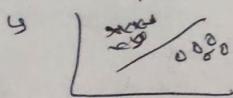


26 May Decision Tree

Story

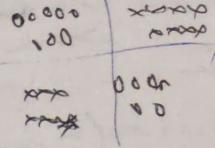
In Class
↳ DT ↳
↳ Reg

→ start with example



→ but if :

↳ by using line



→ why EDC
can't use
here Logistic Regres

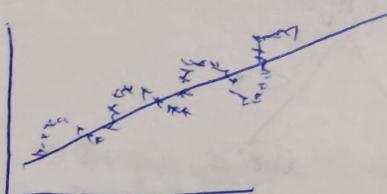
→ Problem with r.t classification

→ Some Cannot use the line to Separate.

→ why? Here data is Non linear

DT SVM

⇒ Problem for Regression



Can use polynomial Regression

↳ because it is not a ML Algo

it is use for Precepe data / OR

DO feature engnor to fit in Linear Regression..

↳ use
DT SVM
DT classifier SVM classifiers
DT regressor SVM Regressor

Decision Tree

Ex:- go out for Ply Cricket

Depth = 3

internal node

→ what do weather check

first node sent

Sunny

Ply Cricket

ye

no

node

node

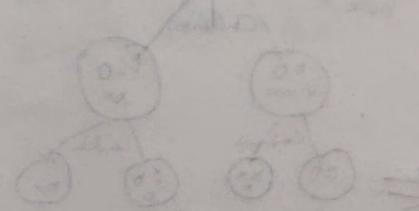
leaf nodes

(no further child node)

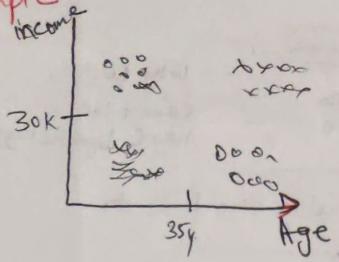
can say it
is look like a
tree

★ So we can say
every node, some
decision is
being made

So this is
called as
Decision tree



Example



→ How to solve this non-linear using DT
think

cond ① Income $\geq 30K$ & Age ≤ 35 yrs
↳ group 1 (fraud)

② Income $< 30K$ & Age ≤ 35 yrs
↳ group 2 (no fraud)

③ if Income $\leq 30K$ & Age ≥ 35 yrs
↳ group 3 → 0th fraudster

④ if Income ≥ 30 & Age ≥ 35 yrs
↳ group 4 → fraudster

So we can write these

Condition Linear regression / Logistic Regression

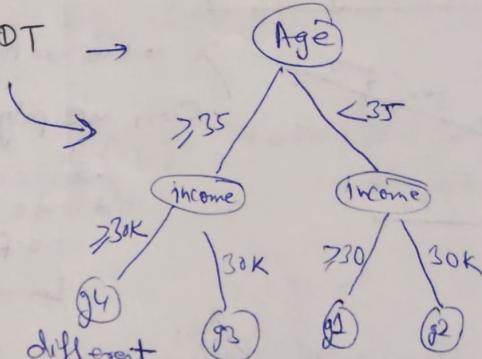
Ans No

Then can we write this in DT →

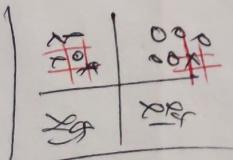
Ans Yes.

So we can say

→ DT divide the data into different Partition or into packets (Nodes)



But still we do by using 2 lines → So ok what if dp are mixed



what if dp like this then what you do

→ keep splitting/dividing until all the dp is of same class.

Adv of DT

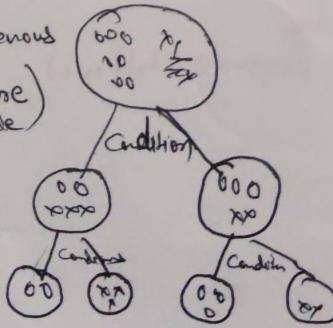
- ① Simple to Understand
- ② Capture the non-linear relationship

→ So what is int of DT

- To create leaf node/pkt where each of the data points belong to one class.

Heterogeneous data (Impure Node)

Pure - Homo (Same Class)
Impure - Hetro (mixture)

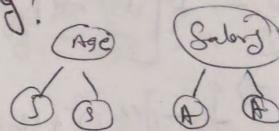


Homogeneous data
(Pure Node)

Now the point is → Let say you split data based on some condition → but when you will know → this will split or not split.

- (Q1) When to split? / How to split
• On what factor condition should be used?
- we know based on Condition then How you fit the conditions.

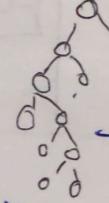
- (Q2) Which feature to be used for splitting?



- (Q3) Till when to split? → who will decide

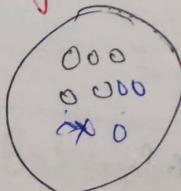
Ans. → Impurity measures

- 3 types
↳ Classification errors X_{10}
↳ Gini
↳ Entropy



Eg:		
Salary	Age	fraud
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	0	0
0	0	0

In ① Classification errors



0:8
X:2

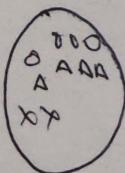
$$CE = 1 - \max p_i \rightarrow \text{probability of i class}$$

$$\begin{aligned} P_X &= \frac{2}{10} = 0.2 \\ P_0 &= \frac{8}{10} = 0.8 \end{aligned} \quad \left. \begin{aligned} CE &= 1 - \max(0.8, 0.2) \\ &= 1 - 0.8 \\ &= 0.2 \end{aligned} \right. \text{ Ans}$$

What it means → you want to pure node → yes

So → To make pure node, if I assign everything to Majority class (0), what is the error rate?

→ For Multiclass



0-4
X-2
A-4

$$\begin{aligned} P_0 &= \frac{4}{10} \\ P_X &= \frac{2}{10} \\ P_A &= \frac{4}{10} \end{aligned}$$

Minority misclassified d.p.

$$\begin{aligned} 1 - \max(0.4, 0.2, 0.4) \\ = 0.6 \end{aligned}$$

$$\begin{aligned} & \cdot [1, 2, 3] \xrightarrow{1 \text{ to } n} \\ & [0, 1, 2] \xrightarrow{0 \text{ to } n-1} \end{aligned}$$

② Gini

$$① G.I = 1 - \sum_{i=0}^{n-1} (p_i)^2$$

↳ Meaning



$$P(0) = \frac{8}{10} > 0.8$$

$$P(X) = \frac{2}{10} > 0.2$$

$i = \text{from class 0 to } n-1$

$$1 - \sum_{i=0}^{n-1} (p_i)^2 = 1 - [(p_0)^2 + (p_n)^2]$$

$$\begin{aligned} & \Rightarrow 1 - [(0.8)^2 + (0.2)^2] \\ & \Rightarrow 1 - (0.64 + 0.04) \\ & \Rightarrow 1 - 0.68 \Rightarrow 0.32 \end{aligned}$$

Same

$$② G.I = \sum_{i=0}^{n-1} p_i (1-p_i)$$

for Two class

$$p_0(1-p_0) + p_X(1-p_X)$$

$$\begin{aligned} & p_0 + p_X = 1 \Rightarrow p_X = 1 - p_0 \\ & \geq 2p_0(1-p_0) \\ & = 2 \times 0.8(1-0.8) \\ & \Rightarrow 1.6(0.2) \\ & = 0.32 \end{aligned}$$

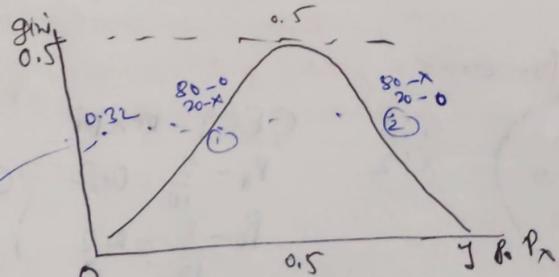


$$\begin{aligned} & \text{Gini} \rightarrow 1 - \sum_{i=1}^2 p_i \Rightarrow 1 - (p_0^2 + p_X^2) \Rightarrow 1 - ((0.5)^2 + 0.5^2) \\ & \Rightarrow 1 - 0.5 = 0.5 \quad \text{Gini} \Rightarrow 0.32 \end{aligned}$$

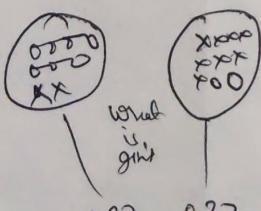
• highest gini $\rightarrow 0.5$ when the same no of dps

which is more impure
long bcz unequal no of dp
of all class

\rightarrow SO



1 more observation



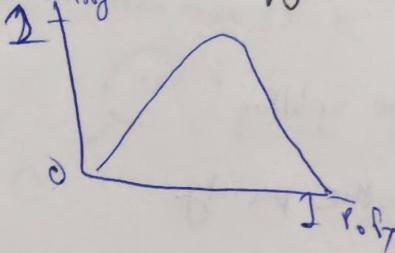
$0.32 \quad 0.32 \rightarrow$ hence both are same

→ Adv
→ Dis Adv
→ use
→ Real world

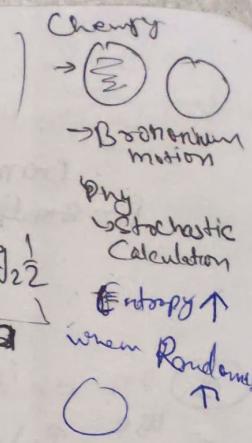
③ Entropy:

$$H = - \sum_{i=1}^n p_i \log_2 p_i$$

So Highest entropy = 1



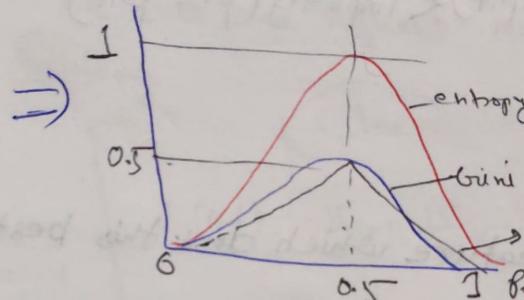
$$\begin{aligned} H &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \\ &= -2 \cdot \frac{1}{2} \log_2 \frac{1}{2} \\ &= -\log_2 \frac{1}{2} \\ &= -\log_2 2^{-1} \Rightarrow -1 \times -1 \\ &= 1 \text{ Jy} \end{aligned}$$



- Entropy - from 0 to 1

Gini - 0 to 0.5

Misclassification error = 0 - 0.5



entropy scaled
↳ $= \frac{\text{Entropy}}{2}$

⇒ For Multiclass classification

$$G.I. = 1 - \sum_{i=1}^n p_i^2$$

$$= 1 - [p_1^2 + p_2^2 + p_3^2 + \dots + p_n^2]$$

$$\text{Entropy} = - \sum_{i=1}^N p_i \log_2 p_i$$

$$- p_1 \log_2 p_1 - p_2 \log_2 p_2 - p_3 \log_2 p_3 - \dots - p_n \log_2 p_n$$

Adv
dis
because

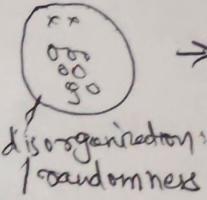
① Which impurity measure to be used?

For small dataset \rightarrow Entropy (for large dataset - bin)

(computations)

② On what feature to split

→ Use Information gain



→ Aim to splitting

↳ We want pure nodes after splitting.

↳ Such a split that maximise the purity

• $\frac{\text{Impurity (Postsplit)}}{\text{Impurity (Presplit)}}$



Split using that feature which does this best

$\text{Post Imp} < \text{Imp Pre}$

$$\Delta \text{Impurity} \Rightarrow \text{pre imp} - \text{post imp}$$

If randomness decreases, purity increase, then
we say Information gain.

