

Ques Overcast & Mild (\rightarrow play or not?)

$$P(\text{yes}) + P(\text{overcast/yes}) + P(\text{mild/no})$$

$$\frac{9}{14}$$

$$P(\text{no}) \cdot P(\text{mild/yes}) \cdot P(\text{mild/no})$$

$$\frac{5}{14}$$

100% will play ✓

{ Data Preprocessing }

In this step the data analyses pipeline

- It involves transforming raw data into a clean & usable format.

- Primary step in this process include -

→ Data Cleaning (Handling missing values, removing duplicates and correcting inconsistencies)

→ Data Integration (combining data from different sources)

→ Data Transformation (normalizing, scaling and transforming data for better analysis)

→ Data Reduction (Reducing the number of attributes or dimensions)

→ Data Rescaling (It involves adjusting the range of data featured to a standard range)

Min-Max Scaling : Rescale the data to fixed range 0 to 1.

$$X_{\text{scaled}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

→ Data Standardization
standardizing transforms the data to have a mean of zero and standard deviation of one.

$$X_{\text{stand}} = \frac{x - \mu}{\sigma}$$

$\mu \rightarrow$ mean
 $\sigma \rightarrow$ standard deviation

→ Normalizing Data : It is scaling the data such that it has a unit known (norm), typically used in algorithms that involves distance like KNN. The most common form is the L_2 normalization.

$$X_{\text{Norm}} = \frac{x}{\|x\|} \quad \text{Euclidean Norms.}$$

Euclidean Distance → It is measured by the straight line distance between two points in Euclidean space for two points $A(x_1, y_1)$ and $B(x_2, y_2)$

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

→ Manhattan distance : It is also known as L_1 distance. It is the sum of the absolute difference of the coordinates for two points $A(x_1, y_1)$ & $B(x_2, y_2)$

$$d(A, B) = |x_2 - x_1| + |y_2 - y_1|$$

→ Handling categorical Data : It needs to transform a numerical format for most machine learning algorithms

One hot Encoder : It converts categorical values into binary matrix. for example "Red", "Green", "Blue"

would be transformed into three columns.

	Red	Green	Blue
1	1	0	0
2	0	1	0
3	0	0	1

Label Encoding \rightarrow It converts into integer (binary values)

Red: 0

Green: 1

Blue: 2

takes a lot of time

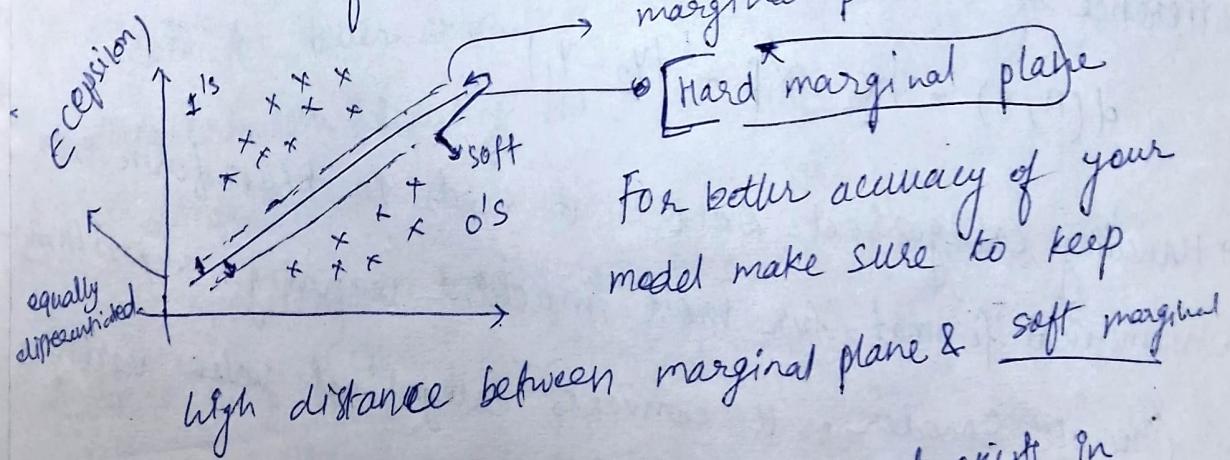
Typically used

F1 score is also called $f-1$ ~~beta~~ score?

- \rightarrow if given option we classify (cls. species)
- \rightarrow if to tell a numeric value b/w a range it's regression

Support Vector Machine

Tasks \rightarrow classification & regression



We don't work on Hard marginal plane - (not exists in real life)

Support Vector Classification

$x \rightarrow$ magnitude or $\sqrt{(\text{direction})}$

$$y = mx + c \quad y = \theta_0 + \theta_1 x \quad y = \beta_0 + \beta_1 x$$

$$y = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

↑ transpose $w^T x$ slope

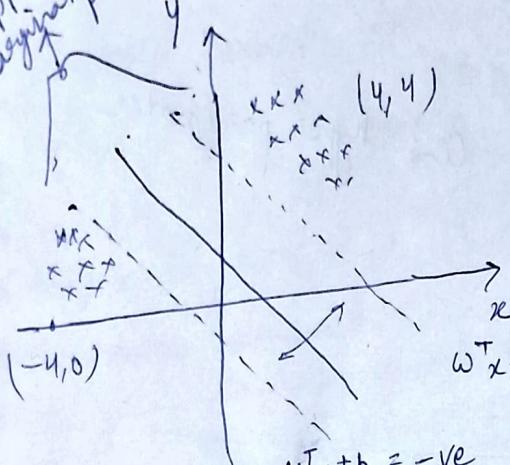
straight line

$$ax + by + c = 0$$

$$y = -\frac{a}{b}x + \frac{c}{b}$$

$$mx + c$$

Support
marginal plane



for ex $\rightarrow 3x + 2y + 4 = 0$
 $-\frac{a}{b}x = -8 \quad \frac{c}{b} = 24$

(w) \rightarrow magnitude of vector

data points will always lie b/w zeroes & ones

$$\frac{w^T x_1 + b}{\|w\|} = +1 \quad w^T x_2 + b = -1$$

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

diff b/w both planes
maximize
maximum $w, b = \frac{2}{\|w\|}$

marginal plane distance

for all alternate data points

$$y_i (w^T x_i + b) \geq 1$$

constraints

Constraints

$$y_i \begin{cases} +1 & w^T x_i + b \geq 1 \\ -1 & w^T x_i + b \leq -1 \end{cases}$$

miss classification (when the maximize value or marginal plane distance is high)

$$\frac{2}{\|w\|}$$

$$\frac{\|w\|}{2}$$

minimise

maximise

cost function (SVM)

$$\text{minimize}(w, b) = \frac{\|w\|}{2} + C \sum_{i=1}^n \xi_i$$

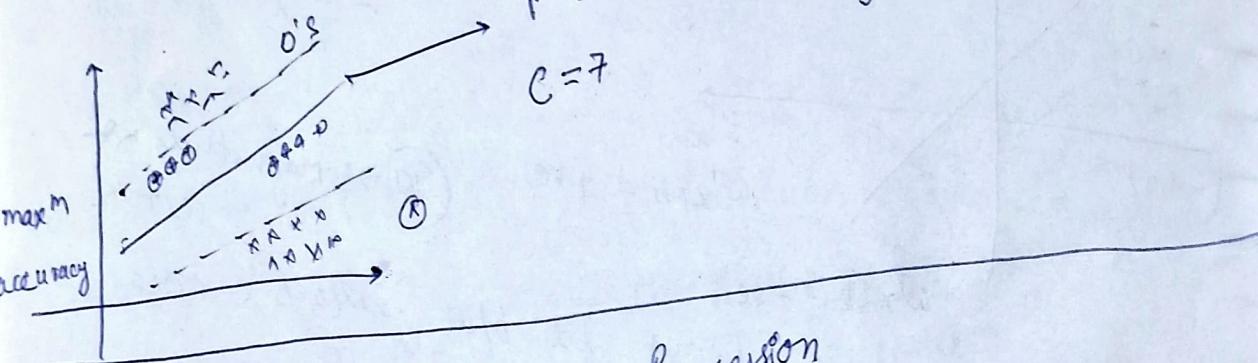
eta

summation of the distance
of miss classification points from marginal plane.

$\xi_i \rightarrow$ How many point we can avoid ~~through~~ miss classification

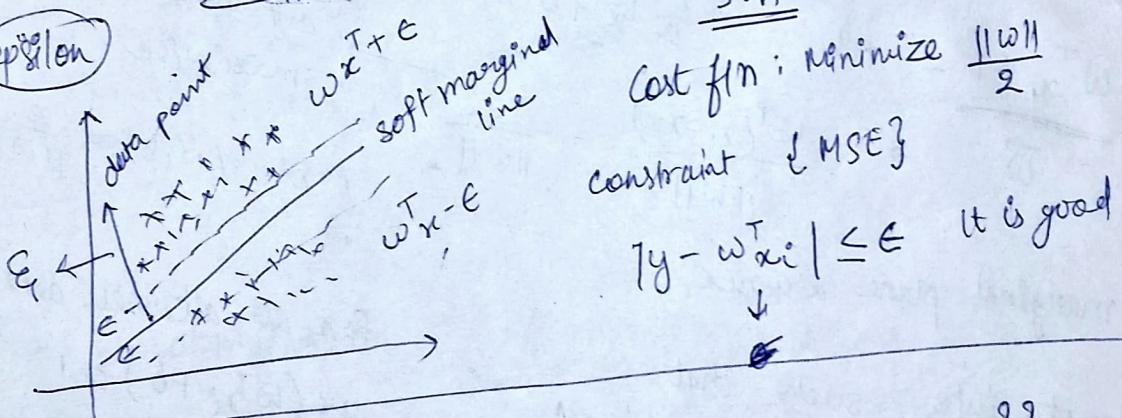
$w \rightarrow$ width of marginal plane

$C \rightarrow$ Hyper parameter



Support Vector Regression

ϵ (epsilon)



SVR

Cost fn: minimize $\frac{\|w\|}{2}$

constraint {MSE}

$$|y - w^T x_i| \leq \epsilon \quad \text{it is good}$$

Cost function

$$\frac{\|w\|}{2} + C \sum_{i=1}^n |\xi_i|$$

penalty
parameter

Hyper parameter

constraint

C in code

$$|y_i - w^T x_i| \leq \epsilon + \xi_i$$

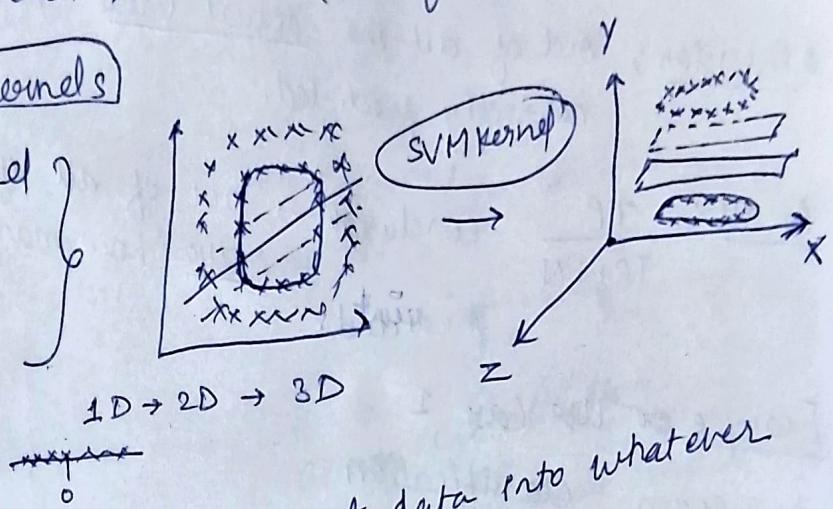
? range less ??

