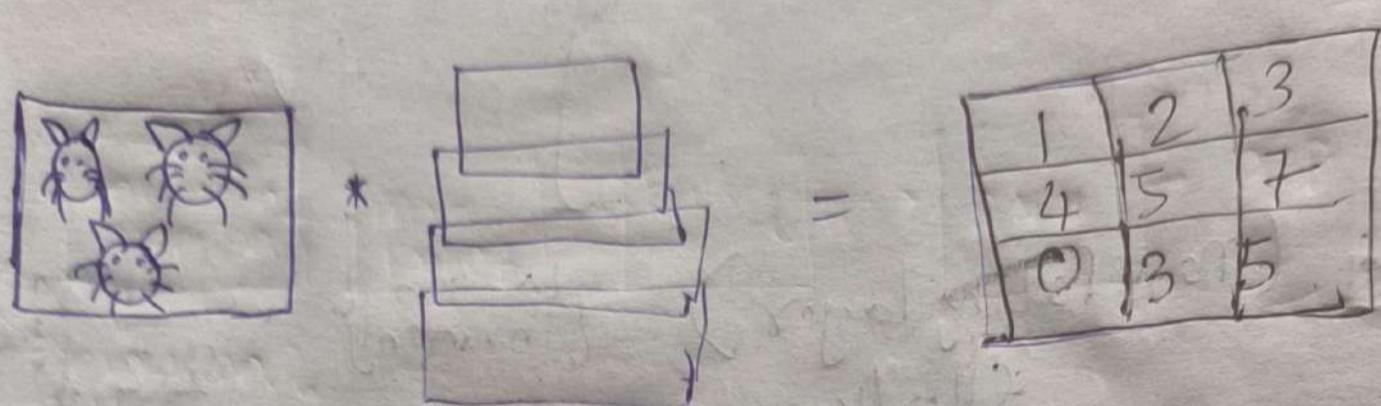
Avg ~~Pooling~~

Pooling

Min
Pooling

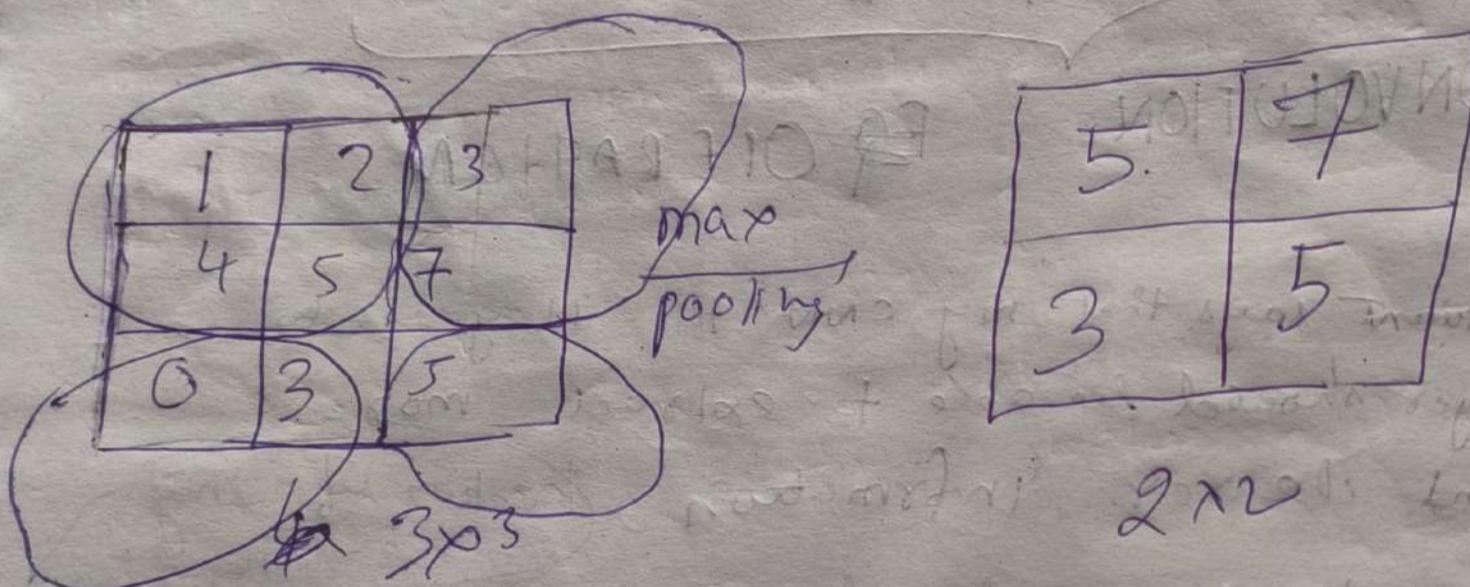
Max Pooling

* In pooling $\boxed{\text{stride} = 2}$ $\star\star\star$



filter layer

- Max pooling helps to extract ~~pooling~~ the best image of the object.
- $\boxed{\text{Stride} = 2 \text{ in max pooling}}$