$\Delta \text{Impurity}$

* I_G is a measure to define degree of disorganisation

+ randomness in a system $\Rightarrow I_G \rightarrow H$

Entropy

$I_G \Rightarrow \Delta \text{Impurity}$

* Whichever feature will have maximum information, that feature will be used split

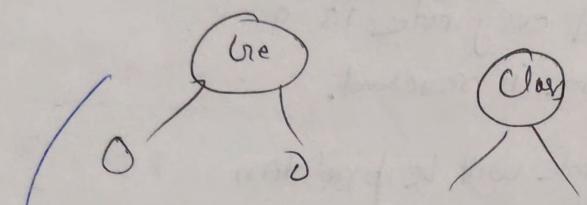
• Resistor /
↳ Semiconductor
at $T = 25^\circ\text{C}$

Example

Breeder	Class	Cricketer (y/n)
M	Y	Yes
F	X	No
N	I*	Yes

Total Std = 30 $\rightarrow 50\% \text{ play cricket}$
 $y_{es} = 15$
 $w_o = 15$
 out of 10 $\rightarrow 2 \text{ play}$
 gender $\rightarrow 10/2 \rightarrow F \rightarrow \text{play ckt}$
 $\rightarrow 20/13 \rightarrow N \text{ (out of } 20 \text{ play)}$

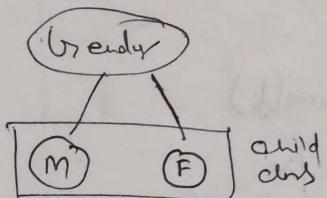
which cases ISL \leftarrow



Class $\rightarrow 1x - 14/43\% \text{ (out of 14, } 1x \text{ play)}$
 $x = 16/57\% \text{ (out of 16 cases, } 57\% \text{ play ckt)}$

Use \rightarrow whichever feature has the highest information gain, used for splitting.

In Entropy



$$I.b_1 = \Delta \text{ Entropy}$$

$$En(\text{Parent}) - En(\text{Child combined})$$

$$\begin{aligned} En(\text{Parent}) &= \frac{15}{30} \log_2 \frac{15}{30} + \frac{15}{30} \log_2 \frac{15}{30} \\ &= 1 \end{aligned}$$

$$\begin{aligned} En(\text{Child}) &= -\frac{2}{10} \log_2 \frac{2}{10} - \frac{8}{10} \log_2 \frac{8}{10} \\ &= 0.72 \end{aligned}$$

$$\begin{aligned} En(\text{Male}) &= -\frac{13}{20} \log_2 \frac{13}{20} - \frac{7}{20} \log_2 \frac{7}{20} \\ &= 0.92 \end{aligned}$$

No of Std is different in both Prob so use weighted

$$En(\text{Parent}) = \frac{15}{30} \log_2 \frac{15}{30} + \frac{15}{30} \log_2 \frac{15}{30} = 1$$

$$En(\text{Child}) = \frac{14}{30} H_1 + \frac{16}{30} H_2$$

$$\begin{aligned} H_1 &= -\frac{2}{10} \log_2 \frac{2}{10} - \frac{8}{10} \log_2 \frac{8}{10} \\ &= 0.72 \end{aligned}$$

$$\begin{aligned} I.G &= E_p - E_c \\ &= 1 - 0.94 \\ &= 0.06 \end{aligned}$$

Since $I.b_1 > I.b_2$
 \therefore we will split base on gender

$$\boxed{\begin{aligned} En(\text{Child/Combined}) &= \frac{\text{No of F}}{\text{Total Std}} \times E_{F \text{ Child}} + \frac{\text{No of M}}{\text{Total Std}} \times E_{M \text{ Child}} \\ &\Downarrow \text{Weighted Entropy} \\ &= \frac{10}{30} \times 0.72 + \frac{20}{30} \times 0.92 = 0.85 \end{aligned}}$$

$$I.G = E_p - E_c \Rightarrow 1 - 0.85 = 0.15$$

Algorithm of DT

Step 1 → Recursive binary splitting / Partitioning the data into smaller subset.

Step 2 → Select the feature to split (I.H)

Step 3 → Apply the split

Step 4 → Repeat the process for the subset of obtained

Step 5 → Continue the process until every node is a pure node / stopping criteria is reached.

Step 6 → Majority Value in the leaf node will be prediction

Eg: Age, Sal, Gender

→ feature is categorical e.g.

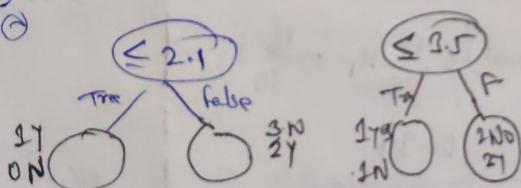
→ What if feature is continuous

f_i	Output
2.1	Y
3.5	N
4.1	N
5.3	Y
6.8	Y

Step: ① Sort the features (f_i)

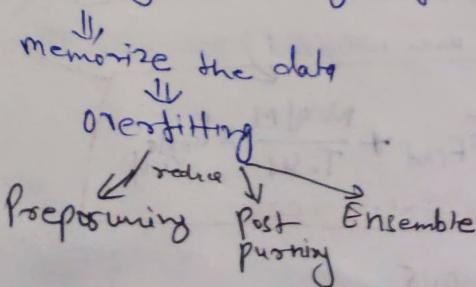
② treat each value / row as a cutoff / threshold and try to build decision tree & whichever threshold gives you the best information gain that will be used for splitting.

- ① Sort
- ②



★ DT is a Greedy Algorithm.

• It keeps on splitting until every leaf node is pure node



1st June DT & SVM

Agenda
↳ pickle
↳ Pre, Post Prnt
↳ Implement
↳ DT Regg
↳ Practical
↳ SVM.
↳ Naive Bayes
↳ Ensemble

Pickle → for use model

↳ import pickle
↳ Pickle.dump(model, open('model.pkl', 'wb'))
↳ Store it

↳ with open("model.pkl", 'rb') as file:

↳ loaded-model = pickle.load(file)

↳ loaded-model.predict()

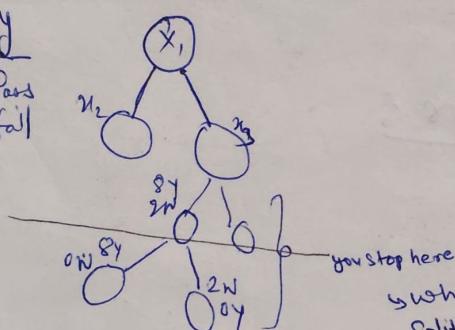
DT post pruning & Pre pruning.

DT is greedy Algorithms. → overfit the Model

In Generalize Training ↑ Test ↑

Training ↑ Test ↓

x_1	x_2	x_3	y
-	-	-	Pass
-	-	-	Fail



Less Node

Parameters to Create DT.

① Criteria

Gini

Entropy

② max_depth → maximum depth of tree

The minimum no of samples needed to be considered as a leaf node default is 1

Visualisation

Hyperparameters tuning using CV

Rising → cutting

③ Splitters → strategy (best, random) → best on information gain

④ min_samples_split → The minⁿ no of sample a node must contain to consider split.

⑤

① Post pruning (Cutting of dt after it is created)

Cutting

↳ Construct the entire dt until leaf node

↳ Prune the DT

↳ works well with small dataset

★ ② Pre-pruning

→ limit the growth of DT

→ Hyperparameter tuning to select best parameter.

Practical

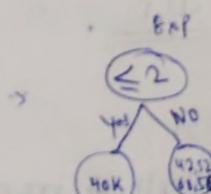
Decision Tree Regressor :

* Regression Classification
↳ $y \rightarrow$ Continuous, $y \rightarrow$ Cat

Calculate Variance ← Entropy

Variance $\leftarrow J(y)$
reduction impure

BP	Gender	Age	Salary
1.	yes	40	40K
2.	yes	42	42K
3	No	52	52K
4	No	60	60K
4.5	yes	58	58K



→ final aim
↳ reduce the Variance

We want feature value give highest

feature reduction

$$\text{Var} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

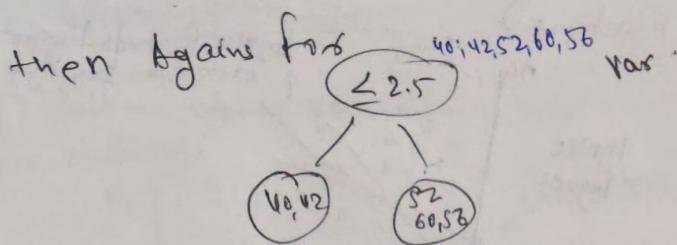
$$\begin{aligned} \text{Var} &= \frac{1}{5} \sum_{i=1}^5 ((40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (58-50)^2) \\ &= \frac{1}{5} (100+64+4+100+36) \Rightarrow 60.8 \end{aligned}$$

$$\text{Var}_{\text{Left Child}} = \frac{1}{2} (40-50)^2 \Rightarrow 100$$

$$\text{Var}_{\text{Right Child}} = \frac{1}{3} (42-50)^2 + (52-50)^2 + (60-50)^2 \Rightarrow 51$$

$$60.8 - \left(\frac{1}{5} (100+51) \right) \Rightarrow 0.004$$

| Variance Reduction $\Rightarrow \text{Var}_{\text{parent}} - \text{Var}_{\text{child combined}}$



$$\text{Var P}_{\text{var}} = 60.8$$

$$\text{Var left child} = \frac{1}{2} (100 + 64) \Rightarrow 82$$

$$\text{Var}_{\text{right child}} = \frac{1}{3} (4 + 36 + 10) \Rightarrow 46.66$$

$$\text{Variance reduction} \Rightarrow 60.8 - \left(\frac{2}{5} \times 82 + \frac{3}{5} \times 46.66 \right)$$

$$\Rightarrow 0.04$$

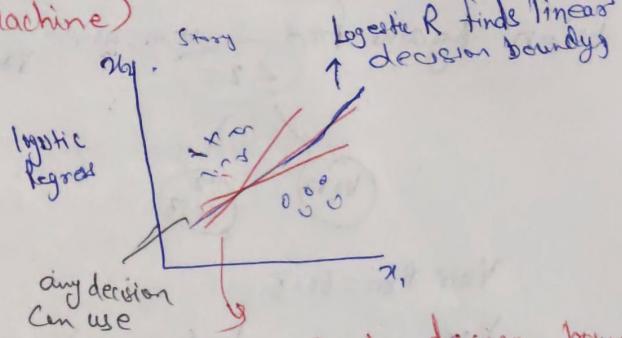
\leq Variance reduction for 2.5 is high so
use this is we 1st for split

considered 7.4.01

SV2 #

SVM (Support Vector Machine)

↳ SVC → classification
↳ SVR → regression



①

Any line can be decision boundary
I give u situation iff it segregates both the class

②

③

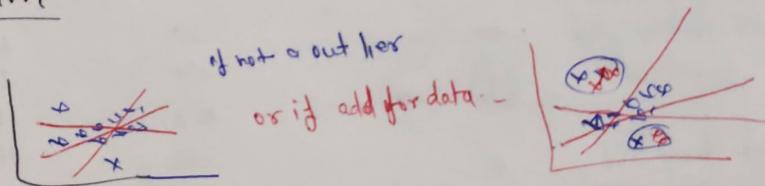
Which scenario is best (1), (2), or (3) yes
3rd one do here come SVM.

Scaling not need in DT

- Logistic Regression Model doesn't care about margin/space across two class.
- If the data point changes slightly Logreg. will give errors.