$\xi_i \rightarrow$ eta

Transformation through SVM kernel is done performed to convert the model into 3D. (for left out data points)

Types of SVM kernels

- ③ Polynomial Kernel
RBF Kernel
Sigmoid Kernel



- Polynomial Kernel → convert data into whatever dimension you need.

convert data into whatever

1D → 2D

Measure Parameters

Confusion Matrix

- Accuracy
- Precision
- Recall
- F1 Beta Score

		0	1
predicted	0	TP	FP
	1	FN	TN

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

Predicts for any new data point.

Precision

Binary classification

Datasets :
1000 data points

900 1's
100 0's

Imbalance datasets

After training the model
How good can it predict more

↓
Precision

$$\text{Precision} = \frac{TP}{TP + FP} \downarrow$$

decreasing FP will give you better precision.

How much is TP & FP
↓
based on this
How much the
model has
predicted.

Recall \rightarrow you have to decrease the FN values.

Precision: Out of all the actual value how many are correctly predicted.

Recall: $\frac{TP}{TP+FN}$ conclusion \rightarrow out of all the predicted value how many are correctly predicted.

Example or Use Case 1
spam classification
Email \rightarrow spam
Model \rightarrow spam
Good classification

Mail \rightarrow not spam
Model \rightarrow spam
Blunder

1	TP	FP
0	TN	FN

\rightarrow reduce it (what if any imp email is received & we are unaware of that)

Use Case 2: To predict whether person has cancer or not
Truth \rightarrow Yes (Cancerous)
Model \rightarrow No

Truth \rightarrow Cancer
Model \rightarrow Cancer

\boxed{FN} (Blunder is cuz of model machine)
Domain expert for medical field

Stock Market analysis \rightarrow Reduce both $\downarrow \downarrow$

F1 Beta Score: Combines both precision & recall
if FP & FN are both imp.

$$1 + \beta^2 \left(\frac{P \times R}{P + R} \right)$$

↳ Harmonic Mean

$$\beta = 1$$
$$F1 \text{ score} \Rightarrow 2 \times \frac{P \times R}{P + R}$$

If FP is more important FN
 $\beta = 0.5$ then $F_{0.5} \text{ score} = (1+0.25) \frac{PXR}{P+R}$

\rightarrow If $FN > FP$
 $\beta = 2$ then $F_2 \text{ score} = 1 + \frac{PXR}{P+R}$

Majorly we calculate F1 Beta score

28th Dec

Decision Tree

DTC — ID] Both
 CART
 (spits binary basis)

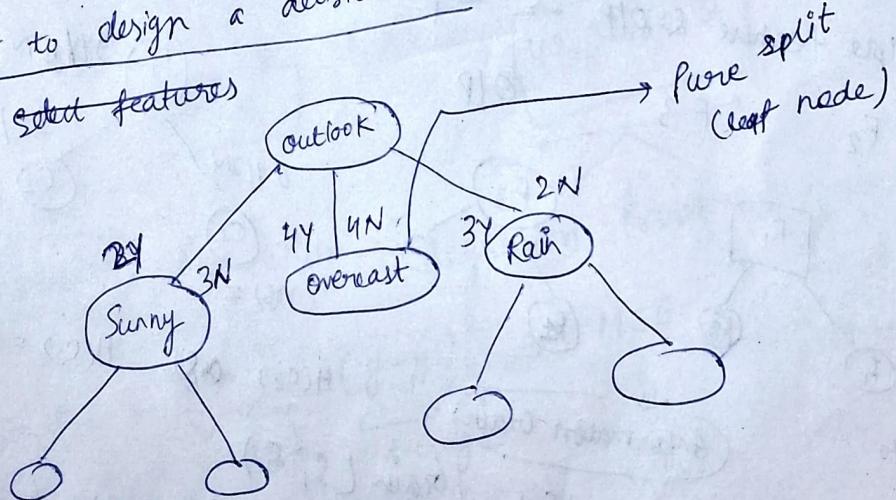
Entropy & Gini Split
 feature to select

Binary classification

- . we have to train our model
- . until ~~we~~ we get a single leaf node.
- . Where is my leaf node we'll check it through entropy count impurity

Steps to design a decision tree

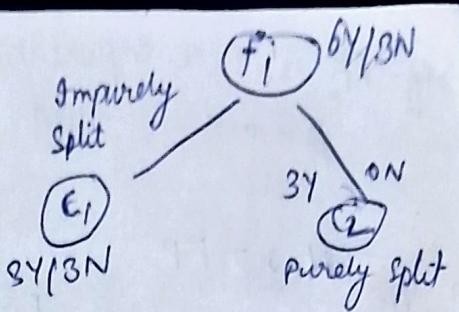
~~Select features~~



① Purity, Pure or Impure

- . Entropy
- . Gini Impurity

Step 2 what feature you need select
for splitting (Information gain)

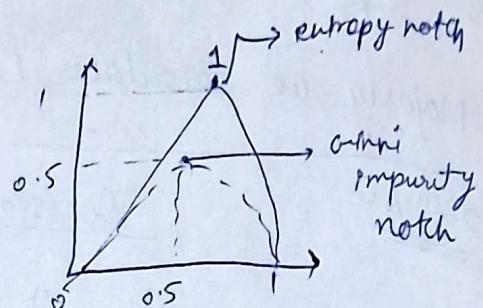


Entropy $(0, 1)$ best suitable

$$H(S) = -P \log_2 P + -(1-P) \log_2 (1-P)$$

$$= \frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$= 0$$

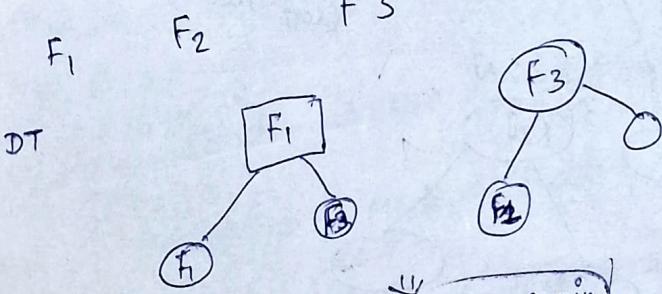


curvi impurity

$$\text{pure } 1 - \sum_{j=1}^m (P_j)^2 \rightarrow 1 - \left(\frac{3}{6}\right)^2 = 0$$

$$\text{impure } 1 - ((P_+)^2 + (P_-)^2) \rightarrow 1 - ((\frac{3}{6})^2 + (\frac{3}{6})^2) = 0.5$$

Decision Tree
multiple feature & OIP \rightarrow Supervised learning

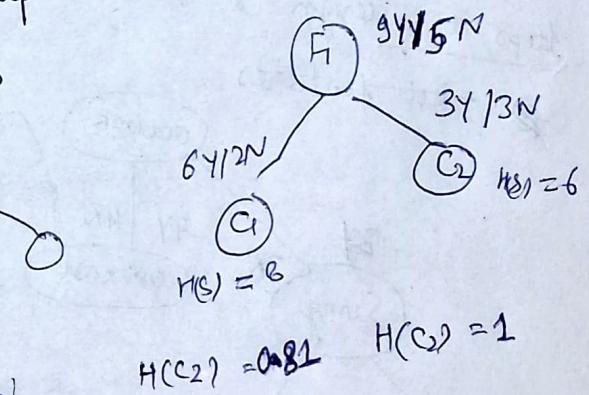


$c \rightarrow$ belongs to

$$\text{Gain}(S, F_2) = H(S) - \frac{\sum |S_i|}{\text{Vtvarious}} H(S_i)$$

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

choose this one



$$\begin{aligned} \text{Gain}(S, F_2) &= 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right] \\ &= 0.049 \end{aligned}$$

$$\text{Gain}(S, F_1) = 0.051$$

$$\text{Gain}(S, F_2) > \text{Gain}(S, F_1)$$

Small dataset

Fixed classification

1, 2, 3, 4

Entropy

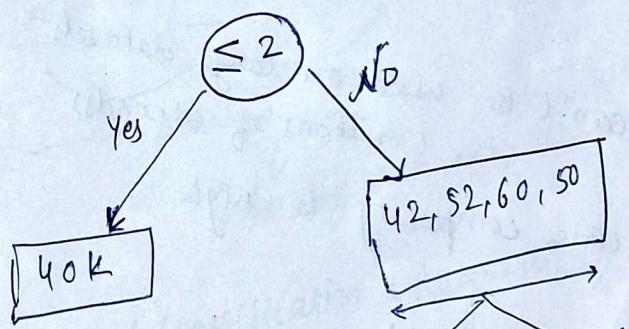
Big dataset

Gini Impurity

DTR (Decision Tree Regression)

	dataset	Salary
Exp	carrier group	
-2	Yes	40 K
2.5	Yes	42 K
3	No	52 K
4	No	60 K
4.5	Yes	56 K

$$C1 = 100 \\ C2 = 42$$



Variance

Reduction {Regression Problem}

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n \underbrace{(y - \hat{y})^2}_{\text{Average}} \Rightarrow \text{MSE}$$

50 → Predicted value

$$\text{Var}(\text{root}) = \frac{1}{5} \rightarrow 60.8$$

$$\text{Var}(\text{root}) = \frac{1}{5} \left\{ (40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2 \right\}$$

DT Parameters

$$\text{Var}(CII) = \frac{1}{2} ((40-50)^2 + (42-50)) = \frac{1}{2} (100+64) = 82$$

$$\text{Var}(C_1) = \frac{1}{3} [4 + 100 + 36] = 46.66$$

child 1 (c1)
Classification

Variance Reduction

$$60.8 : \left(\frac{2}{5} \times 82 + \frac{3}{5} \times 46.66 \right) = 0.04$$

$$60.8 - \left(\frac{1}{3} \times 100 + \frac{4}{5} \times 42 \right)$$

$$60.8 - [20 +] = 0$$

if we have to make multiple choices through regression
 0 ≤ 2.5 $0.04 \rightarrow$ variance
 Left side Right side High
 targets

Disadvantages → can't be used on large datasets
(millions of records)
→ time complexity is high

Ensemble Modelling

Types → 1) Bagging

2) Boosting

(L & R)