avg
pooling

6	5
1.5	2.5

min pooling

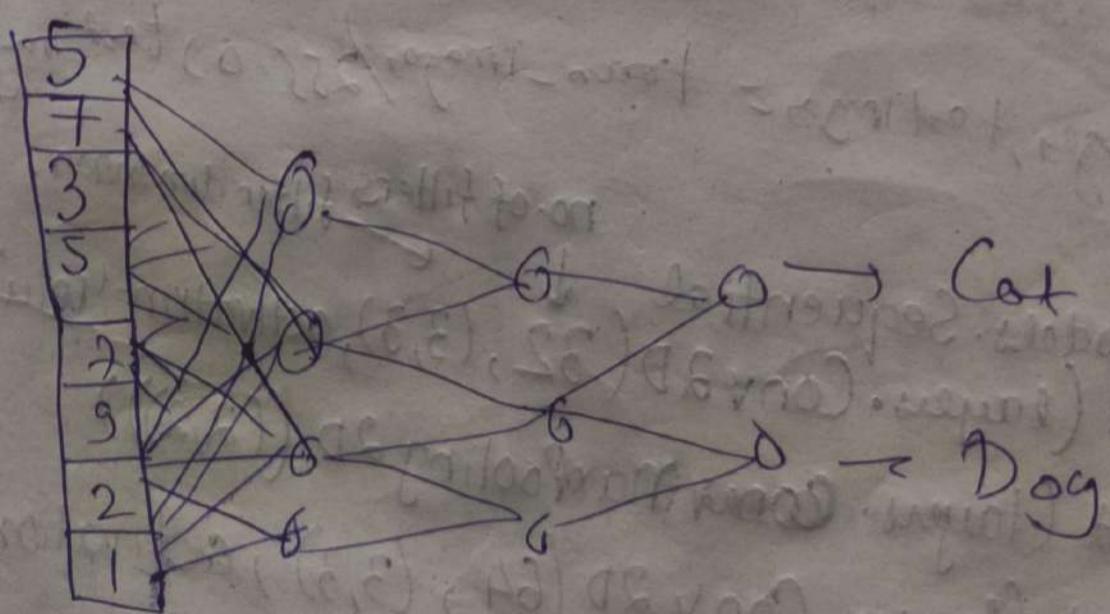
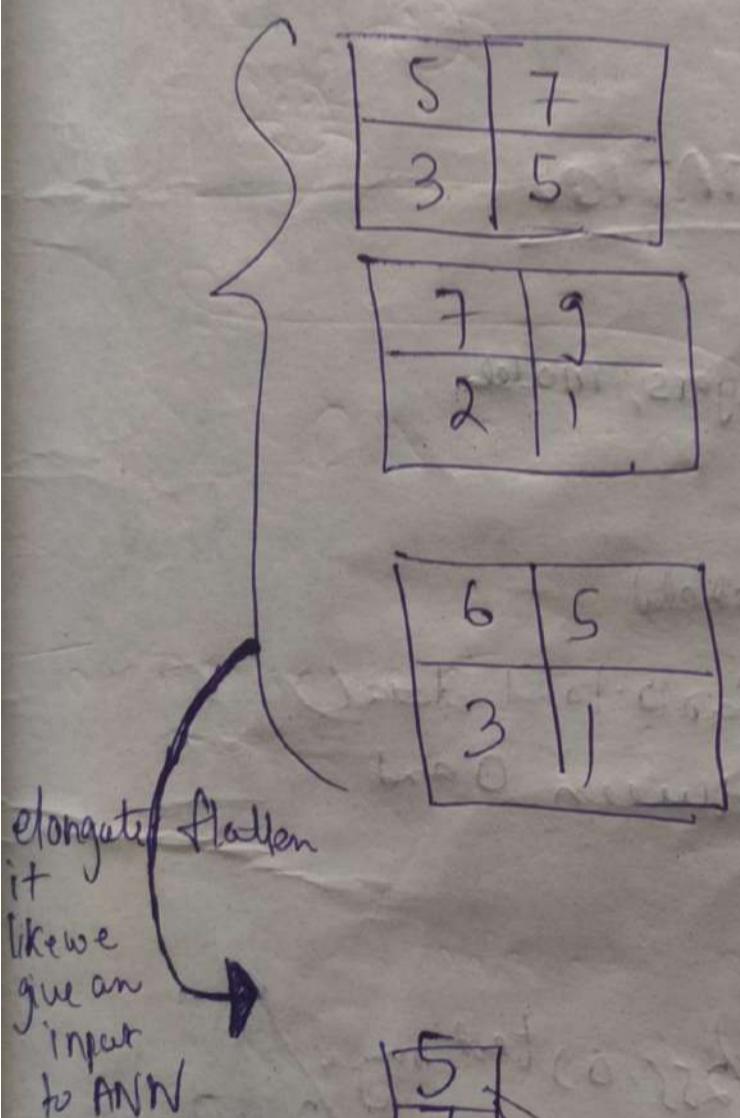
1	3
0	5

- * According to different problem statement we can use max pooling, min pooling and average pooling.
- * Max pooling solves the problem of local invariant.
- * With respect to different filter operation we get different max pooling value.
- We can combine CNN and maxpooling.