So for can we solve with DT → Ans Yes (And we have multiple option)
but SVM also solve this

• W.R.T Regression



↳ In SVM not impact by outliers but in Reg.LR it reflect.

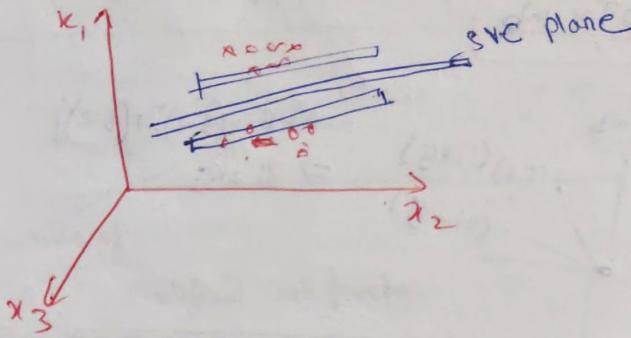
SVC

Marginal line 1 nearest to class 1
Marginal line 2 nearest to class 2
Margin → The gap b/w nearest dp of both class.
best fit line / decision boundary
And what is the best line → is between the Mline1 & line 2
To make marginal line for both class you need atleast 2 Support Vectors.

Support Vectors

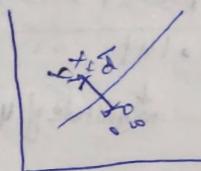
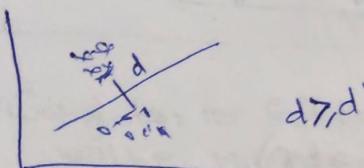
COMMON QUESTIONS

① Tell me about yourself?
→ for Multiclass



- There are no limit on SV.
- at least no of SV is 2
- SVC is also known as Margin Classifiers.
- SV use for

Now give you 2 scenarios



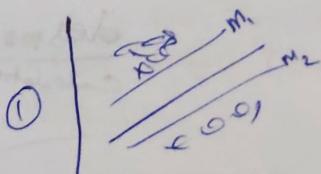
- Choose model whose marginal distance should maximum.

find max \leftarrow Step 1 : find out all possible classifier

~~Step 2~~ : find out margin of each classifier

~~Step 3~~ : Select the classifier which has maximum.

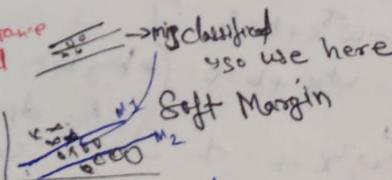
Maths behind
↳ SVM, C



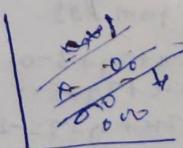
Hard Margin : None of the d.p.s are misclassified / overlapping



so Marginal distance will be small



- You are ready to misclassify some of the d.p.s in order to have good testing accuracy its a soft margin classifier.



$C = 3$ decide by CV

No of d.p.s so Sacrificed so that the best fit line is not overfitting.

→ eq. of line $y = mx + c, 0, + \theta_0 x$

$$ax + by + c = 0$$

$$by = -ax - c$$

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

$$\frac{c}{b} = M$$

$$y = Mx + C$$

In 3D $y = \theta_0 + \theta_1 x_1$

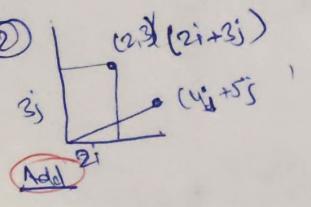
$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \rightarrow \text{hyperplane}$$

$$y = b + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \rightarrow w \text{ weights}$$

$$y = b + \sum_{i=1}^n \theta_i x_i \rightarrow \theta_i \text{ bias}$$

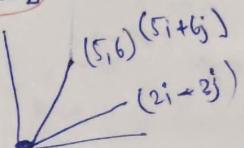
$$W = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}, X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \Rightarrow W^T \cdot X \Rightarrow$$

① $[w_0 w_1 w_2 w_3 w_4] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = w_0 x_0 + w_1 x_1 + \dots + w_4 x_4$ use in SVM like this
 $\Rightarrow y = \vec{w} \cdot \vec{x} + b$

② 
 Add $\vec{A} + \vec{B} = (5i+3j) + (2i+2j) = (7i+5j)$

③ Vector Sub

$$\vec{A} - \vec{B}$$



$$\vec{A} - \vec{B} = (5i+6j) - (2i+3j) = 3i+3j$$

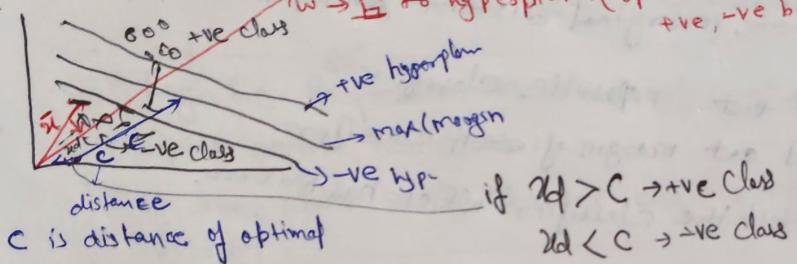
④ Dot Product of two vectors



$w \cdot u \rightarrow$ dot product means projection of u on w
 $\Rightarrow \text{proj}_{w \perp} u = |u| \cos \theta = |u|$

\downarrow magnitude of w \uparrow magnitude of u

$\vec{w} \rightarrow$ hyperplane (optimal classifier)
 +ve, -ve both



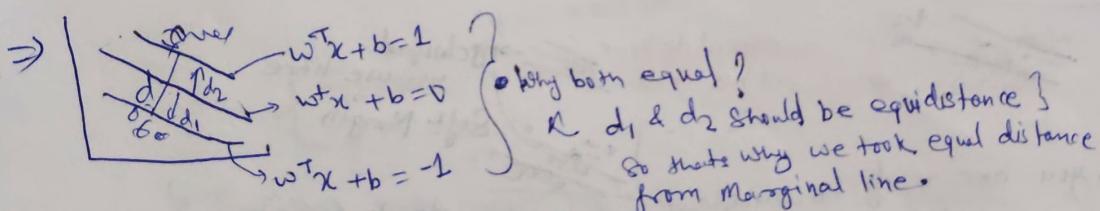
we know
 projection is dot product
 $\vec{w} \cdot \vec{x} = c$

⑤ projection on w

$$w^T \vec{x} + b = 0$$

$$y = \begin{cases} +ve & \text{if } \vec{x} \cdot \vec{w} + b > 0 \\ -ve & \text{if } \vec{x} \cdot \vec{w} + b \leq 0 \end{cases}$$

desmos
 calculator



why only 1? not (2, 3, 4, 5)

because the 1 is not change from LHS

$$\begin{aligned} y_1(\vec{w} \cdot \vec{x}_1 + b) &\geq 1 \\ 2x_1 + 3x_2 + 3 &\geq 1 \\ 2x_1 + 2x_2 + 3 &= 1 \\ 2x_1 + 3x_2 + 2 &= -1 \end{aligned}$$

$$y_2(\vec{w} \cdot \vec{x}_2 + b) \geq 1$$

$$2x_1 + 3x_2 + 3 \geq 1$$

$$2x_1 + 3x_2 + 2 = 5$$

$$2x_1 + 3x_2 + 2 = 5$$

$$2x_1 + 3x_2 + 2 = 5$$

So minimised such that

$$y_1(\vec{w} \cdot \vec{x}_1 + b) \geq 1$$

$$\text{diff } \vec{x}_2 - \vec{x}_1$$

Projection on w

$\vec{x}_2 - \vec{x}_1$

So get shortest distance b/w the marginal plane, we need projection of $\vec{x}_2 - \vec{x}_1$ on unit vector ' w ' to get d

$$d = \frac{(\vec{x}_2 - \vec{x}_1) \cdot \vec{w}}{\|\vec{w}\|}$$

$$\frac{\vec{x}_2 \cdot \vec{w} - \vec{x}_1 \cdot \vec{w}}{\|\vec{w}\|}$$

where x_1 & x_2 are support vectors, they lie on marginal hyperplane
 for +ve class $y_1 = 1 \Rightarrow \vec{x}_1(\vec{w} \cdot \vec{x}_1 + b) = 1 \Rightarrow \vec{w} \cdot \vec{x}_1 = 1 - b$

Putting ① & ②

$$\frac{(1-b) - (b-1)}{\|w\|} - \frac{2}{\|w\|} = d$$

for -ve class

$$y_i = -1$$

$$-1(wx_2 + b) = -1$$

$$wx_2 = -b - 1 \quad \text{--- (2)}$$

So Cost function

$$\text{maximize}_{w,b} \frac{2}{\|w\|}$$

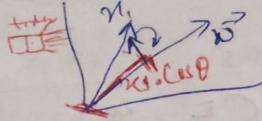
$$\text{such that } y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1$$

magnitude of coefficient

Modified Cost function

$$\text{minimize}_{w,b} \frac{\|w\|}{2} \text{ by varying } w \& b$$

$$y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1$$



→ Cost function for Soft Margin

$$\text{will be minimize}_{w,b} \frac{\|w\|}{2}$$

$\underbrace{\|w\|}_{\text{SVM error}}$ Collected

bring loss
Misclassification error / Margin error

$$+ C \sum_{i=1}^n \epsilon_i$$

$\underbrace{\epsilon_i}_{\substack{\text{no. of} \\ \text{misclassified} \\ \text{d.p.}}}$ Zeta

$\star \star$

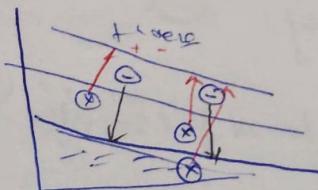
It is the distance of all misclassified d.p.s due to new marginal plane

$$C = \frac{1}{\lambda}$$

in sklearn

$$\sum \epsilon_i$$

$$\epsilon_i = zeta$$



$$\text{Such that } y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1$$

+ve d.p. from their own plan
& -ve d.p. distance from their own plan

→ So summation of all misclassified d.p. from other plan

CS of AI work with and explain it in simple words

→ first go into what problem they solve → few simple lines
Implementation gives scope of future increment
- the scope of reinforcement

For CS • GPT-3 for many models (GPT-3, DALL-E, etc.)

Page 23

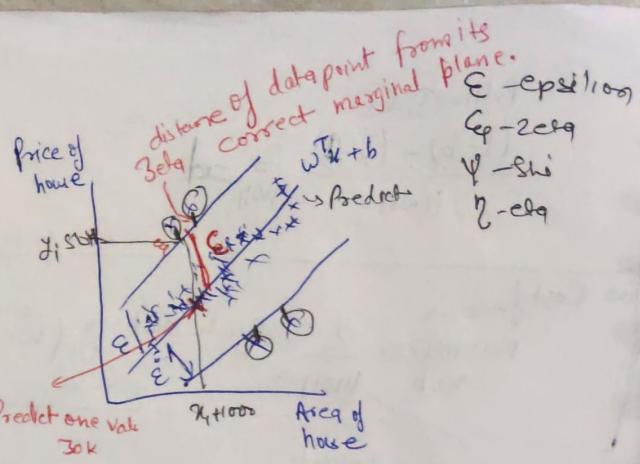
Support Vector Regressor

SVR

→ We want all dp's in Surrounding of best fit line.