(L & R)
Random forest

Ada Boost

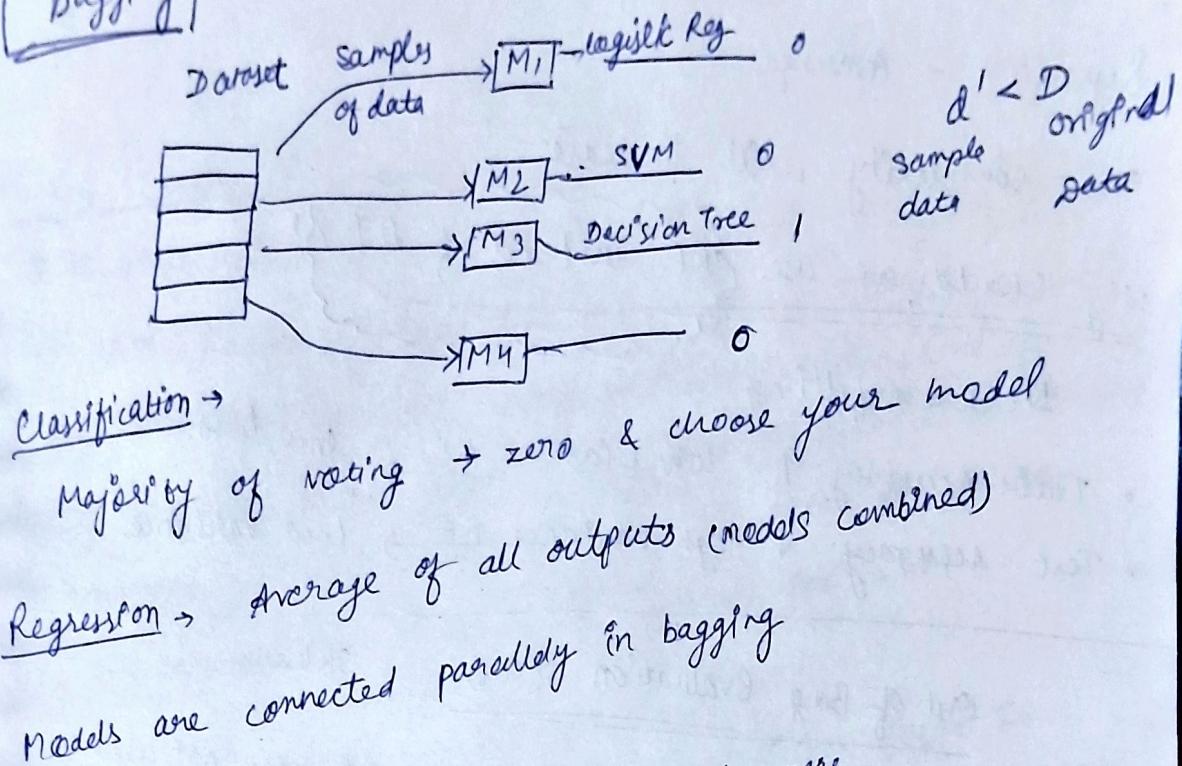
CAT Boost

CAT Boost

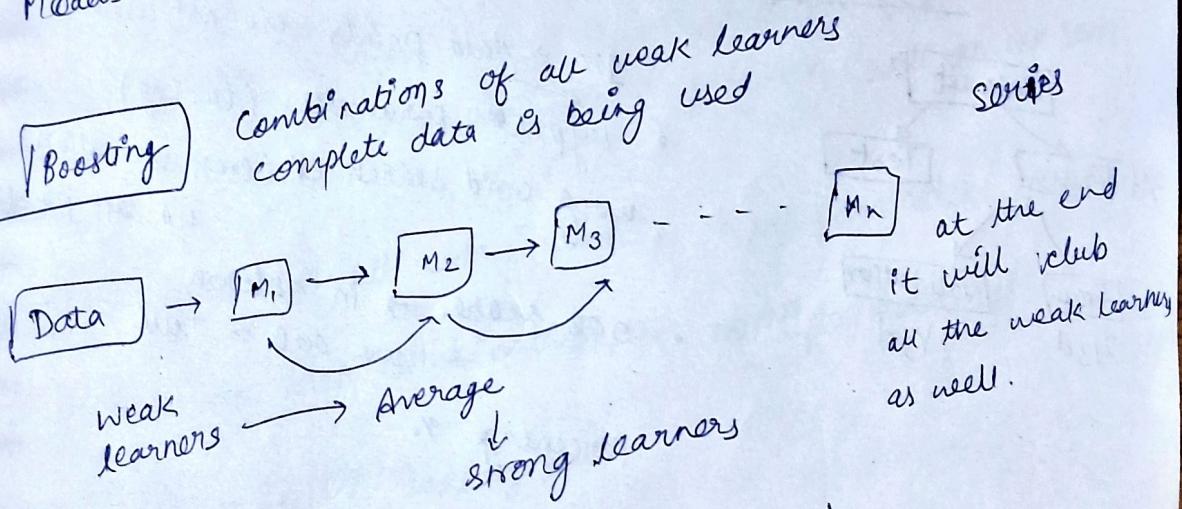
CATBoost
Gradient Boost
XGBoost

↓
multiple
decision
free

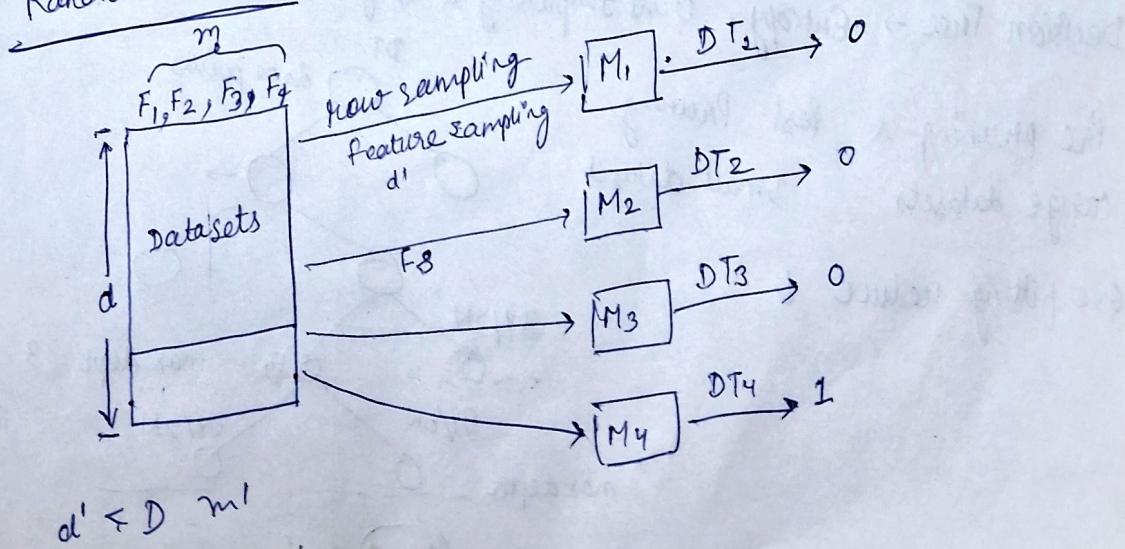
Bagging



Boosting



Random Forest (Classification & Regression)



Classification \rightarrow majority voting classifier

Regression \rightarrow Average

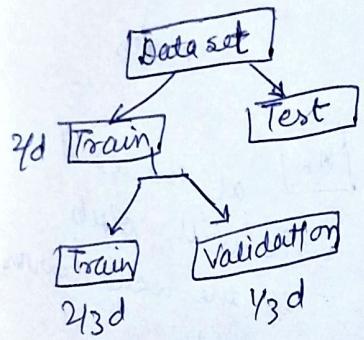
Time complexity will increase

why should we use RF instead of DT?

DT \rightarrow overfitting

- Train Accuracy \uparrow low Bias $\xrightarrow{\text{RF}}$ low bias
- Test Accuracy \downarrow high variance $\xrightarrow{\text{RF}}$ low variance

Out of Bag Evaluation (OOB) sklearn



- 100 \rightarrow data points were not used
- hyperparameters (tuning)
using Grid Search CV (cross validation)
- no data points
- OOB score \rightarrow in sklearn
use OOB = True

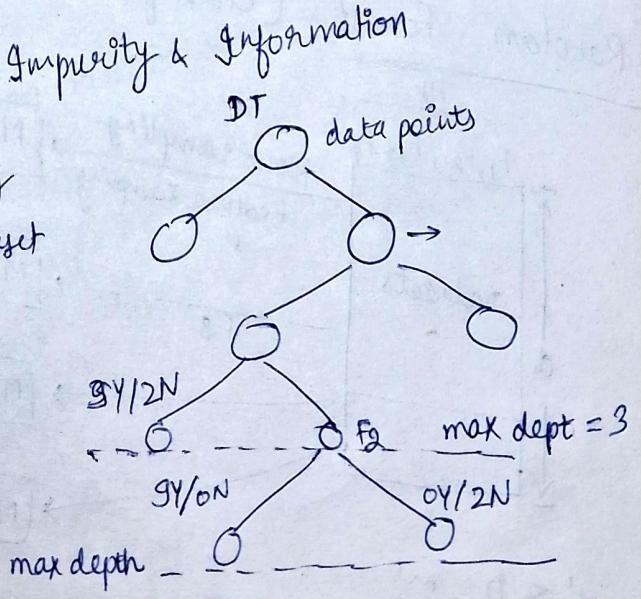
Percentage %

Random forest (classification & regression)

Decision Tree \rightarrow Entropy, Gini Impurity & Information

Pre pruning & Post Pruning
large datasets small dataset

overfitting reduce \downarrow



Grid Search CV

- It is a technique in machine learning used to find the best combination of hyperparameters for a model.
- It systematically ~~try~~ tries every possible combination of hyperparameter values you specify and evaluate the model performance for each combination using cross validation.
Ex: Hyperparameter C values & kernels for SVM (code)

Cross Validation → Split the data into multiple folds to ensure the model performs well on unseen data.

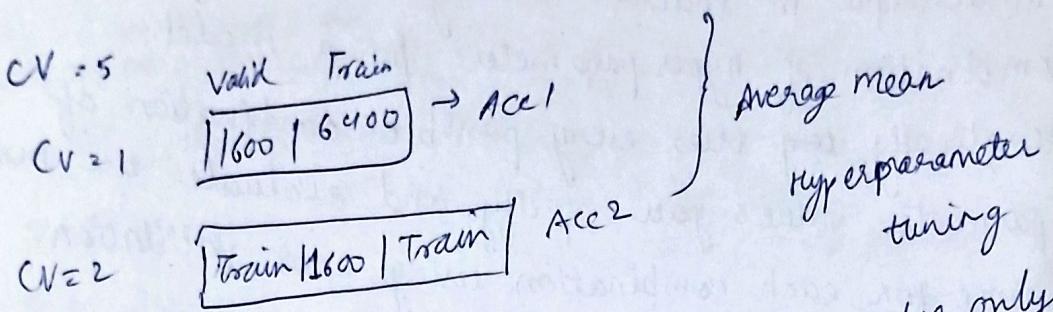
Grid → A list of all combinations of hyperparameters to try.

CV → Evaluate each combinations on multiple data split.

deadlog (decision tree) ?? Dead logram.