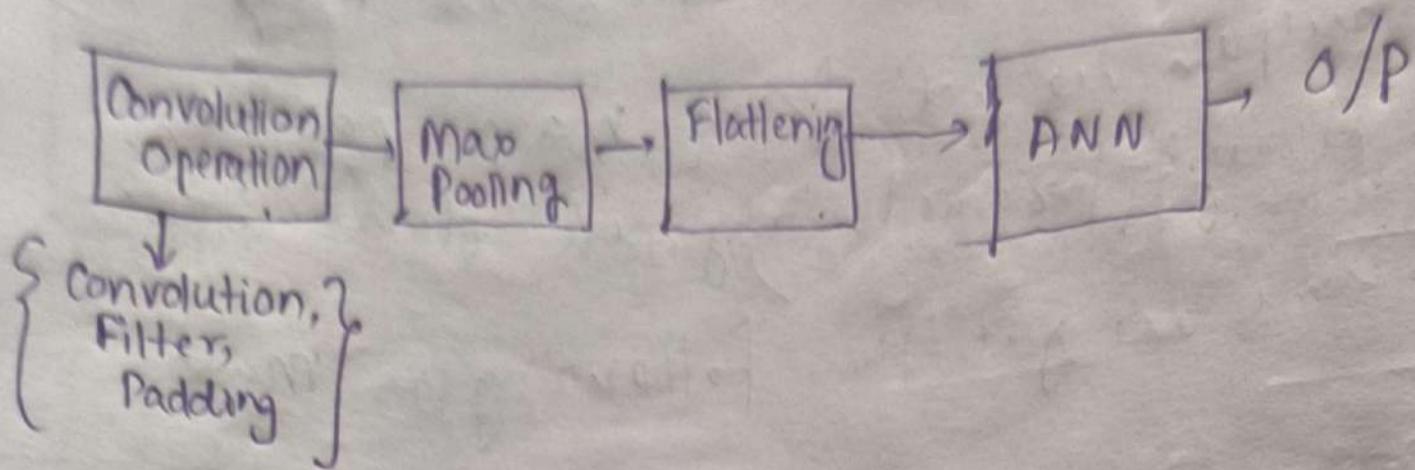
④ Flattening layer (looks like an ANN Dense Layer)

Whatever outputs we get from max-pooling we will straightforwardly just flatten it.

Let us assume we get the following as max pooling:-



Summary:



CNN PRACTICAL IMPLEMENTATION

```
- import tensorflow as tf  
import  
from tensorflow.keras import datasets, layers, models  
import matplotlib.pyplot as plt  
(train_images, train_labels), (test_images, test_labels)  
= datasets.cifar10.load_data()  
# Normalize pixel values to be between 0 and  
1  
train_images, test_images = train_images/255.0, test_images/255.0  
#  
model = models.Sequential  
model.add(layers.Conv2D(32, (3,3), activation='relu', input_shape  
= (32, 32, 3))  
model.add(layers.MaxPooling2D((2,2)))  
model.add(layers.Conv2D(64, (3,3), activation='relu'))  
model.add(layers.MaxPooling2D((2,2)))  
model.add(layers.Conv2D(64, (3,3), activation='relu'))
```

```
model.summary()
```

```
model.add(layers.Flatten())
```

```
model.add(layers.Dense(64, activation='relu'))
```

```
model.add(layers.Dense(10)) ← output layer
```

```
model.summary()
```

```
model.compile(optimizer='adam',
```

```
loss=tf.keras.losses.SparseCategorical  
Crossentropy(from_logits=True)  
metrics=['accuracy'])
```

```
history = model.fit(train_images, train_labels,  
epochs=10, validation_  
data=(test_images),  
validation_labels))
```

```
# Visualizing:
```

```
plt.plot(history.history['accuracy'], label='accuracy')
```

```
plt.plot(history.history['val_accuracy'], label='val-accuracy')
```

```
plt.xlabel('epoch')
```

```
plt.ylabel('Accuracy')
```

```
plt.ylim([0.5, 1])
```

```
plt.legend(loc='lower right')
```

~~iter~~ - loss , test_acc = model.evaluate(test_images,
test_labels
(~~subset~~ = validation, verbose = 2))
verbore = 2)

print (test_acc))