↓

$$w^T x + b + \epsilon \text{ to } w^T x + b - \epsilon$$



for regression

$$CF = \underset{w, b}{\text{Minimize}} \frac{\|w\|}{2} + C \sum_{i=1}^n (\epsilon_i + \epsilon)$$

Constraints

$$|y_i - y_{\text{pred}}| \leq \epsilon + \epsilon_p$$

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

$$\text{Yact} - y_{\text{pred}} = \text{Error} = \epsilon + \epsilon_p$$

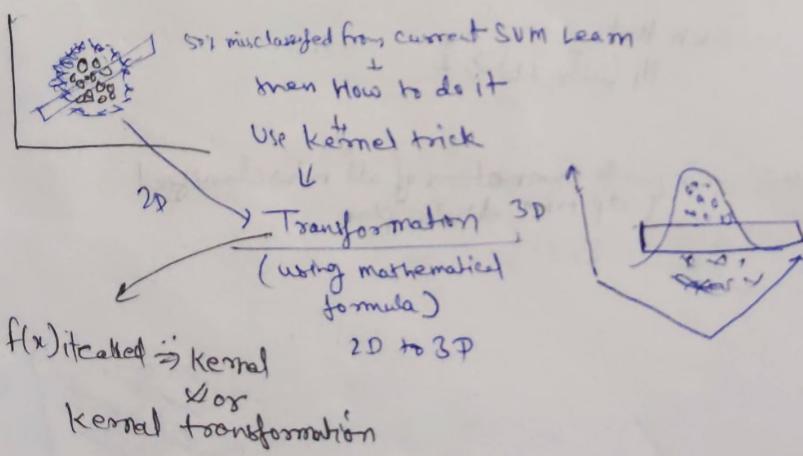
minimize this

→ it done in backed by gradient descent

- All the dp's in b/w of the marginal plan.

ϵ_p - is the distance of dp from its correct marginal plane.

SVM → Kernel trick



eg	X	y
	1	yes
	2	no
	3	no
	4	yes
	5	no
		y-

6 Idea → To change the data from 1D to 2D or increase the dimension

→ Kernel trick → ~~way~~ → you are not sending the data in higher dimension ~~using~~ if you are using mathematical transformation to achieve it.

① Polynomial Kernel:

$$f(x_1, x_2) = (x_1^T \cdot x_2 + c)^d$$

d is degree of polygon
constant

$x_1, x_2 \rightarrow$ features

$$x_1^T \cdot x_2 = [x_1^T] \cdot [x_2]$$

$$[x_1^T]^T \cdot [x_2]$$

Electric Vehicle Data Analysis

(1) Problem Statement

or orally

(2) radial basis feature (rbf)

- Create non-linear combination of features to bring feature higher dimension.

$$* f(x_1, x_2) = -\frac{\|x_1 - x_2\|^2}{\sigma^2}$$

$\|x_1 - x_2\| \rightarrow$ Euclidean distance b/w x_1 & x_2

$\sigma \rightarrow$ Variance.

(3) Sigmoid kernel

$$f(u) = \frac{1}{1 + e^{-u}}$$

- * How to choose right kernel

↓
Hyperparameters tuning

Which one choose

① distinguishable \rightarrow linear kernel



②
we \rightarrow Rbf Kernel

③
 \rightarrow use polynomial kernel

④ Sigmoid. \rightarrow most for binary

* Bessel kernel $\rightarrow f(x, y) = J_{(v+1)} \frac{(\sigma \|x - y\|)}{\|x - y\| - \sigma(v+1)}$

* ANOVA kernel $\rightarrow f(x, y) = \sum_{n=1}^N \exp(-\sigma (x_n - y_n)^2)$

8 June Naive Bayes & SVM

2 hr SVM implement

Naive Bayes (only for classification)

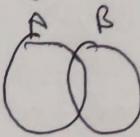
- Prob: how likely something is to happen $P = \frac{\text{No of favourable outcome}}{\text{Total no of outcome}}$
- $P(M) = \frac{1}{2}$

- Independent Events & dependent events.
 - \downarrow outcome of an event is dependent on another event.
 - outcome of one event is not affected by another Events

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

→ Conditional probability

Bayes theorem



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

$$\text{So } P(A|B) * P(B) = P(B|A) * P(A)$$

Bayes theorem

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

$P(A|B)$ = Prob of event A given

$P(A)$ = Prob of event A, B has occurred

$P(B)$ = " " event B,

$P(B|A)$ = Prob of event B given A had occurred.

example

$$P(y|x_1, x_2, x_3) = \frac{P(y) * P(x_1|y) * P(x_2|y) * P(x_3|y)}{P(x_1, x_2, x_3)}$$

$$P(y|x_1, x_2, x_3) = \frac{P(y) * P(x_1|y) * P(x_2|y) * P(x_3|y)}{P(x_1) * P(x_2) * P(x_3)}$$

y has two possibility

$$P_{\text{yes}} = P(y|\text{yes}) = \frac{P(\text{yes}) * P(x_1|\text{yes}) * P(x_2|\text{yes}) * P(x_3|\text{yes})}{P(x_1) * P(x_2) * P(x_3) - \text{constant}}$$

$$P_{\text{no}} = P(y|\text{no}) = \frac{P(\text{no}) * P(x_1|\text{no}) * P(x_2|\text{no}) * P(x_3|\text{no})}{P(x_1) * P(x_2) * P(x_3) - \text{constant}}$$

NB $\Rightarrow P_{\text{yes}} = P_{\text{no}}$

→ NB
 → Ensemble
 ↳ bagging / Boosting
 → KNN
 ↳ PCA
 ↳ unsupervised learning
 ↳ Anomaly detection

* for multi class Classification

$$P(c_k|x) = \frac{P(x|c_k) \cdot P(c_k)}{P(x)}$$

$$P(c_1|u_1, u_2, u_3) = \frac{P(c_1) + P(u_1|c_1) \cdot P(u_2|c_1) \cdot P(u_3|c_1)}{P(u_1) P(u_2) P(u_3)}$$

$$P(c_2|x_1, u_2, u_3) = \frac{P(c_2) + P(u_1|c_2) \cdot P(u_2|c_2) + P(u_3|c_2)}{P(u_1) P(u_2) P(u_3)}$$

$$\text{Max Probabilities} \\ \text{Predicted class}$$

Example

Outlook Data Sets

Q. 10

outlook	yes	No	$P(E y)$	$P(E N)$
Sunny	2	3	2/9	3/5
Overcast	4	0	4/9	0
Rain	3	2	3/9	2/5
	9	5		

for

Temp	Y	N	$P(E y)$	$P(E N)$	$P(y N)$
Hot	2	3	2/9	2/5	9/14
Mild	4	1	4/9	2/5	4/14
Cold	3	2	3/9	1/5	3/14
	9	5			1/5

$$P(y/\text{sunny, hot}) = P(y) \cdot P(s|y) \cdot P(H|y)$$

$$\therefore \frac{9}{14} \cdot \frac{2}{9} \cdot \frac{2}{9} = \frac{2}{63} \approx 0.031$$

$$P(N/S, H) = P(N) \cdot P(s|N) + P(H|N)$$

$$= \frac{5}{14} \cdot \frac{3}{4} \times \frac{2}{5} = \frac{3}{35} = 0.085$$

→ No Normalized the data because both has no number of equal datapoint
 So → Normalized the prob to get in percentage for $P(y/\text{sunny, hot})$

$$P(y/\text{sunny, hot}) = \frac{0.031}{0.031 + 0.085} = 0.27 = 27\%$$

$$P(N/\text{sunny, hot}) = \frac{0.085}{0.031 + 0.085} = 0.73 = 73\% \quad \text{Prediction} \Rightarrow \text{The person will not go to play football}$$

Why NB called → because here we use Bayes theorem

Why Naive → It assumes that features are independent

→ No Multi-collinearity

Adv:

- based on probability
- fast Algorithm (only calculation)
- work well with less training data
- performance is good

→ go NLP use

NB → because
of independence

Dis Adv

- Assumes features are independent
- It doesn't work for Regression Problem

⇒ Variant of NB

① Bernoulli NB

data: $\begin{bmatrix} f_1 & f_2 & f_3 \\ M & F & M \end{bmatrix}$ $\begin{bmatrix} y_1 & y_2 & y_3 \\ 0 & 1 & 0 \end{bmatrix}$ $\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix}$ $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$

→ Independent Variable follow Bernoulli,
Bernoulli NB.

→ it is called - Sparse Matrix (where most 0 and 1)
also it leads Overfitting
(because it not proper algos
it only memorize thing)

② Multinomial NB

when inputs are → text data

e.g.: mail classification → spam / ham

e.g.: I am going to play football

you think play & football dependent
Answer is No why because

I also possibility to playing played
any possibility.

(but also include context)

e.g.: Email body Output

you got a discount spam

millions of Rs spam

offer lets ham

↓

Change (Numerical Value) → NLP technique
Vectors

- ① Tf - Idf
- ② BOW
- ③ Word2Vec

Multinomial and Bernoulli is use for textual data in NLP

③ Gaussian NB

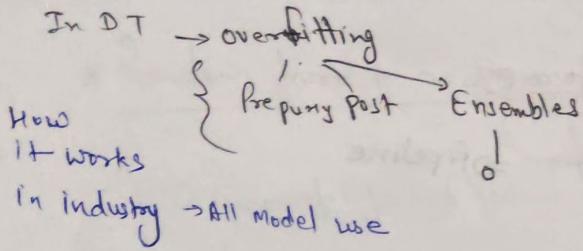
→ If features are continuous & follow Normal Distribution

- if some feature use Bernoulli distribution or some Normal Distribution use that Variant which ever feature is higher.

1 EC-9 Random Forest & Bagging 9 June

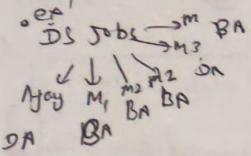
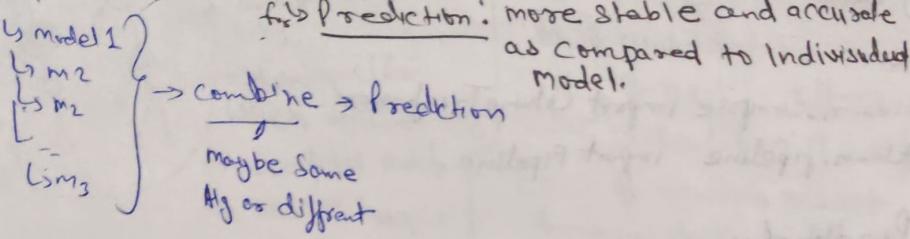
Which Model use

- Linear data: MLR, Logit
- Non linear \rightarrow DT, SVM



Ensembles and Its techniques

data \rightarrow Combine multiple Models



- One person might give you wrong advice
- multiple mentors \rightarrow chance of getting wrong less/minimizing

of same Algo

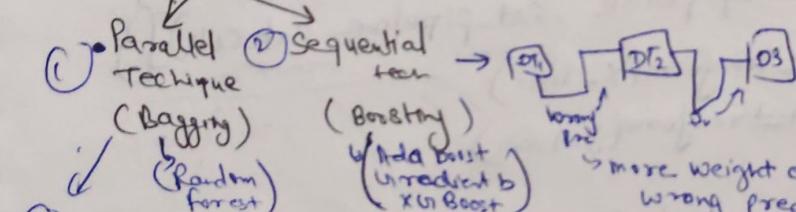
- ① DT₁ \rightarrow max depth = 5
- DT₂ = 7, 8 = 10
- ... = 12

different algos also

- logistic Reg
- SVC
- DTC(5)
- DT(5)
- NBC

NOTATION
• Delite
• ZS

- So Ensemble \rightarrow No necessarily only one type of Algo



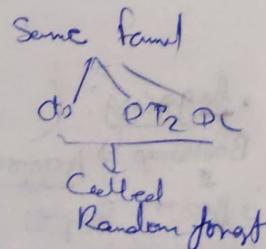
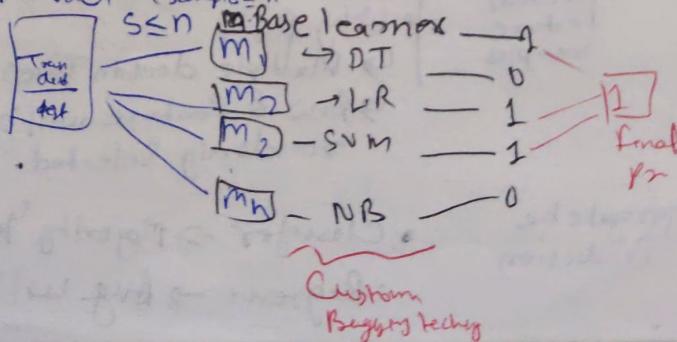
parallel technique of Ensembles

data $\xrightarrow{\text{train}} \text{DT}_1, \text{DT}_2, \dots, \text{DT}_n$ Here all the models here are built parallelly and independent of each others.