Train, Test & Validation

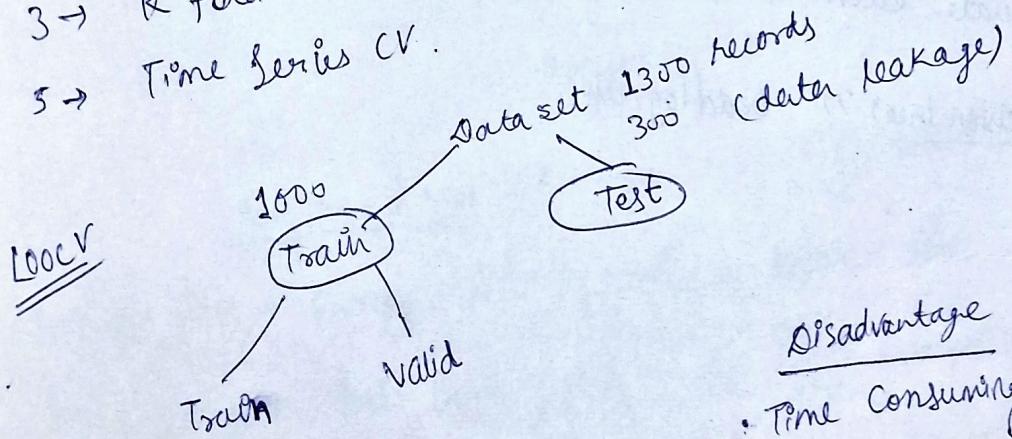
dataset
 f_1 f_2 f_3 f_4 O/P



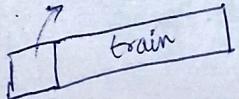
random_state → divisible by 2 (pick that value only)

Types of CV

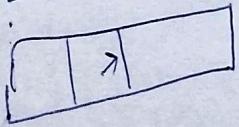
- 1 → LOOCV (Leave one out CV) → not recommended for using
- 2 → HOOCV (Hold out CV)
- 3 → K fold CV
- 4 → Stratified K fold CV
- 5 → Time Series CV.



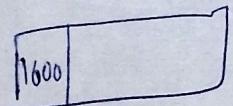
Leave One out $\rightarrow CV=1$



$CV=2$



HOOCV $CV=1$



Disadvantage
 • Time Consuming Task

• overfitting

Same for HOOCV

No validation dataset

Train dataset accuracy

HOCV → holds the data points

70% & 30% partition → Train's accuracy will be ↑

With the help of random state it will work for some time.

200	800
-----	-----

200 800 → Experiments.

K fold CV

$k=10$ experiments

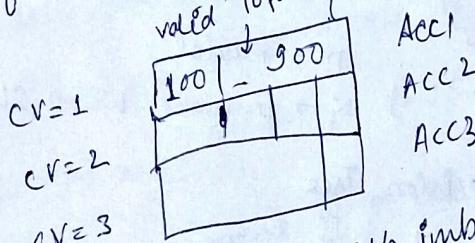
$k \rightarrow$ hyperparameter

$$\frac{1000}{11} = 100$$

→ These cannot work on imbalanced dataset. ↗ HOCV
↗ LOOCV

K fold CV

$$\frac{1000}{11} = 100$$



K → value for k is 10

iterations → cv values

10 folds

[]

data points

K fold doesn't work with imbalanced dataset.

Stratified K fold

Train & validate

Equal or \approx Equal
all the categories of output

Imbalanced dataset upto some extent

Time Series cr → Days (Sequence) or Series (for better accuracy use hyperparameter)

day1 | day2 | day3 | day4 | day5 |
.....
validate

Boosting

Bagging

RFC

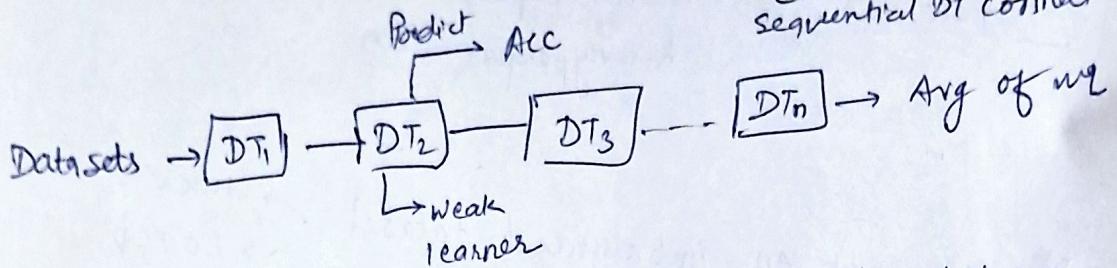
RFR

Boosting

Ada Boost → Adaptive

Gradient → XG Boost

CAT Boost R
sequential DT connect



weak learners → Haven't learnt much from training dataset

Ada Boost → Assignment a weight to the weak learner

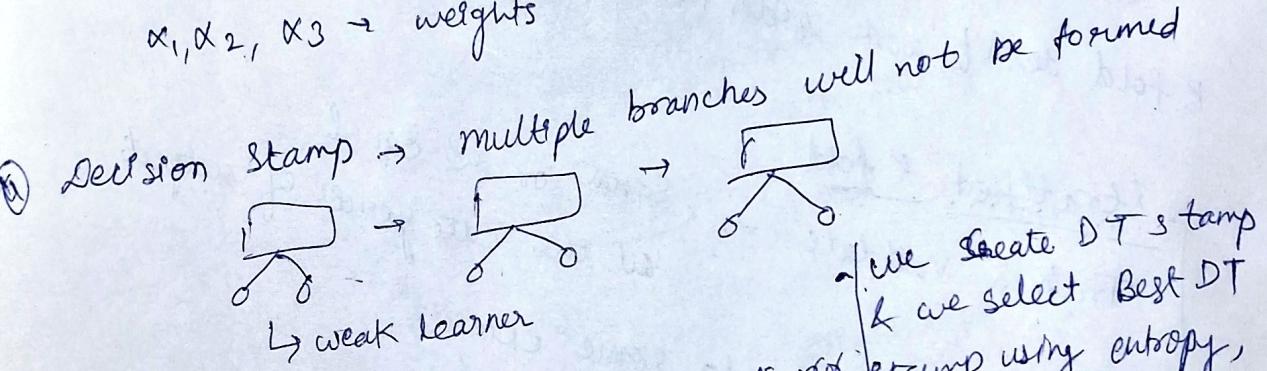
$$f = \alpha_1 m_1 + \alpha_2 m_2 + \alpha_3 m_3 + \dots + \alpha_n m_n$$

m = model
 α → to make WL → SL

classification, regression

Decision Tree

$\alpha_1, \alpha_2, \alpha_3$ → weights



	Normal wt 0.08	Updated wt 0.058	Normal wt 0.08	Salary	Credit
"				$\leq 50K$	B
"				$\leq 50K$	G
"				$\leq 50K$	G
"				$> 50K$	B
"				$> 50K$	G
"				$> 50K$	N
"				$\leq 50K$	N

Appeared	Root	
No	$1/2 \downarrow$	
Yes	$1/2 \downarrow$	
Yes	$1/2 \downarrow$	
No	$1/2 \downarrow$	
Yes	$1/2 \downarrow$	
Yes	$1/2 \uparrow$	
No	$1/2 \uparrow$	

$S \leq 50K$

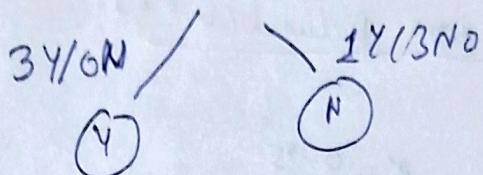
24/12N

ND

which condition is not satisfied

$$C = G_1$$

> 50 N (yes) \rightarrow error



Step one \rightarrow sum of the total error & performance of stamp.

$$\text{Total sum of error} = \cancel{17} \quad 17$$

$$\text{Performance} = \frac{1}{2} \ln \left\{ \frac{1 - TE}{TE} \right\} = \frac{1}{2} \ln (6) \approx 0.896$$

$$x_1 = 0.896 \rightarrow \text{weight}$$

3) update the weights correctly & incorrectly classified points
- performance of stamp (increase wt of errors)

a) $wt \times e^{-0.896} = 0.058$

b) $wt \times e^{-0.896} = 0.349$

4) Normalize wt computation & assigning bins

Assigned bins

$$0 - 0.08$$

$$0.08 - 0.16$$

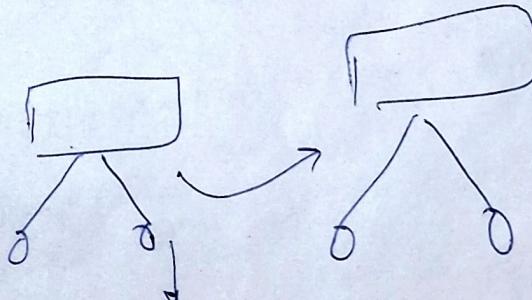
$$0.16 - 0.24$$

$$0.24 - 0.32$$

$$0.32 - 0.40$$

$$0.40 - 0.48$$

$$0.48 - 0.56$$



To prepare
datapoints

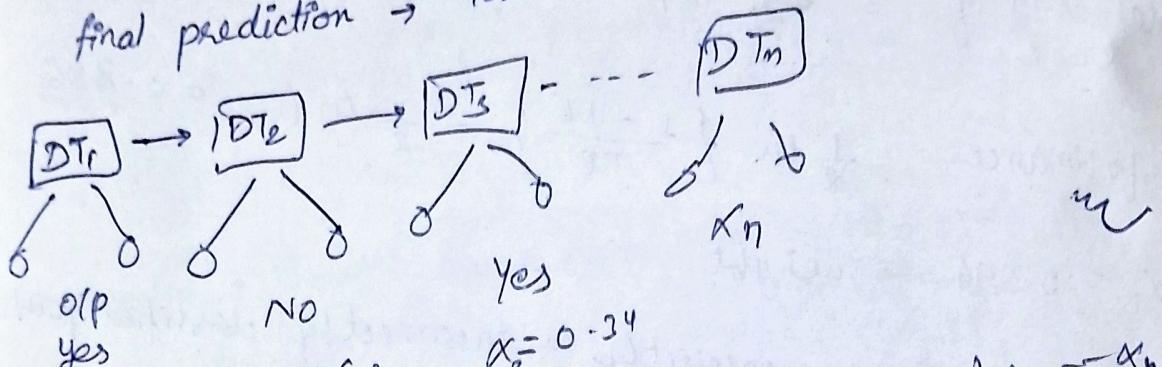
5) Select datapoint to send to next stamp
 ↗ Iteration process selecting random value b/w 0 & 1.

Random

$$f = \alpha_1 m_1 + \alpha_2 m_2 + \dots \quad \alpha = 0.95$$

$$f = 0.896$$

final prediction → Test data $\approx 50K$, 61



$$\alpha = 0.896$$

$$f = 0.896(\text{Yes}) + 0.62(\text{No})$$

$$f = 1.136(\text{Yes}) + 0.350(\text{No})$$

previously weak learners → weight & assign bins values
 ↗ adapt → strong learners.

