

★ IMP TERMINOLOGIES:

- (1) Training and Testing
- (2) Bias and Variance
- (3) Overfitting and Underfitting
- (4) Evaluation metric : Train / test accuracy, precision, recall, F1 score , Confusion matrix - TP, TN, FP, FN
- (5) Complexity \Rightarrow No of parameters , categorical/numerical data, relation between features, size of dataset
- (6) Curse of dimensionality : Dimension reduction by CNN, PCA, etc.

- (Q) What are dif types of datasets used in ML?
- (Q) How the validation & cross validation is dif from train and test accuracy?
- (Q) Why is it imp to split data for ML algos to verify the performance?
- (Q) What are the ethical considerations when one develops and deploys ML algos?
- (Q) What are the issues related to privacy and data ethics?
- (Q) What are the future trends of ML?
- (Q) Does the ML domain impact industries and employment?
- (Q) Bard vs Gemini Pro Vs Gemini Ultra

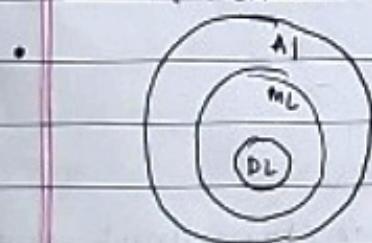
- Reference book : Bishop (2nd one in syllabus)

> Objectives of ML:

- Process Data
- Analyze
- Conclude from data
- ~~Sensible~~ ^{Sensible} output
- Netflix uses ML for recommendations
- Facebook uses ML for targeted ads and auto-tagging which

studies personal features using Deep Mind concept. (ML+NLP)

- Spam filtering: Classification Problem
- Need for ML:
 - Increase in data generation
 - Improve decision making
 - Uncover patterns & trends
 - Solve complex problems
- ML: Field of study that gives the computer capability to learn without being EXPLICITLY programmed. In ML you know both i/p and expected o/p. In traditional programming, ~~you~~ or i/p generated output.
- Train phase and Test phase.
- Objective is to lessen error margin : Ideal o/p=actual o/p.
- Train phase should have the higher portion of samples so that it can tune more parameters ~~so~~ to make the model achieve test accuracy almost equal to train accuracy.



→ Branch of AI, concerned with design and development of algorithm.

→ ML is based on automation and acquired knowledge.

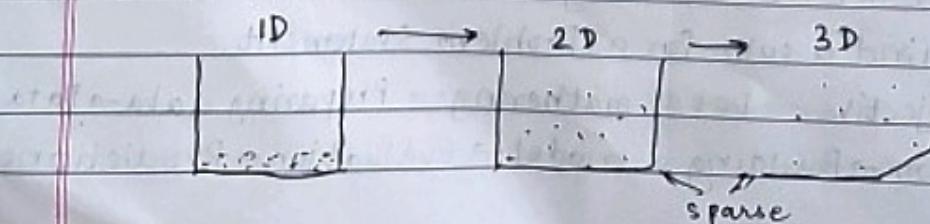
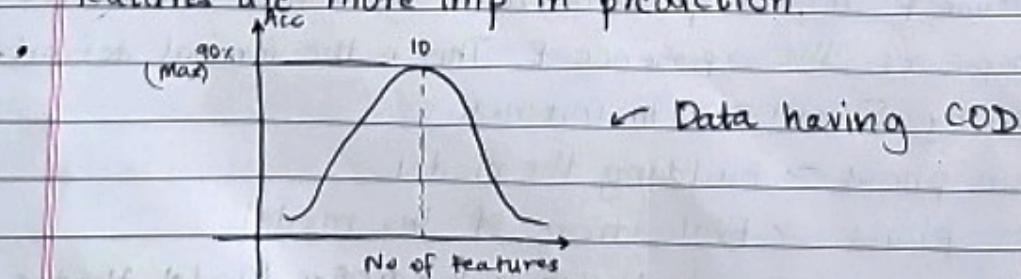
- ML: A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves the experience E. This is the formal definition given by Samuel Tom m. Mitchell.
- Train phase → Building the model
Test phase → Evaluation of the model
- ML process involves building a Predictive Model that can be used to find a soln for a Problem Statement.
Define Objective → Data Gathering → Preparing data → Data Exploration → Building a model → Evaluation → Predictions

- Types of ML algos:
 - Supervised learning
 - Unsupervised learning
 - Reinforcement learning
 - Supervised : Technique in which we teach ~~the~~ or train the machine using data which is well-labelled.
 - Unsupervised : Training of machine using info that is unlabeled and allowing the algo to act on that information without guidance.
 - Reinforcement : An agent is put in an environment and he learns to behave in this environment by performing certain actions & observing rewards for those actions.
 - Association Rule : Eg Market Basket Analysis
- Curse of Dimensionality (COD)

M ₁	M ₂	M ₃	M ₄
2f Acc ₁ =60%	3f Acc ₂ =70%	5f Acc ₃ =80%	10f Acc ₄ =90%
M ₅		M ₆	
50f Acc ₅ =85%		100f Acc ₆ =80%	

Acc \Rightarrow Accuracy, f \Rightarrow features

- As features start increases beyond certain limit, accuracy will start decreasing rapidly. This is curse of dimensionality.
- Overfitting occurs and model gets confused as to which features are more imp in prediction.