② S

* Bagging (Custom Bagging):

- Parallel Model ($n \text{ samples} \leq n'$)



→ Custom Bagging class
use Pipeline & Column transformer

② A Sequence of data transformation \Rightarrow Pipeline

Set of all feature

Column Transform (group all the pipeline steps
for each of the column)

- handling missing Value
- data encoding
- Feature scale

→ from sklearn.compose import ColumnTransformer

→ from sklearn.pipeline import Pipeline

⇒ Num-pipeline = Pipeline(steps = [("imputation", SimpleImputer(strategy = "median")),
("Scaling", StandardScaler)])

⇒ Cat-pipeline = Pipeline(steps = [("imputation", SimpleImputer(strategy
= "most_frequent"))])

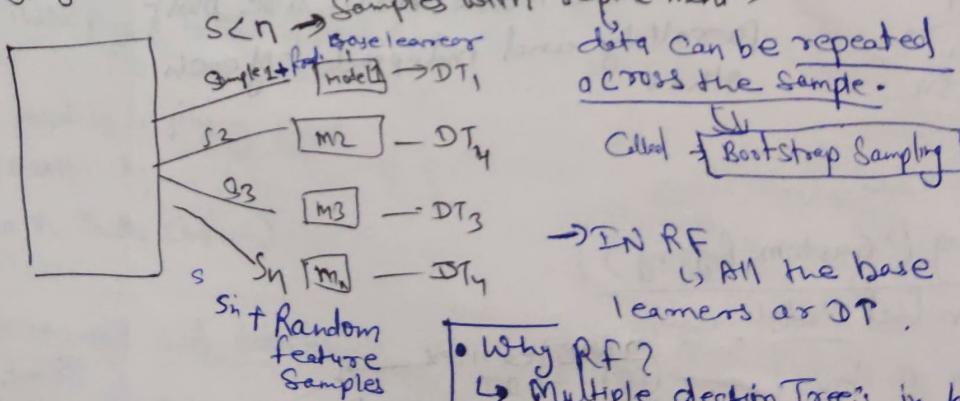
⇒ ColumnTransformers([("Num-pipeline", num_pipeline, num_cols),
("Cat-pipeline", cat_pipeline, cat_cols)])

X-train = preprocessor.fit_transform(X-train)

X-test = preprocessor.fit_transform(y-test)

Random forest classifier and regressor

(Bagging technique)



5ⁿ + Random
feature
Samples

→ IN RF
↳ All the base
learners are DT.

• Why RF?

↳ Multiple decision Trees in parallel
↳ Row & feature will be randomly Selected.

- Classifier \rightarrow Majority 'K' the prediction
- Regressor \rightarrow Avg. will be the pred.

• Bagging:
Bootstrap \rightarrow Aggregating

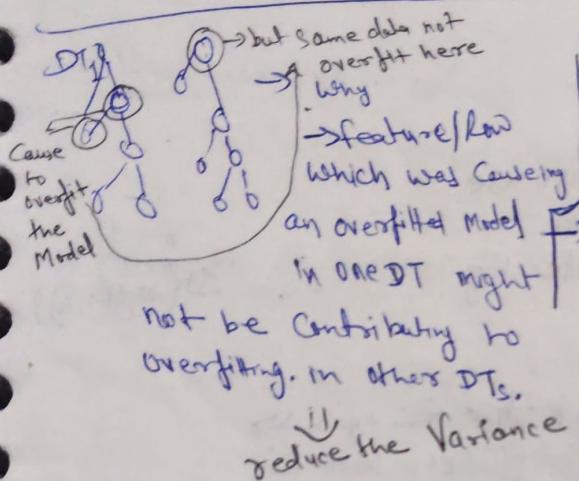
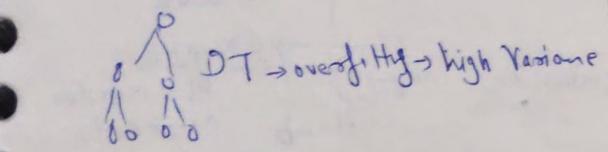
different Samples
with replacement

Aggregate the
Predictions

Bootstrap Sampling

* Random forest

* Random forest reduces the Variance



Random forest

↳ random sampling of rows & features

→ data is splitted into small chunks (randomly)

→ higher variance is being reduced

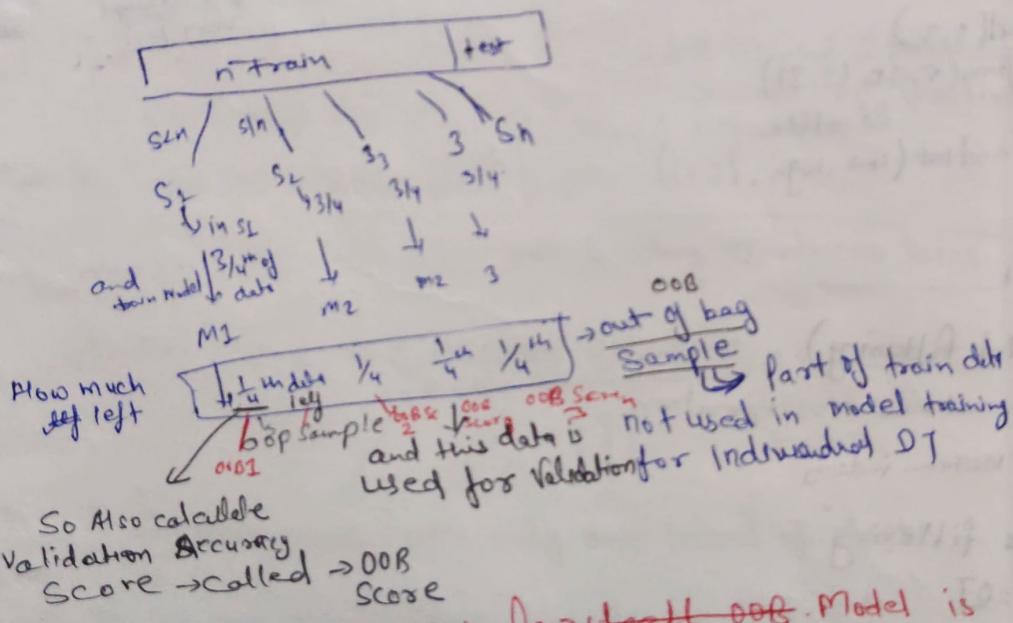
→ Reduce Bias via boosting

↳ Due to random subset of rows and features each of the subset is different representation of itself.
Each of random subset trains a diff. model.

june 14 data to date

⇒ OOB Score :

out of bag



- In training itself oob is low if self oob. Model is not performing well on train data.

14 June toolkit

→ Use Numpy in Real World

→ reshape

b → go in depth

- 3 different way to create Array
- if you have list → of 2D [C]

→ for loop vs Vectorized operation, max operation?

>> list_2D = [[1, 2, 3], [4, 5, 6]]

>> arr = np.array(list_2D)

>> type(list_2D)

>> type(arr) or arr.dtype → we know data type

• np.full() → np.full((n, m), Val, dtype="float")

np.ones((n, m)) * 5

→ Concept of vectorization use
also use broadcasting

Use Valid broadcasting

list([day(1, 10)])

⑧ Generate an array of 100 evenly spaced

np.linspace(1, 100, 100)

⑨ np.asarray & np.array, np.asarrayarray

→ If copy take same & then again we

array to find off start space left ho

→ No copying array when we use adv obj

when we use

adv obj & also maintain their functionality (Called Sub Clasering)

asarray Usefull हीजी

Properties

np.expand_dims() uses np.asarrayarray

directly not use by

you always
very less.

⑩ np.random.rand(3, 3)

>> np.random.uniform(5, 20, (3, 3))

>> np.random.randint(10, high, (5, 6))

→ Compare all

arr[arr > 2]

→ This type of thing is called filtering

broadcasting → element wise → == operation

>> arr[1, 2] == 0

← Boolean indexing

Now then Use filtering

arr[arr > 2 == 0]

>> np.where

↙ Not use index
True False False
→ ETC

1D, 1D

[True] * 3 + [False] * 2

83 → 3D → np.argmax → np.random.randint(1, 10, (3, 3, 3))

argmax → index return

argmax() → np.argmax(arr=3D, axis=2)

Column

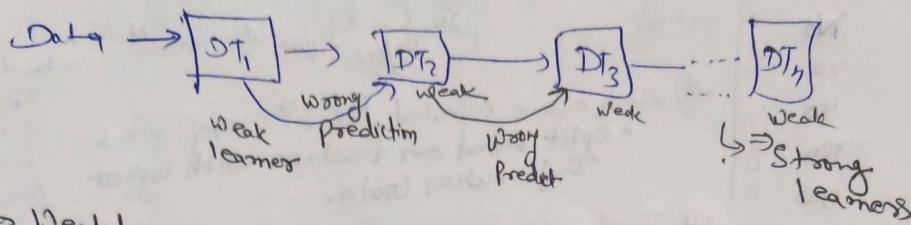
axis=1 → Row level

axis=0 → ?? find

→ Building
In n
Classification
Models
→ May
Prob

⇒ Boosting

- ① Ada Boost
- ② Gradient Boosting
- ③ XG Boosting

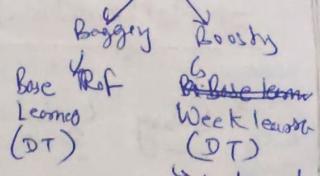


→ Weak learners:-

Model is very simple and not properly from data.

- Principle of Boosting :- Building a first model on training data and build a second model to rectify the errors present in this first Model.

Ensembles



What do you mean by base learner and weak learners?

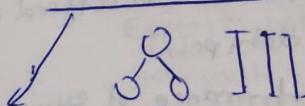
- # AdaBoost (Adaptive Boosting) (Yael Freund & Robert Schapire, 1995)
- Concept of Boosting
 - work both for Classification & Regression
- you know the Boosting work flow

before SVM we

→ Bagging → reduce variance, but boosting reduces bias

In Ada's concept → Decision Stumps :- 1 parent node &

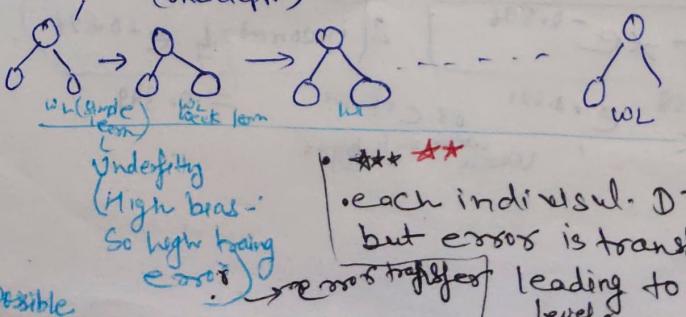
2 child nodes
that is.