Gradient Boosting

Regression

Classification

Reg Dataset

	Degree	Salary
1	B.E	50K
2	Master	70K
3	Master	80K
4	Ph.D.	<u>100K</u>
		mean 75K

	Predicted	Residual	Predicted
	\hat{y}	$y - \hat{y} (R_1)$	R_2
1	75K	-25K	-23K
2	75K	5K	-3K
3	75K	5K	3
4	75K	25K	23K

Assumed from previous error

Step 1 → Create Base Model

$$\frac{50K + 70K + 80K + 100K}{4} = 75K$$

75K

Step 2 Compute residual error

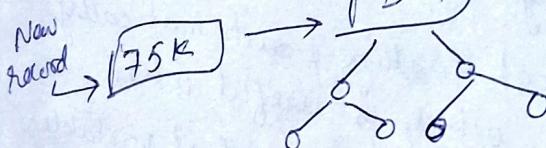
learning rate [0, 1]
 α

Step 3 Construct a DT
consider HYP x_i & O/P R_2

Predict O/P =

$$75 + (-23) = 52$$

Boosting
week
learner



generalized mode

$$f(x) = x_0 h_0(x) + x_1 h_1(x_1) + x_2 h_2(x_2) + \dots + x_n h_n(x_n)$$

$$f(x) = \sum_{i=0}^n x_i h_i(x)$$

XG Boost

extant extension gradient boosting

- It is optimised gradient boosting library designed to be highly efficient, flexible and portable. It provides parallel tree boosting (also known as gradient boosting decision tree, GB DT, GBM(maching))

Key features

- Regularisation → Help to reduce overfitting
- Parallel processing → utilize multiple CPU cores to train model faster.
- Handling missing values → automatically learn the best way to handle missing data.
- Cross validation → Built in CV support.
- ~~CAT Boost~~

CAT Boost → This strong gradient boosting algo handles categorical features automatically, without the need for explicit pre processing such as one hot encoding.

Key features, handling categorical feature; automatically process categorical feature which is useful when dealing with dataset that contains many categorical variables.

Robust to overfitting → provide better generalization model on test set

Efficient → faster training due to order boosting

boosting.

Comparison of XGBoost, Gradient Boosting & catBoost

features	XGBoost	Gradient Boost (Scikit-learn)	catBoost
Categorical data handling	Requires preprocessing	Requires preprocessing	Built-in support
Speed	fast with II processing	moderate	fast, optimized for categorical data
Model Tuning	Highly customizable	Simple	very customizable
Visualization Support	Tree & feature importance	feature importance	Faster
Training Time	Moderate	longer	

Principal Component Analysis (PCA)

Dimension Reduction						
x	y	$x - \bar{x}$	$y - \bar{y}$	$x^T y_1$	x_i^2	y_i^2
2.5	2.4	0.63	0.4	0.33	0.4781	
0.5	0.7	-1.31	-1.26	1.585	1.7111	
2.2	2.9	0.39	0.99	0.3866	0.1522	
1.9	2.2	0.09	0.29	0.026	0.0081	
3.1	3.0	1.29	1.09	1.4061		
2.3	2.7	0.49	0.79	0.3871		
2	1.6	0.19	-0.31	-0.0589		
1	1.1	-0.81	-0.81	0.6561		
1	1.6	-0.31	-0.31	0.0961		
1	0.9	-1.01	-1.01	0.7172		
MEAN		1.81	1.91	5.539	5.545	6.449

$$\text{cov}(x_i, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$|A - \lambda I| = |A - \lambda I|$$

characteristic equation.

$$\begin{bmatrix} 0.5549 & 0.5539 \\ 0.5539 & 0.64499 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0.5549 & 0.5539 \\ 0.5539 & 0.64499 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = \begin{bmatrix} 0.5549 - \lambda & 0.5539 \\ 0.5539 & 0.64499 - \lambda \end{bmatrix}$$

$$\det(A - \lambda I) = (0.5549 - \lambda)(0.64499 - \lambda) - (0.5539)$$

$$= (\lambda - 0.04908)(\lambda - 1.2848) = 0$$

$\lambda_1 = 0.04908$ $\lambda_2 = 1.2848 \rightarrow$ Eigen values

Eigen vector \rightarrow direction

$$\lambda_1 = 0.04908,$$

$$B = \begin{bmatrix} 0.50582 & 0.5539 \\ 0.5539 & 0.59502 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$Bx = 0$ Augmented matrix

$$0.50582x_1 + 0.5539x_2 = 0$$

$$0.5539x_1 + 0.59502x_2 = 0$$

$$-0.735178 \\ 0.67787$$

$$\left\{ \begin{array}{l} \text{Eigen vector } v_1 = 0.04908 \\ \text{Eigen vector } v_2 = 1.2848 \end{array} \right. = \begin{bmatrix} -0.735178 \\ 0.67787 \end{bmatrix}$$

$$v_2 > v_1 \rightarrow \text{PC1}$$

PC1 PC2

larger λ value is PC1

smaller λ value is PC2

$$\begin{bmatrix} 0.5549 & 0.5539 \\ 0.5539 & 0.64499 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0.5549 & 0.5539 \\ 0.5539 & 0.64499 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0.5549 - 1 & 0.5539 \\ 0.5539 & 0.64499 - 1 \end{bmatrix}$$

$$\text{det}[A - \lambda I] = (0.5549 - \lambda)(0.64499 - \lambda) - (0.5539)$$

$$= (\lambda - 0.04908)(\lambda - 1.2840) = 0$$

$\lambda_1 = 0.04908$ $\lambda_2 = 1.2840 \rightarrow \text{Eigen values}$

Eigen vector \rightarrow direction

$$\lambda_1 = 0.04908, \quad B = \begin{bmatrix} 0.50582 & 0.5539 \\ 0.5539 & 0.59502 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$Bx = 0$ Augmented matrix

$$0.50582x_1 + 0.5539x_2 = 0$$

$$0.5539x_1 + 0.59502x_2 = 0$$

$$\left\{ \begin{array}{l} \text{Eigen vector } v_1 = 0.04908 = \begin{bmatrix} -0.735178 \\ 0.67787 \end{bmatrix} \\ \text{Eigen vector } v_2 = 1.2840 = \begin{bmatrix} -0.67787 \\ -0.735178 \end{bmatrix} \end{array} \right.$$

$$v_2 \gg v_1 \rightarrow \text{PC1}$$

PC1 PC2

larger value is PC1

smaller value is PC2

feature selection vs feature extraction (Dimensionality Reduction)

3D data \rightarrow 2D data

100d \rightarrow 3D to 2D

- Prevent \rightarrow curse of dimensionality
- Improve \rightarrow The performance of the model
- Visualize the data \rightarrow understand the data.

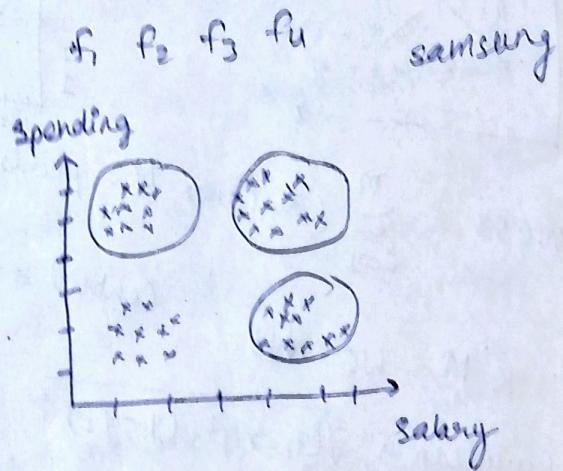
After calculating the covariance (it's used for feature selection) we do use feature extraction for dimensionality reduction.

Unsupervised Learning

no labels only features.

supervised learning

$f_1 f_2 f_3 f_4 f_5$ O/P



K-means clustering

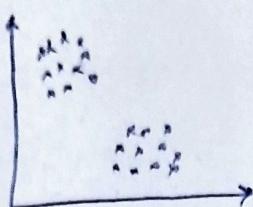
Hierarchical clustering

DBScan clustering.

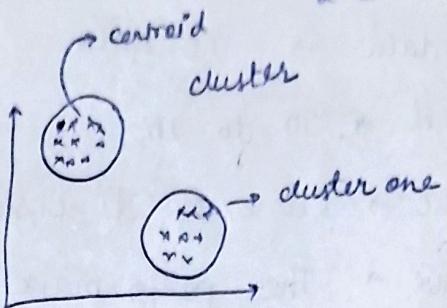
Unsupervised Learning

12.2

k-means clustering algo



$\xrightarrow{\text{k-means}}$

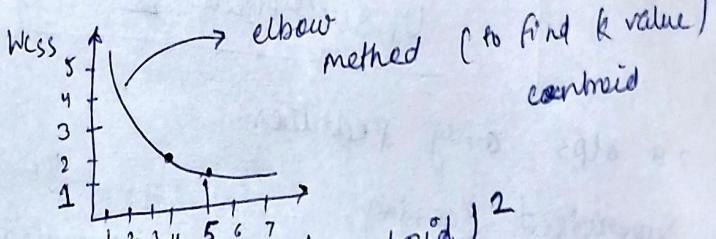
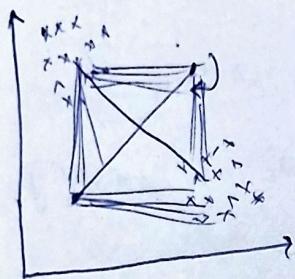


Step one: Initialize some K centroids

Step two: Points that are nearest to the centroid \rightarrow group

Step three: Move the centroids \rightarrow Average.

How do we select the K -value?
 $WCSS = \text{within the cluster sum of square}$



$$WCSS = \sum_{i=1}^m (\text{distance b/w points to nearest centroid})^2$$

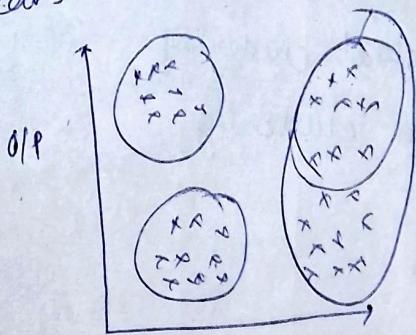
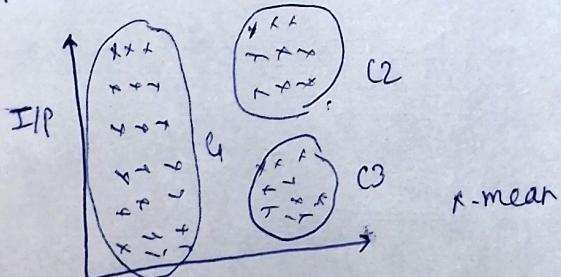
Centroid $\rightarrow K$

Euclidean & Manhattan Distance.

$K \uparrow \uparrow \quad WCSS \downarrow \downarrow$

Euclidean $\rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$ Manhattan $\equiv |x_2 - x_1| + |y_2 - y_1|$

Random Initialization TRAP (K -means $^{++}$)

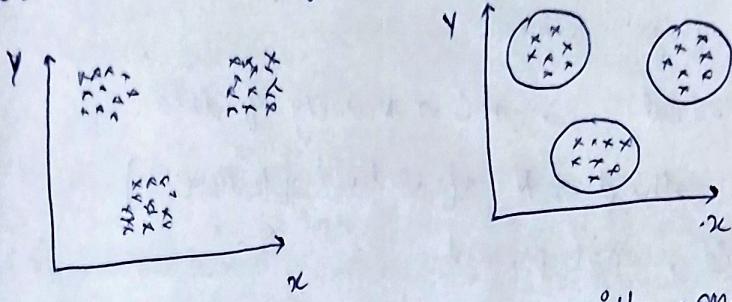


Takes wrong \rightarrow corrects them

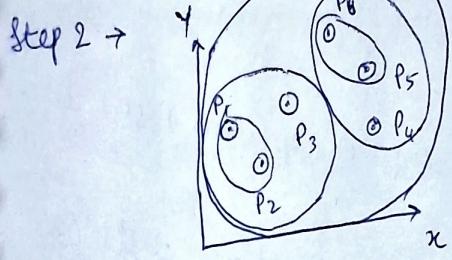
(Predicted)
values

Unsupervised Learning

Hierarchical clustering \rightarrow No centroid available

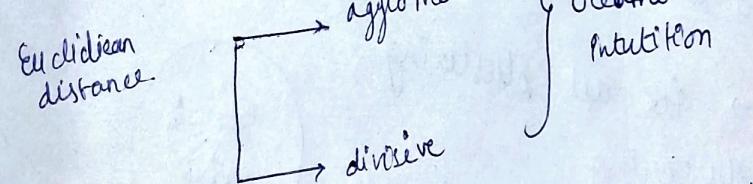


Step 1 \rightarrow for each point initialize we will consider it as a separate center.



find the nearest point create a new center.

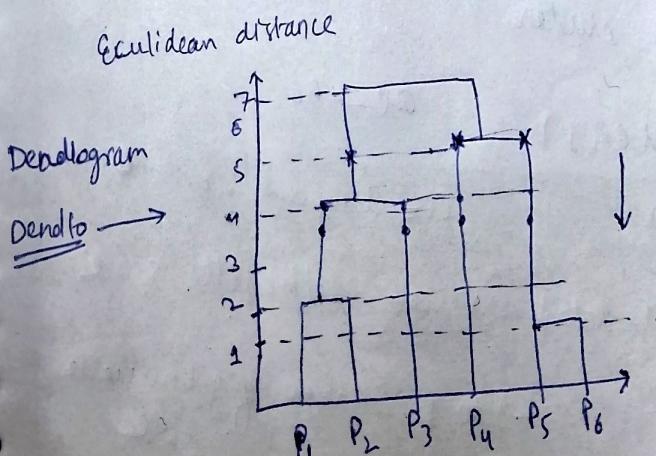
Step 2 \rightarrow



Top to bottom
(agglomerative)

Bottom to Top
(divisive)

Threshold value \rightarrow select the longest vertical line such that no horizontal line passes through it.



agglomerative
such vertical line from which no horizontal line is passing through.

threshold = 5

K-means vs Hierarchical

- Stability and flexibility
- Dataset size - { small HC
huge K-means }
- K-means → Numerical . HC → variety of data
- centroid → Elbow method : No. of centroid [K means]
- HC Centroid : not present

DB Scan

- Core point
 - Border point
 - outlier
- min points = 4
 ϵ = radius.

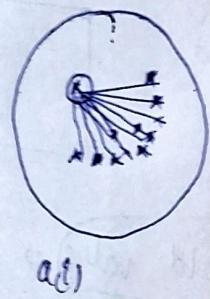
Core point → No of points within the ϵ (epsilon) should be greater than & equal to four.

Border point → No of data points within the radius will be less than minimum points.

* DBScan is robust to the outliers

Method for validation → for all clustering

using silhouette clustering.

$$\text{a}(i) = \frac{1}{|C_i|-1} \sum_{j \in C_i \cup j} (d(i, j))$$


a(i)

i^* = datapoint i in the cluster

$$2) b(i) = \min_{j \neq i} \frac{1}{|C_i|} \sum_{j \in C_i} d(i, j)$$

$a_i < b_i$

$i^* = \text{datapoint } i \text{ in the cluster}$

In clustering score is calculated through silhouette score

$$\text{Score } s(i) = b(i) - a(i)$$

-1 to 1 Range of silhouette score.

$$\max \{ a(i), b(i) \}$$

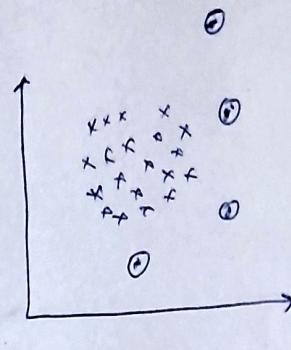
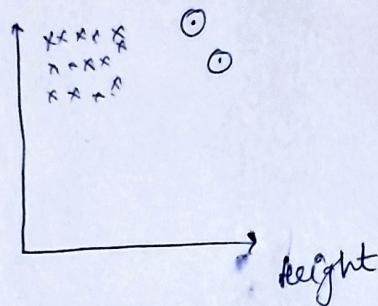
$$S(\delta) = \begin{cases} 1 - \frac{a_i^{\delta} b_i^{\delta}}{a_i^{\delta} + b_i^{\delta}} & \text{if } a_i^{\delta} < b_i^{\delta} \\ 0 & \text{if } a_i^{\delta} = b_i^{\delta} \\ \frac{a_i^{\delta} b_i^{\delta}}{a_i^{\delta} + b_i^{\delta}} & \text{if } a_i^{\delta} > b_i^{\delta} \end{cases}$$

if $a_i^{\delta} > b_i^{\delta}$ = good cluster

if $a_i^{\delta} < b_i^{\delta}$ = bad cluster

Anomaly Detection : To detect outliers

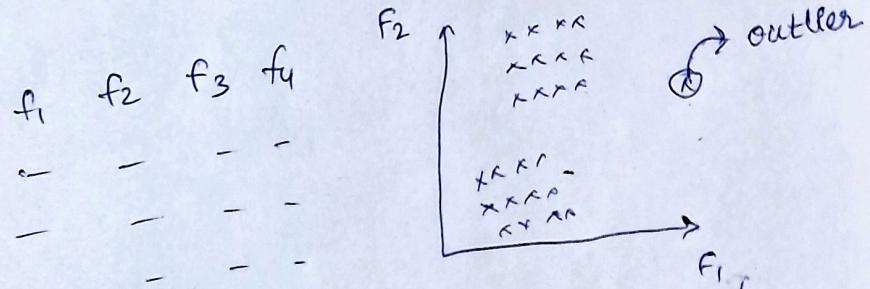
weight



It is important for problem statement

• Isolation Forest : Decision Tree (making different type of isolation tree)

lets consider



will work until it gets single leaf node not found.

$$\text{math } S(x, m) = \left\{ \frac{E h(x)}{C(m)} \right\}$$

Since m = total no. of data points

x = data points (isolation data points)

$E(h(x))$ = Average search depth of x from isolated tree.

$C(m)$ = Average depth of all data points

$E(h(x)) \ll C(m) \Rightarrow S(x, m) \approx 1$ or near about one.

$E(h(x)) \gg C(m) \Rightarrow S(x, m) \approx 0.5$

anomaly score
normal datapoints
if it is near one than its outlier

first we define threshold 0.5

Associate Rule Mining (Learning)

- Market - Basket Analysis \rightarrow work on transactions
 - Aim \rightarrow interest patterns [extended frequent patterns]
 - Terminologies \rightarrow item set k /item {bread}, {milk} and etc.
 - \rightarrow Support count (Σ) frequency of dataset given set of transactions
 - \rightarrow frequent item set greater than and equal to minimum support threshold.
 - \rightarrow Support ("s") s is the percentage of transactions in dataset that contains:
- ~~A & B~~ $A \& B \quad (A \cup B) \Rightarrow s = \text{Prob}(A \cup B)$
- \rightarrow confidence (c) c is the percentage of transaction containing A that also contain B.
- $\Delta c = \cancel{\text{Prob}} P\left(\frac{s}{P(A)}\right) \quad \text{or} \quad P\left(\frac{P(A \& B)}{P(A)}\right)$

Concept: lift $\vdash A \rightarrow B$ lift $\dagger = \frac{s(A \& B)}{s(A)s(B)}$ $\text{lift}^2 = \frac{\text{independent items}}$
support $A \& B$

if lift is > 1 {positive effect of item on each other}
lift is < 1 {negative effect of item on each other}

Strong Association Rule (which item to be picked)

calculate s & c equal. we define a minimum threshold.
less number of association rule in dataset when increase
satisfied both minimum threshold both support and decrease
minimum threshold coeff confidence also. (e.g am sleepy ??)

2:36 3rd Jan
pm
So much sleepy!! UGHH...

- Mining process →
- find all frequent item sets
 - generate strong associate rules from the frequent item set.
 - Support count :- frequency of item sets

Algorithm used → Apriori

uses K item set, $(K-1)$ → used to find frequent item set
Properties of Apriori → All non empty subset of frequent item set must also be frequent. for example

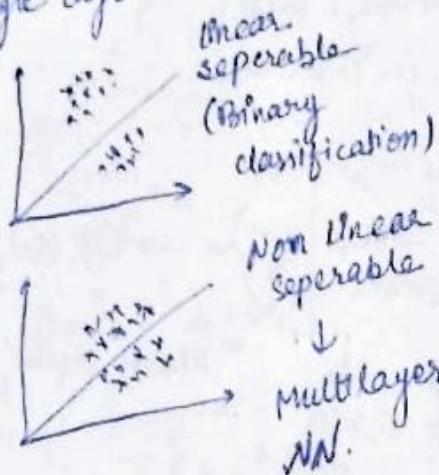
$$\{A, B\} \rightarrow \{A\} \{B\}$$

Anti monotonicity Property → If item set i doesn't satisfy the minimum support threshold than i is not frequent. If an item A is added to the item set i . $(i \cup A)$ can't occur more frequently than i .

Conclusion → If our confidence is approx. equal to 1 we say it's a strong association analysis.