A Decision tree with only one level of depth

↓
Here why we use binary not a fully grow tree (one depth)

In full grow tree we can't boost, not possible



→ to Adaboost Solve the problem of high bias.

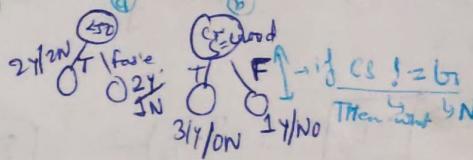
Mathematical Explanation of AdaBoost

Create decision tree Stump & predict
so based on it
This predict wrong because Actual if yes &

	Salary	Credit Score	Approved	Wt
0.058	$\leq 50K$	B	No	1/7
0.058	$\leq 50K$	G	Yes	1/7
0.058	$\leq 50K$	G	Yes	1/7
0.058	$> 50K$	B	No	1/7
0.058	$> 50K$	G	Yes	1/7
0.058	$> 50K$	N	Yes	1/7
0.058	$\leq 50K$	N	No	1/7

Step 1

Create decision tree Stump & predict



- ④ Which one is selected for 1st Stump ④ or ⑥
- Split based on feature with highest Information Gain.

All equal Weights.

Step 2 → Assign equal weights to all the datapoints

$$\sum \frac{1}{7}$$

different in wt
so we normalize
the wt

Normalized wt = $\frac{\text{Update wt}}{\text{Total updated wt}}$

$$0.058/0.69$$

$$0.058/0.07$$

$$0.500/0.03$$

$$0.500/0.03$$

Step 3 → Total error & performance of Stump.

$$\text{TE} = \sum \frac{1}{7}$$

(Sum of weights of misclassified dpt)

So what is the performance of Model?

Performance of Stumps =

$$\frac{1}{2} \ln \left[\frac{1 - \text{TE}}{\text{TE}} \right]$$

base e

Boosting Prediction function

$$f = \alpha_1 M_1 + \alpha_2 M_2 + \dots + \alpha_n M_n$$

$$0.896$$

α_i = weight of Model
 α_i = $e^{-\text{Performance of Model}}$
 α_i = $e^{-\text{Performance of Model}}$

Performance of individual Stump (Model) will be considered as weight for that specific model.

$$\Rightarrow \frac{1}{2} \ln \left[\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right]$$

$$\Rightarrow \approx 0.896$$

Steps : So Next step is to give more priority to wrong prediction
Step 4 Give more weights to incorrect classified dpt's & lesser weights to correct data point

For correct class wt_{new} = $w_t \times e^{-\text{Performance of Stump}}$

for incorrect classif dpt
wt_{new} = $w_t \times e^{+\text{Performance of Stump}}$

$$\text{So } \boxed{\text{Correct} \rightarrow \frac{1}{7} \times e^{-0.896}} \\ = 0.058 e^{-0.896}$$

$$\boxed{\text{incorrect} \rightarrow \frac{1}{7} e^{+0.896}} \\ = 0.349$$

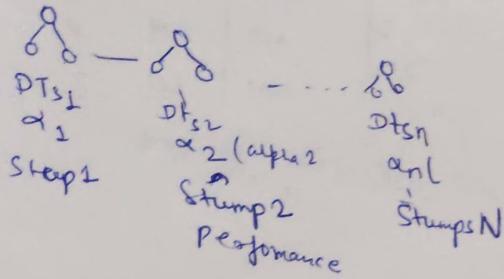
which weight is bigger

Forwards transfer to next DTs to Reduces Bias

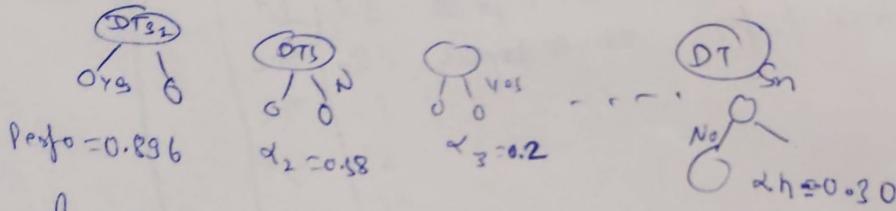
Node	Bin Assignment	Random Hos (0 to 0.98)
0.02	0 - 0.088	
0.03	0.08 - 0.16	
0.04	0.16 - 0.24	
0.05	0.24 - 0.32	
0.06	0.32 - 0.40	
	0.40 - 0.48	
	0.48 - 0.56	
	0.56 - 0.64	
	0.64 - 0.72	
	0.72 - 0.80	
	0.80 - 0.88	
	0.88 - 0.96	
	0.96 - 1.00	

Note
The row with higher bin range due to misclassified will be repeated.

It means more priority on misclassified.



Final Prediction



$$\begin{aligned}
 f &= \alpha_1(m_1) + \alpha_2(m_2) + \alpha_3(m_3) + \dots + \alpha_n(m_n) \\
 &\Rightarrow 0.896(\text{Yes}) + 0.58(\text{No} + 0.2\text{Yes}) + \dots + 0.30(\text{No}) \\
 &\Rightarrow 1.18(\text{Yes}) + 0.25(\text{No})
 \end{aligned}$$

So final prediction $\Rightarrow \text{Yes}$ ✓

When to stop tree \rightarrow by define the no of tree)

AdaBoost Regressor: All the steps will be same except instead of Information gain, you will use Variance reduction to Select the decision stump.

6

total \$1000.00
Capital 500.00
Expenses 100.00
20% commission 100.00
Contingency expense 100.00
Total 1000.00
Paid on Hastings for 1000.00
not remitted to business traps
or whatever amount did Mrs. Mc-
Gowen have on hand?

Gradient Boosting

↳ Regression
↳ Classification

Exp	Degree	Act Salary	\hat{y} Predict
2	Btech	50K	75
3	M	70K	75
4	M	80K	75
5	Phd	100K	75

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

Step 1: Create a base model [75]

Step 2: Compute the Residual / error from y_{actual} .

Step 3: Make a dt considering input as the independent Variable (Exp, degree) and target variable as R_1 .



Using this decision tree make prediction for input variables.

	Act-Salary	\hat{y} Predict	Residual Prediction
R_1	R_2	as R_1 is taken as target variable.	
-25	-23		
-5	-3		
5	3		
25	20		

1st Gradient Called

Gradient Boosting Classifier

• (1) → Base prediction eg: 0.5

• Instead of Variance Reduction Using Information Gain in Decision tree.

1920 - 1921

1921 - 1922

1922 - 1923

1923 - 1924

1924 - 1925

1925 - 1926

1926 - 1927

1927 - 1928

1928 - 1929

1929 - 1930

1930 - 1931

1931 - 1932

1932 - 1933

1933 - 1934

1934 - 1935

1935 - 1936

1936 - 1937

1937 - 1938

1938 - 1939

1939 - 1940

1940 - 1941

1941 - 1942

1942 - 1943

1943 - 1944

1944 - 1945

1945 - 1946

1946 - 1947

1947 - 1948

1948 - 1949

1949 - 1950

1950 - 1951

1951 - 1952

1952 - 1953

1953 - 1954

1954 - 1955

1955 - 1956

1956 - 1957

1957 - 1958

1958 - 1959

1959 - 1960

1960 - 1961

1961 - 1962

refined and to standard
in the following year. It
was adopted with considerable reluctance by the students
and teachers in

XGBoost (Extreme Gradient Boosting)

Classifier / Regressor

Step 1: Create a Base Model

XGBoost Classifier

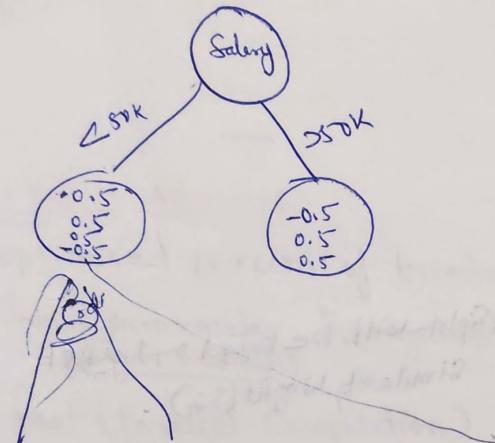
Salary	Credit/Score	Appraiser	Base model Residual	
			P _r	R _r
<=50K	B	0	0.5	0.5
<=50K	G	1	0.5	0.5
<=50K	G	0	0.5	0.5
>50K	B	0	0.5	-0.5
>50K	G	1	0.5	0.5
>50K	N	1	0.5	0.5
<=50K	N	0	0.5	-0.5

Step 2: Construct a decision tree

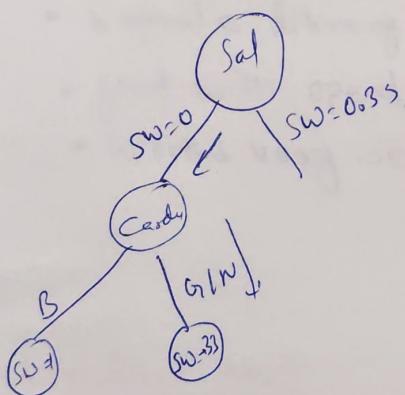
Step 3: Calculate Similarity weight

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

Base Model Residual
 $\boxed{0.5} \rightarrow (0.5, -0.5, 0.5, -0.5, 0.5, 0.5)$



$$\Rightarrow SW_{Root} = \frac{[0.5 + 0.5 + 0.5 + -0.5 + \\ 0.5 + (-0.5)]^2}{0.5 + (1-0.5)^2} = 0.5 \\ = 0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5) + \\ 0.5(-0.5) + 0.5(1.5) + 0.5(1-0.5) \\ \Rightarrow \frac{(0.5)^2}{0.25 \times 7} = \frac{0.25}{1.75} = 0.14$$



$$SW_{Left\ Child} = (0.5 + 0.5 + 0.5 - 0.5)^2 = 0 \\ SW_{Right\ Child} = \frac{(-0.5 + 0.5 + 0.5)^2}{0.5 \times 3} = \frac{0.25}{0.75} = 0.33 \\ SW_{Gain} = (SW_{LC} + SW_{RC}) - Root \\ = (0 + 0.33) - 0.14 \\ = 0.19$$

$$\rightarrow O/P = \sigma(\text{base} + \alpha \cdot SW) \Rightarrow \sigma(0 + 0.19 \cdot 0.33)$$

* Whichever feature has the highest similarity weight gain, the split will happen on that feature.

XGBoost Classifies $\rightarrow O/P$

Output

$$\rightarrow \sigma(\text{Base Model} + \alpha \cdot (SW))$$

$$\frac{1}{1+e^{-x}} \quad \frac{\log P}{P-FP}$$

$$\Rightarrow 0.508 \approx$$

XGBoost Regressor

Exp	Drop	Salary	Base Prediction	Residual
2	yes	40	51	-11
2.5	yes	42	51	-9
3	No	52	51	1
4	No	60	51	9
4.5	Yes	62	51	11

$$\text{Base Prediction} = \frac{(40+42+52+60+62)}{5} \\ \Rightarrow 51.2 \Rightarrow 51$$

$$SW_{root} = \frac{(-11+9+1-9+11)^2}{5+1} = \frac{1}{6}$$

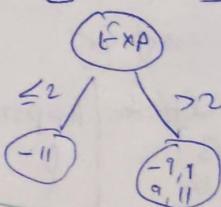
$$SW_{leftchild} = \frac{(-11)^2}{1+1} = \frac{121}{2} = 60.5$$

$$SW_{rightchild} = \frac{(-9+1+9+11)^2}{4+1} = \frac{144}{5} = 28.8$$

$$SW_{gain} = (SW_{child}) - SW_{root} \\ = (60.5 + 28.8) - 0.166 \\ = 89.133$$

here calculate for 2

$$R_i = [-11, -9, 1, 9, 11]$$



(Similarity weight)
* SW gain_i = $\frac{(\sum R_i)^2}{\sum P_i(1-P_i) + \lambda}$

P_i = Probability

* SW_{reg} = $\frac{(\sum R_i)^2}{\text{No of Residuals} + \lambda}$

d = 0

* Split will be based on highest Similarity Weight (SW).