Perception Model

single layer



MultiLayer (Artificial Intelligence)

- Forward Propagation
- Backward Propagation
- Loss function
- Activation function
- Optimizers (Gradient Descent)

on Jan
forward propagation, backward propagation, vanishing gradient, clear chain rule propagation

test dataset	IQ	Study Hr	Play Hr	O/P
95	4	5	2	1
100	5	2	7	1
98	2	7	0	0

Pass/Fail

bias \rightarrow always applied on the output of hidden layer

Input \rightarrow O/P : forward propagation
Let us $\{w_1, w_2, w_3\}$ $\{b\}$ $\{0.001\}$

$$\begin{aligned} z &= \sum_{i=1}^n x_i w_i + b \\ &= 95 \times 0.01 + 0.001 \\ &\quad 100 \times 0.02 + 0.001 \\ &\quad 95 \times 0.03 + 0.001 \end{aligned}$$

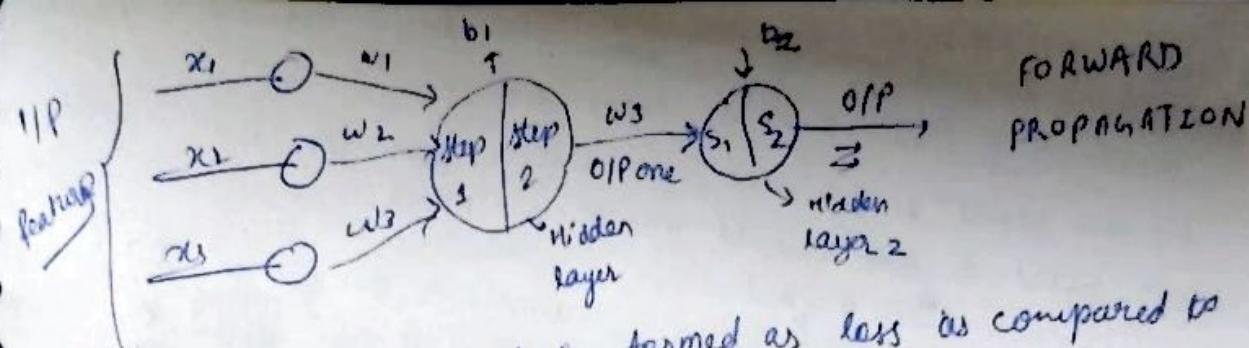
$$z = 1.151$$

$$z = [95 \times 0.01 + 100 \times 0.02 + 98 \times 0.03 + 1 \times 0.001]$$

$$= 1.151$$

$$\text{Step 2: Sigmoid fn} = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-1.151}} = 0.759$$

$w_i w_i + b \rightarrow$ add this to hidden layer's O/P to get the final output



whatever final O/P we get is termed as loss as compared to real data the value will be lesser. (loss Error)

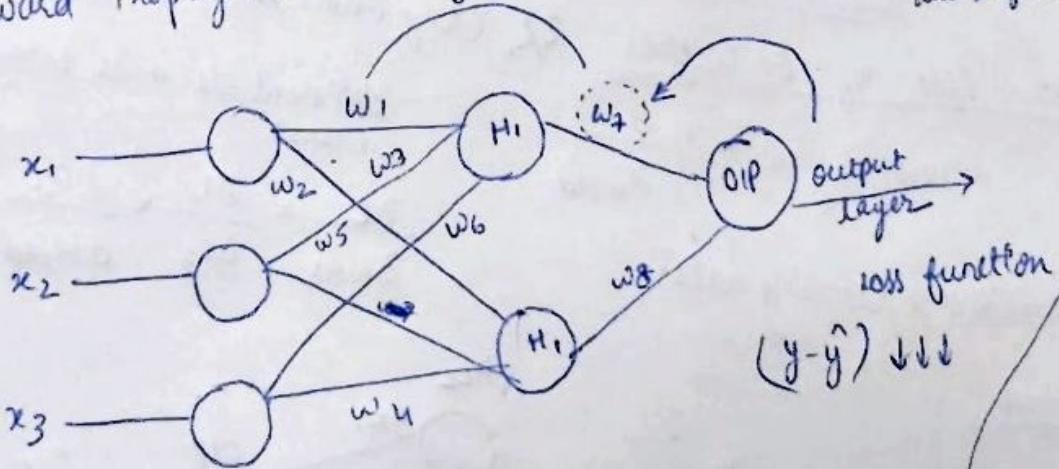
- Loss fn
- 1. Regression
 - MSE
 - MAE
 - Huber loss

- 2. classification
 - Binary cross Entropy
 - Categorical cross entropy

Vanishing → sigmoid function (0 to 1) sometimes it reduces too much so it makes tough to identify the global minima. [Gradient Descent]. (Replace with tanh)

middle value
0 = 0.25
sigmoid reduce
too much the value
will shift jump

Backward Propagation & weight updation

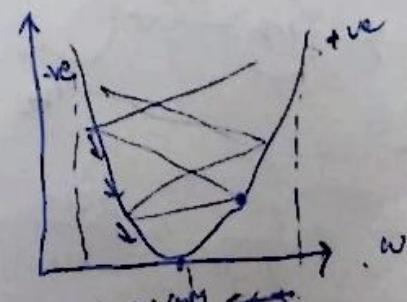


Updation formula

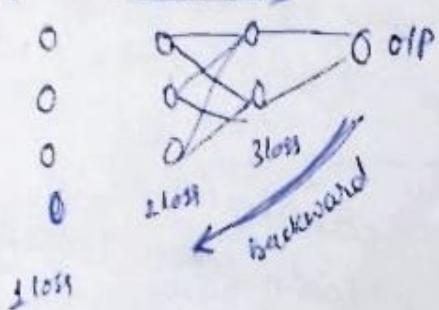
$$w_{j,\text{new}} = w_{j,\text{old}} - \eta \frac{\partial L}{\partial w_{j,\text{old}}} \quad \begin{array}{l} \text{slope} \\ \text{loss} \end{array}$$

↳ learning rate

$$\eta = 0.002$$



Back propagation → whenever it comes moving backwards & reduce thousands of hidden layers & drop out layers O/P loss should be minimize, but there's some loss within your hidden layers so that should also be minimized.



whatever loss you get try to reduce it near about 0.

continuous cycle
global minima

update weight & bias
to reduce loss*

optimized (Gradient Descent) → to reduce the loss value

$$W_{j, \text{new}} > W_{j, \text{old}}$$

vanishing
when sigmoid function stops giving O/P.

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

$$w_{\text{old}} + \eta (+ve)$$

* When w reaches global minima :-

$$W_{\text{new}} = W_{\text{old}}$$

Chain Rule of Derivatives

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

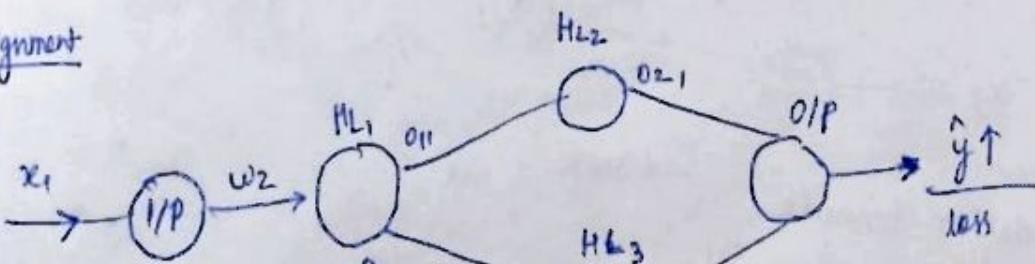
movement → learning rate

2. (why maths ??)

differentiate → at every point (loss)

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial o_3} \times \frac{\partial o_3}{\partial w_{\text{old}}}$$

Assignment



$$\frac{d}{d w_{\text{old}}} = 20$$

$$\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{\text{old}}} \rightarrow \frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial o_{12}} \times \frac{\partial o_{12}}{\partial w_{\text{old}}}$$

Activation function

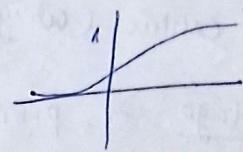
① Tan h

② Relu

③ Prelu

④ Swish

majorly used
in place of sigmoid



Sigmoid fn [0 to 1]

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

$$\sigma(x) = \frac{1}{1+e^{-x}} \Rightarrow 0 \text{ to } 0.25$$

differentiable

$$\sigma(z) = 0, 1$$

$$\frac{\partial \sigma(z)}{\partial z} = 0 \text{ to } 0.25$$

$$\phi(z)$$

Advantages \rightarrow it's suitable for binary classification
clear predictions that are very close between 0 to 1

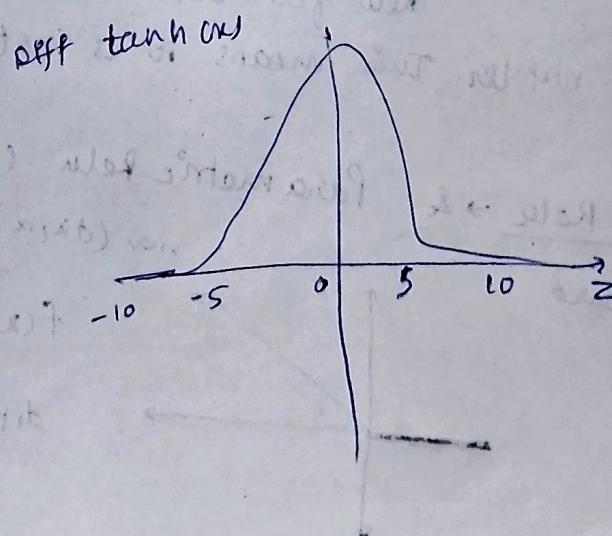
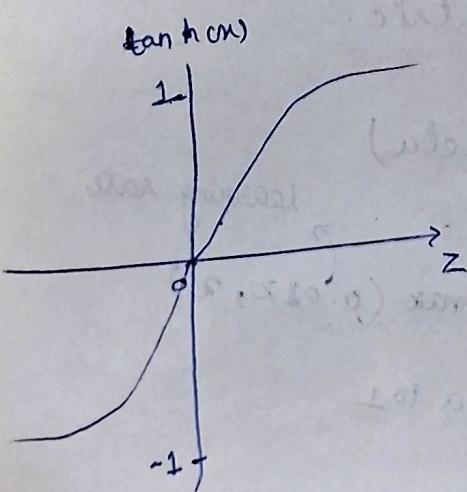
DPS disadvantages \rightarrow prone to vanishing gradient descent

- function output is not zero centered
- Mathematical operations are relatively time consuming.

Tan h activation function [Tangent hyperbola]

forward propagation

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad [-1 \text{ to } 1]$$



Advantages

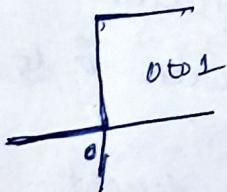
zero centric (weight update efficient)

Disadvantage → prone to vanishing gradient problems.
Time complexity.

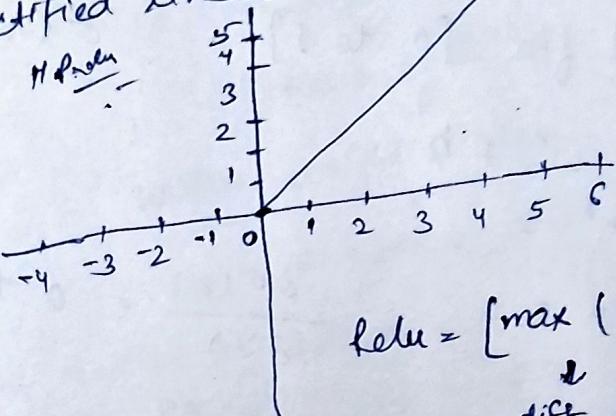
Relu Activation fn (Rectified linear unit)

• tanh → -1 to 1
 $d\theta \rightarrow 0$ to 1

• sigmoid → 0 to 1
 $d\theta \rightarrow 0$ to 2.5



$x = 0$
dead neuron
 $w, x_1 = 0$



$$\text{relu} = [\max(0, x)]$$

diff relu

Advantage → solving vanishing gradient problems
 $\max(0, x) \rightarrow$ calculation is super fast.