~~Same for the Experience feature calculate.~~

No calculate with Exp. 2.5

$$\begin{aligned} & \text{EXP} \\ & \leq 2.5 \quad > 2.5 \\ & (-11, -9) \quad (1, 9, 11) \\ & \frac{(-11-9)^2}{2+1} = \frac{400}{3} = 133.3 \end{aligned} \rightarrow SW_{root} = 0.166$$

$$\rightarrow \frac{(1+9+11)^2}{3+1} = \frac{211}{4} = 110.25$$

$$SW_{gain} \Rightarrow 133.3 + 110.25 - 0.166$$

≈ 244

so highest of SW is here
So use this to split.

$$SW = \frac{\sum_{\text{positive}}}{\sum_{\text{positive}}(1-p_s) + \sum_{\text{negative}}} \rightarrow \text{hyperparameter}$$

$\uparrow \uparrow SW \downarrow \downarrow$

Chances are there
Sagain can become
negative

\downarrow
No Sagain \Rightarrow Stop Splitting

* $\leq p_s(-p_s) \Rightarrow$ Cover Value

\hookrightarrow any SW less $p_s(1-p_s)$, we stop splitting

SDD:

\Rightarrow XG Boost + Advantage

- Optimised version of梯度 boosting
- Inbuilt parameter F (regularization: λ in gradient)
- Fast (Parallel Computation)
- Separate library itself
- Good with speed & accuracy
- Works very well with large dataset.

Gradient Boost

- Computationally expensive
- prone to overfitting (better than AdaBoost)

* what is significance
of lambda in XGB
 \downarrow
what is the mean

- regularization
- Parameter
- Remember

- LR/MLR \rightarrow Ridge, lasso
- Log $\rightarrow L_1, L_2$
- SVM $\rightarrow C$
- SVR $\rightarrow C$
- DTC/DTR $\rightarrow CCP_A$
- Bagging $\rightarrow \alpha$

• Boosted Random forest

- Black box & white box model
- Model which has explainability power
- $y = B_0 + B_1x_1 + B_2x_2$
- White Box Model

\rightarrow option time complexity

$N \rightarrow$ No of training sample

$M \rightarrow$ No of features

$T \rightarrow$ no of tree

$d \rightarrow$ max depth

\bullet $O(T \cdot NM \log N)$

\bullet Gradient Boosting by A
 $O(T \cdot NM \log N)$

in practice this is very very fast

- Separate library
- Parallel Computation

\bullet Ada Boost $\approx N, M, T$
 $O(T \cdot Nm)$

Ada Boost

- for small dataset
- Sensitive to noise & outliers (bec. 1 depth used)

length

length of the organ or

length of (q-1) + 1

(length of q-1) + 1 + 0

length

length of each part

length of each horizontal

horizontal part of the organ

(length of each part)

length of each part + 1

length of each part + 1

length of each part + 1

length of each part

length of each part

(length of each part) + 1

length of each part

length of each part

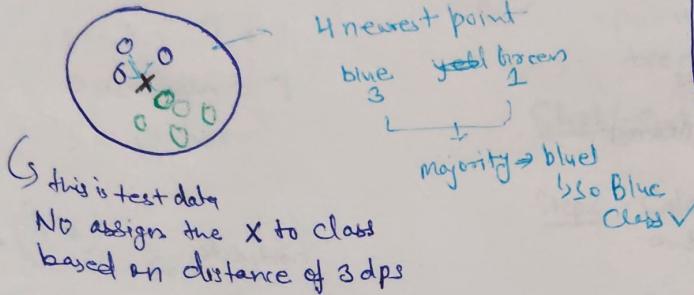
22 June KNN

K-Nearest Neighbour

- ① KNN Classifier
- ② KNN Regressor

KNN Classifier

Start with example.



Scenario 2
instead of 4 nearest dp
take 7 nearest dp.
→ turn the → green
K = nearest data point

• So As K changes
the class of new
dp might also change

↓
K is a hyperparameter

-
- Sort the distances
 - $\{ \begin{matrix} 0-1 \\ 0-2 \\ 0-3 \\ 0-4 \\ 0-7 \\ 0-9 \\ 0-9 \end{matrix} \}$ 4 nearest pts

binary class → K as odd value
3, 5, 7, ...

Algorithms

- ① Plot the dp in n-d space
- ② Initialize the K-Value
(No of neighbours you want to consider)
 $K \in 1 \rightarrow \infty$ (generally $K > 3$)
- Calculate the distance of new dpts wrt to all dpts
- sort the distance
- based on K find class of that K nearest dp's
- ③ find the mode of the class
- ④ Assign the class.

Distance

① Euclidean distance

$$(x_1, y_1) \quad (x_2, y_2)$$

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

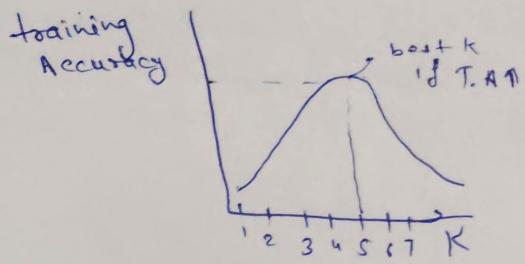
② Manhattan distance

$$(x_1, y_1) \quad (x_2, y_2)$$

$$d_1 \quad d_2$$

$$\text{distance} = d_1 + d_2$$

• How to decide K



KNN Regressor

Given → Predictions for training
→ Predictions for testing
→ Predictions for new data
→ Prediction for new data
will be avg(price) of nearest
K d.p.
No of train



points of X are given and
point for which we have to find

Advantages

- Easy to understand/very intuitive,
- Performance of model in terms of additive metric is good

Dis Adv

- Lazy learners

→ You are calculating
parameters of model
in real time - know
(Lazy learners)

- Computation time is more.

before fit best fit line used
fit → Coeff's
Predict.

$$R_0 + \beta_1 X_1 + \beta_2 X_2$$

- All the model parameters are calculated while training & used for prediction.
(Eager learners)

ML
↓
eager
learns
lazy
learns
↓

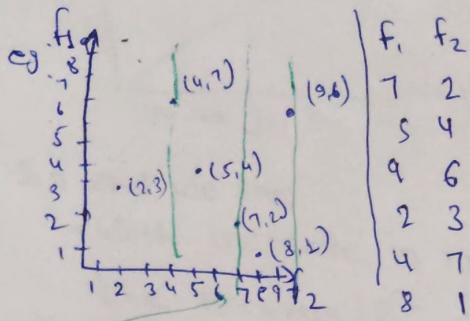
- ① Locally Weighted Learning
- ② Case-Based Reasoning
- ③ Lazy Bayesian Rules

if small dataset \rightarrow use KNN (No industry use \rightarrow because lazy learns
why KNN lazy)

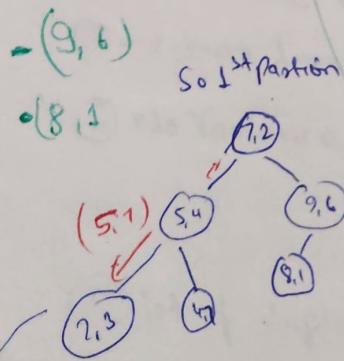
Variants of KNN (to optimize the time)

- K-D Tree
- Ball tree

In real time the distance of test-data from each of dp is calculated \Rightarrow Brut force



$f_2 = 2, 4, 5, 7, 8, 9$
median = 7
So partition



Q) What is the advantage here

If new dp come, (5, 7)
So directly come when

You know the data tree

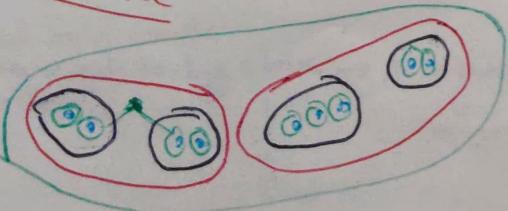
\hookrightarrow easily track the data points

Neighbours

• So instead of calculating Using Searching

1. 1st

② Ball tree



Advantage:

- only need to calculate distance of the nearest group element. (This also use hierarchical)



Cosine Similarity
range (0 to 1)
1 - similar
0 - dissimilar

KNN as missing Value Imputer

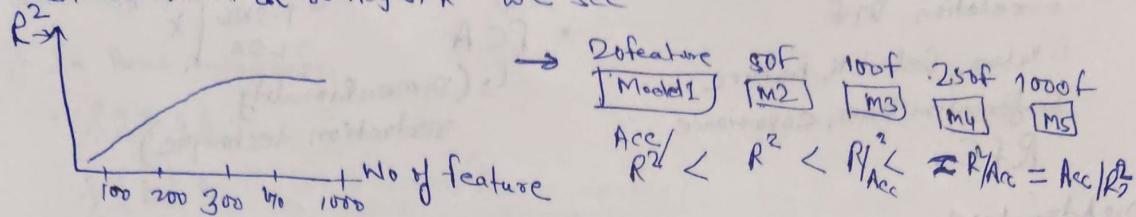
KNNImputer (n_neighbors=2)

from sklearn.impute import KNNImputer

23rd June : PCA (Principle Component Analysis)

Example

When we train R^2 to Adjust R^2 we see



So what we see

- With increase in no of features, after one point of the Accuracy / R^2 Stop increasing.

→ So

Q) Why it happen so?

↳ maybe

- Multi Collinearity ($f_1 f_2 f_3 \approx f_4$)
- All the entries of features are exactly same
- two feature are same

↳ No Variance

$$\begin{pmatrix} f_1 \\ 1.01 \\ 1.02 \\ 1.03 \\ 1.04 \\ \vdots \end{pmatrix} \quad \left| \begin{array}{c|c} f_1 & f_2 \\ \hline 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ \vdots & \vdots \end{array} \right| \quad \left\{ \begin{array}{l} f_1 \\ f_2 \\ \vdots \end{array} \right\} \text{No Variance in data}$$

(and we know
ML all
about learn
Pattern)

↳ lot of duplicate entries

$$\begin{pmatrix} f_1 \\ 1 \\ 2 \\ 2 \\ \vdots \end{pmatrix}$$

So what point I made that the feature we using might be not that much useful.

We know [with increase in no of feature performance of model degrades]

↓ When increase the no of feature why performance decrease?

Execution time
Validation time

Real world example.

You want to buy a house → go to broker → buy 3BHK → 60L

then you say → near Airport → 70L
again add

like that model also confuse when
no of feature.

$$\begin{matrix} Beach & \rightarrow 80L \\ Golf & \rightarrow 90L \\ Mkt & \vdots \\ Hospital & \vdots \end{matrix}$$

→ If you add More
Brokers confuse what
will be price

It is Course of dimensionality with increase
in no of feature the Performance of Model decreases

⇒ To Remove Curse of Dimensionality:

① Feature Selection

- Correlation, DIF.
- P-value, Select K feature
- Information Gain, Covariance
- RFE.