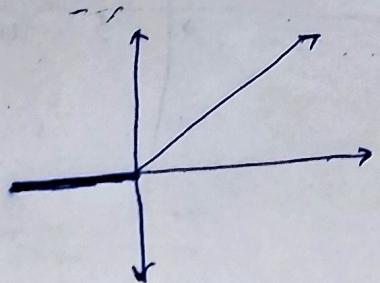
The ReLU function has a linear relationship
it is much faster than sigmoid or tanh

Disadvantage → dead neurons

ReLU function output $0, x$ is equal to zero or
positive number. This means it is not centric.

Leaky ReLU & Parametric ReLU (P-relu)

forward



$$\max(dx_1, x)$$

$$f(x) = \max(0.01x, x)$$

$$\text{diff} = \alpha \text{ to } 1$$

learning rate

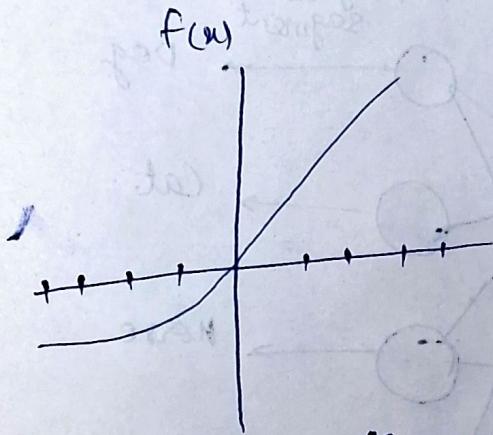
ReLU \rightarrow dead neuron \rightarrow dead neuron problem

- Leaky ReLU has all the advantage of ReLU.
- It removes dead ReLU problems

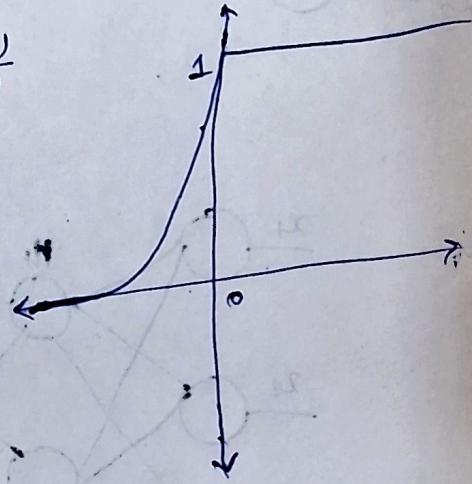
\hookrightarrow neuron

- Parametric & Leaky ReLU are same (working based)

Exponential Linear Unit (ELU)



$$\frac{d f(x)}{d x}$$



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

It is used to solve ReLU problems

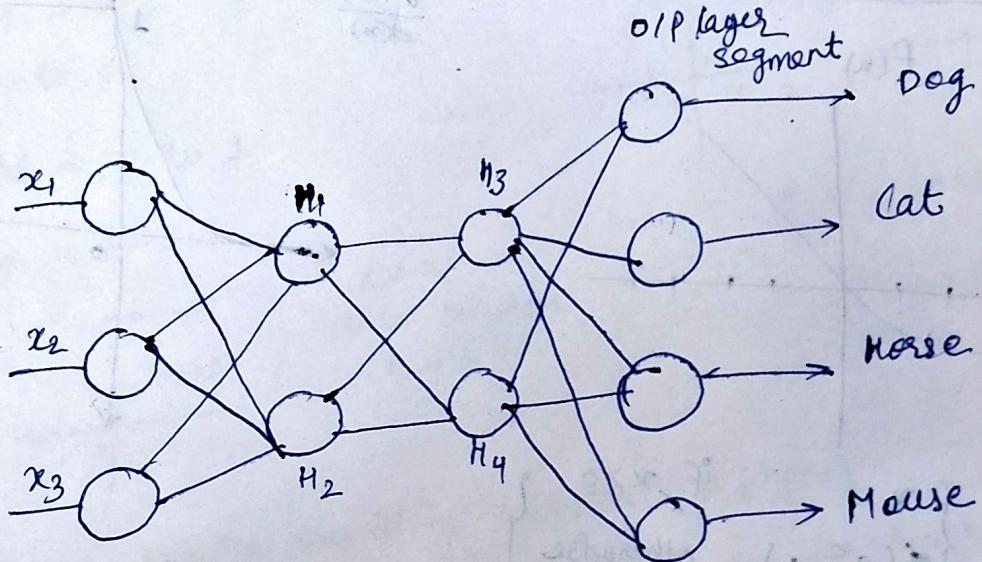
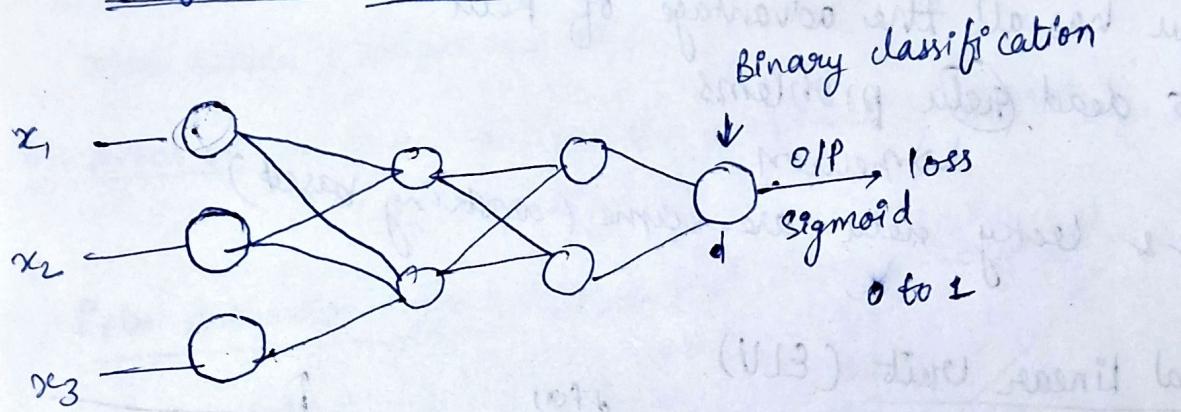
Advantage \rightarrow No dead neuron issue
zero centric

Disadvantage \rightarrow slightly more computationally intensive

All activation fn used for binary classification
For multiclass classification, we use Soft-Max activation function.

all Activation fn \rightarrow Diff \rightarrow Backward

Soft Max Activation function



Activation formula

$$\frac{e^{y_i}}{\sum_{k=0}^n e^{y_k}}$$

$$y_i = \mathbf{O} \cdot \mathbf{x} \cdot \mathbf{w} + b$$

↳ output

- Optimizers
 - Adagrad and RMS PROP
 - GD
 - SGD
 - mini batch SGD
 - Adam optimizer

GD O/P loss &

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

↓
iteration

1000 datapoints

all in forward

all in backward

weight updation

1 epoch = 1 iteration

Loss function $\frac{1}{n} \sum_{i=1}^n (y_i - g)^2$ MSE

Advantages of GD \rightarrow convergence will happen
Disadvantages \rightarrow Huge usage of RAM & GPU

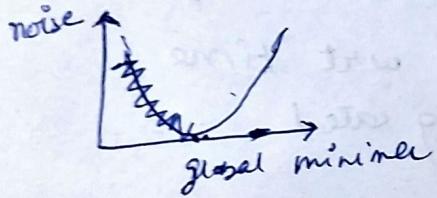
Stochastic GD \approx Sleepy head \rightarrow Ayush

3:03 PM

* worst time
after lunch

Advantage \rightarrow Solves resource issue -

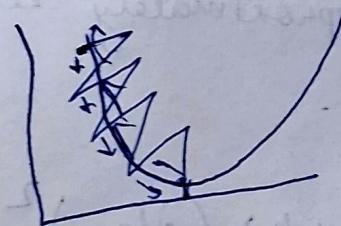
Noise: (Bias?)



disadvantage \rightarrow time complexity is increased.
convergence or convergence will also

noise get introduced..

Mini Batch



Advantages
convergence speed will increase
noise will be then when compared to SGD.

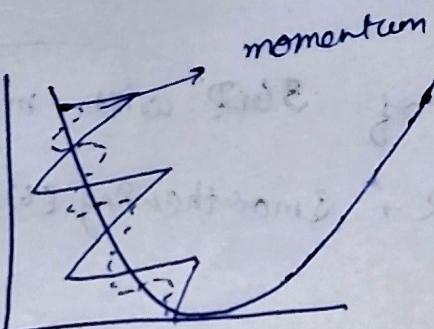
Advantage \rightarrow Efficient resource usage (RAM, Ram)

disadvantages. Noise still exist
 \rightarrow noise still exist

Experimental wt Avg.

Moment

cost



smoothening
for updating batches
increase every data point
no dp should be missed
used in ARIMA +
SARIMA

Advantage :- Reduce the noise and quick conversion

As the conversion happens the learning rate should be changed (AdaGrad)

• adaDelta, RMS Prop ?

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

$$\eta' = \frac{n}{\sqrt{x_t + \epsilon}} \rightarrow \text{epsilon}$$

whenever,

$$\therefore \eta' \propto x_t \downarrow$$

↳ changes wrt time
(learning rate)

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

Disadvantage $\rightarrow \eta'$ = possibility to become a very small value approximately zero.

AdaDelta And RMS Prop

$$S_{dw_t} = \beta S_{dw_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

$$\beta = 0.95 \quad S_{dw_t} = 0$$

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

dynamic LR + Smoothening
EWA
Exponential weight
Average

ADAM

→ Combinations of SGD with momentum
+ RMS Prop [Dynamic LR + Smoothening EWA]

(Best optimizer)

Smoothening
Gradient
Descent

Weight Initialization Techniques

1. Uniform Distribution
2. Xavier / Glorot Distribution / Initialization
3. Kaiming He Initialization

Key points why we use ??

- weight should be small
 - weight should not be same
 - weight should have good variance
- Uniform Distribution $w_{ij} \approx \text{uniform}$

$$\left[\frac{-1}{\sqrt{\text{input}}}, \frac{1}{\sqrt{\text{input}}} \right]$$

$$-\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}$$

Xavier & Glorot

Normal
 $w_{ij} \approx N(0, \sigma)$

 $\sigma = \sqrt{\frac{2}{I_1 P + O_1 P}}$

→ Research based

uniform

$$\left[\frac{-\sqrt{6}}{\sqrt{I_1 P + O_1 P}}, \frac{\sqrt{6}}{\sqrt{I_1 P + O_1 P}} \right]$$

Kaiming He Initialization

He normal
 $\sigma = \sqrt{\frac{2}{I_1 P}}$

He uniform

$$\left[\sqrt{\frac{-6}{I_1 P}}, \sqrt{\frac{6}{I_1 P}} \right]$$