DisAdvantage

- Dropping data
- Time Consuming

② Feature Extraction

- PCA
 - ↳ T-SNE
 - ↳ LDA
- also put not see
- ↳ (Dimensionality reduction technique)

Q. Why we should remove Curse Dimensionality?

① To Improve performance of Model

② Visualize the data → for insight (for human eyes)

③ Prevents from overfitting

④ Better Interpretation

⑤ To Remove Curse of Dimensionality

* Feature Selection

→ CORR, COV

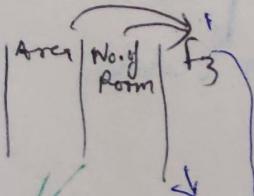
eg	Area	No. of Room	Price of house
1000	1000	1000	1000

if we choose or see corr b/w & justify...
then what we do with rest of the feature maybe
we drop it.

So what we do if not drop

Do

* Feature Extraction : Do data transformation to extract new feature which represents both the feature.



So How we can do this
How you do House Area Price
Because you have domain knowledge (Common Sense)

using PCA under
then come PCA

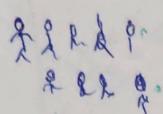
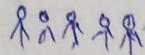
So when more feature domain expert also confuse

lets	# No of bathroom	distance	Distance from University
		Airport	

Why PCA come

⇒ PCA

Eg: In Batch Photography



Total

Cameras on
this side

So how many

person will
capture → All person

How in Real world works

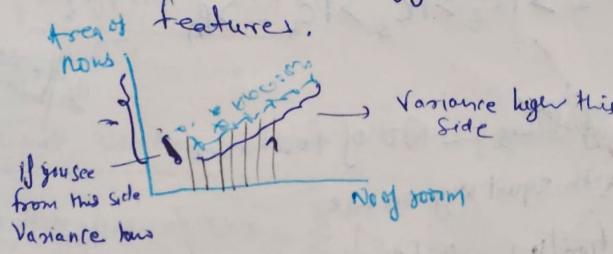
So how many person captured
in camera only 4

So what's conveyed is

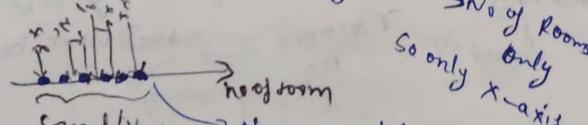
Variance / Spread is more compare to this side
you understand why



So same analogy to three
features.



- Now 1 more point
you calculate corr & decide to drop Area of house
then what is one thing left.



Now are data point in
X-axis and lost
this variance

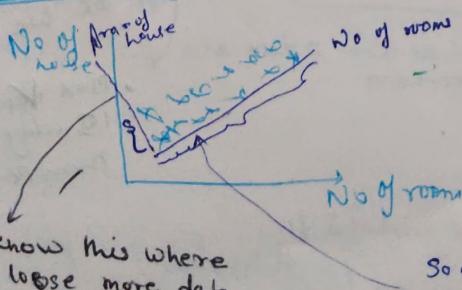
No. of rooms
So only X-axis

And we know
ML all about
pattern
So Avoid losing
this information

What is Conclusion

Remove feature loose the data

↳ So Now what do

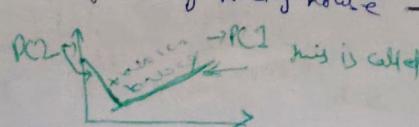


So know this where
we loose more data

So know if use No. of rooms
we cover 80% of area of house

So what we do to capture
or variance without losing
data

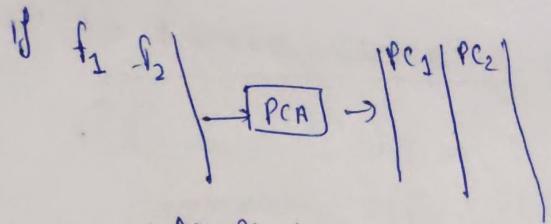
So can we change the axis
like this



→ 2D to 1D → (No. of rooms only)

without losing much of
Variance / data))

only Change
the axis
Coordinate
Rotation



• As PC_1 has maximum variance you will select only PC_1 for model training.

④ So why not we use PC_2 also → if we both then why we reduce.

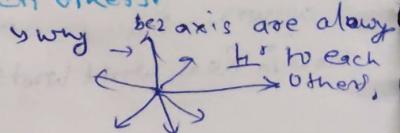
⑤ Sir if we have 1000 feature what we do

$$f_1 | f_2 | f_3 | \dots | f_{1000} \rightarrow PC_1 | PC_2 | PC_3 | \dots | PC_{1000}$$

but always remember those 1000 PC has 2 Characteristics

- ① All of these PC's will be $\perp\!\!\!\perp$ to each others.
- ② PC The Variance Capture by

$$PC_1 > PC_2 > PC_3 > PC_4 \dots > PC_n$$



- ③ No of PCs = No of features

④ (if both equal why we use)

$$\begin{array}{c|c} f_1 | f_2 | \dots | f_{1000} | y \\ \hline \end{array} \rightarrow PC_1 | PC_2 | \dots | PC_{1000}$$

$\frac{\gamma \times 1000}{d}$ $\gamma \times 1000$

use 2nd point
 $PC_1 > PC_2 > PC_3 > \dots$
 then what benefit

$$PC_1 | PC_2 | PC_3 | PC_4 | \dots | 100, 200, 300, 400, \dots$$

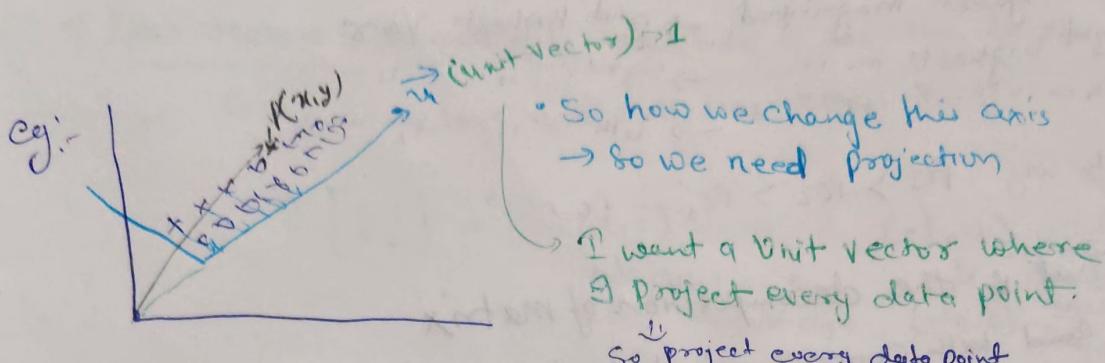
4D \rightarrow 3D
 use PC_1, PC_2, PC_3

So How many PC we Select

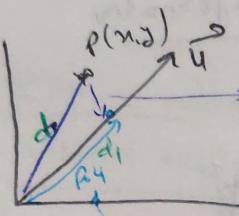
• Max Variance is along $PC_1 > PC_2 > \dots$

PCA maths

* Axis transformation \rightarrow Eigen decomposition of covariance matrix



Ques



you want to projection in \vec{u} using dot product

$$\text{Proj of } \vec{p} \text{ on } \vec{u} = \frac{\vec{p} \cdot \vec{u}}{\|\vec{u}\|} \quad \rightarrow \text{Also see in SVM}$$

$$\|\vec{u}\| = 1$$

$$\text{Proj of } \vec{p} \text{ on } \vec{u} = \vec{p} \cdot \vec{u}$$

d, d1, d2, ... are distance from origin in unit vector

So what is projection here
are distance

If you want that unit vector where variance / spread is max.

So you know the Variance formula: $\Rightarrow \text{Var} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

Aim: to find that unit vector which captures maximum variance after projection.

These are understanding

So know what is Covariance

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

what we do

if $x_1, x_2 \rightarrow$ we make \rightarrow co-variance matrix

$$\begin{matrix} x_1 & x_2 \\ \hline x_1 & \text{Cov}(x_1, x_1) & \text{Cov}(x_1, x_2) \\ x_2 & \text{Cov}(x_2, x_1) & \text{Cov}(x_2, x_2) \end{matrix} \text{ or } \begin{matrix} \text{Var } x_1 & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Var } x_2 \end{matrix}$$

So why we calculate this \leftarrow if 3 feature we 3 by 3 matrix so on. Covariance matrix or

\downarrow whole idea is to \downarrow to capture the maximum variance. Variance / Covariance matrix

Theorem if you decompose a covariance matrix of feature, you will get eigen value & eigen vector and eigen vectors

with highest magnitude of eigen value capture
the maximum variance/spread.

- ↳ 2nd highest magnitude → 2nd highest Var
- ↳ 3rd highest mag., → 3rd " "
- 4th " " → 4th " "

So that's why $\text{PC}_1 > \text{PC}_2 > \text{PC}_3 \dots \text{PC}_n$.

Q: Then what is the decomposition of matrix?

↳ So go ahead with concept → that is Linear transformation of Matrix

$$\begin{bmatrix} 7 & 9 & 2 \\ 8 & 5 & 1 \\ 9 & 0 & 3 \end{bmatrix}$$

or

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 3 & 0 \\ 7 & 8 & 9 \end{bmatrix}$$



'Analogy'

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

So as coordinate

So linear Trans
Say then same
vector/matrix
looks like this

(changing the
orientation)

* Linear transformation

↳ A matrix transformation that brings changes in the coordinate system.

$$A\vec{v} = \lambda \cdot \vec{v}$$

Matrix ↓ Eigen Value Eigen Vector Eigen Value Eigen Vector

(When we use co-variance matrix
as A, eigen decomposition
of covariance matrix)

How it clearer.

Grobner.org

↳ Linear algebra
Visualization

↳ show magnitude

$$\begin{bmatrix} 6 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

*

Linear transformation pointers

- Many vectors were changing both magnitude & direction.
- Few vectors changed only magnitude & not the direction
 - These vectors are called as eigen vector

eg $\begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix} \rightarrow \text{eigen } (1, 0)$

↑
takk direction same
(3, 1) magn. change

↓
direction same
(-4, -2) magn. change

eigen vector =

Eigen Value are change in magnitude for eigen Vector.

$$A\vec{v} = R\vec{v}$$

Co-variance Scalars ####

- The vector which only changes magnitude and not direction will be equals to dimensions of matrix / no. of features.

No of feature = No of PC

$$AV = R V \rightarrow \text{vector of } p_1, p_2, p_3$$

Stack / Srink

Step to calculate Eigen Value / Eigen vector (PC's)

np.linalg.svd.

Step-1 :- Standardize the data (make the data mean centered)

Step-2 :- Covariance matrix (df.cov())

Step-3 :- Eigen decomposition (np.linalg.svd) -

where you decompose $AV = HV$ $H \rightarrow$ eigen Value $V \rightarrow$ eigen Vector (direction not change)

1000 features \rightarrow 1000 PCs

$PC_1 > PC_2 > PC_3 \dots PC_n$

* Max var will be captured by first few PCs

if more than 50 var to go

\rightarrow PCA explained - Variance ration

for choose how much PCA.

sum(%) \Rightarrow np.cumsum (percentage)

Implementation is data (load data)

PCA.fit_transform(x_train)

PCA.components_ \rightarrow covariance

⇒ DisAdvantage of PCA

- Variance spread across all axis is same \rightarrow PCA will fail.
- PCA will loose the pattern $\cancel{\text{if}} \rightarrow \text{for}$

MNIST → handwritten digits (0-9)

Image: 28×28 pixels

PCA

→ before/After
PCA do perform
then ML

Mnist dataset with PCA + random forest

$$VA = VA$$

After

before

(128) rotation/center crop selected of 90%
(other 100% data are zeros) plot with scatter plot

(100, 100) center crop 100x100

(crop from center) 100x100

100x100

100x100