- Data points get scattered as no of features increases.
- Variance and bias increase.
- Underfitting: In a box of balls put an orange if model is trained only on 2f i.e., shape=spherical and color=any then orange will get identified as ball. This is underfitting.
- C_{OD} \Rightarrow Overfitting OR underfitting, not able to generalize $\xrightarrow{\text{underlying pattern}}$
not utilize imp features.
- Objective of ML is keep the model as simple as possible but good prediction (not in cost of complexity and not to confuse the model). so C_{OD} generally misleads the model in predicting ~~no~~ output.
- ML algo is a function.

$$y = f(x)$$

dependant \downarrow independent variable

- We get sweet pt (middle pt of over and underfitting) by regularization techniques.
- When you give data, for prediction to ML model, they try to find out underlying function $y=f(x)$ present in data.
- Parametric ML algo: Make assumptions about the function of data to give output. Eg: Linear regression. We assume relation between data is linear i.e., $y=mx+c$. So it is 1D with 2 parameters (m and c).
- No of parameters does not grow based on no of samples.
- Eg: Logistic regression, SVM, Naive Bayes, Simple NN.
- Non-parametric ML algo: No assumptions about the function of data. Eg: Decision Tree: Conditions on parameters are given wrt which decisions are made.
- No of parameters grow wrt no of samples.
- Eg: Random forest, SVM with dif kernels, Complex NN, Xboot

- Parametric

- Assumptions made about the function of dataset
- No of parameter does not grow wrt no of samples
- Less data reqd to train the model.
- Due to assumptions, chances of mistakes, i.e. Underfitting.
- Examples

Non-Parametric

- No assumptions about the function.
- No of parameters grows wrt no of samples
- No assumption, so more data reqd for training.
- No assumption so we try to consider each data point which can lead to Overfitting
- Examples

- Linear Regression:

$$y = mx + c$$

DEFV slope IND V → intercept

slope \Rightarrow rate of change of y wrt x

intercept \Rightarrow value of y at $x=0$.

Eg.

square of Dev(x)	student score (x)	student score (y)	Mean (\bar{x})	Mean (\bar{y})	Dev(x)	Dev(y)	Prod of Dev	SOP
+	8	10	10	13	-2	-3	6	12
0	10	13	10	13	0	0	0	0
+	12	16			+2	+3	6	

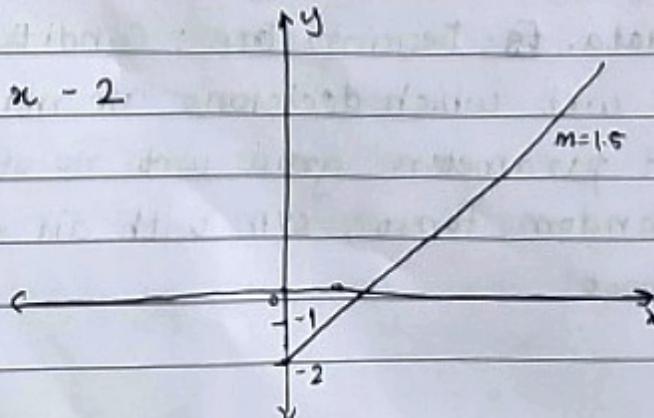
$$m = \frac{\text{Sum of product of Dev}}{\text{Sum of squares of Dev}x} = \frac{12}{6} = 1.5$$

$$b = \bar{y} - (m * \bar{x}) = 13 - (1.5 * 10) = -2$$

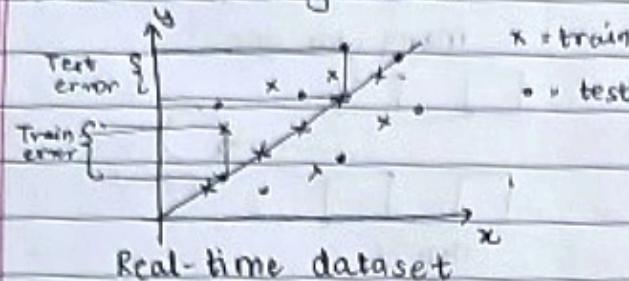
$$\text{Eqn: } y = 1.5x - 2$$

$$\text{Eq: } x = 20$$

$$y = 30 - 2 = 28$$



- Least error Regression.



It is the mathematical method used to find best fit line that represents relationship between dependent and independent variables in such a way that error is minimized.

Eg:

Price of Tshirt(x)	Tshirt sold(y)	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_{dev})(y_{dev})$	$y = mx + c$	Error (p)
2	4	-3.2	-4.2	13.44	3.3424	0.657
3	5	-2.2	-3.2	7.04	4.8604	0.139
5	7	-0.2	-1.2	0.24	7.8964	-0.896
7	10	1.8	1.8	3.24	10.9324	-0.932
9	15	3.2	6.8	25.84	13.9684	1.031
$\bar{x} = 5.2$	$\bar{y} = 8.2$			SOP = 49.8		$y_{actual} - y_{pred}$

$$m = \frac{49.8}{32.8} = 1.518$$

$$b = 0.3064$$

HW $x = 2 \ 4 \ 6 \ 8$ } Both linear and least error.
 $y = 3 \ 7 \ 8 \ 10$

- Multivariable Regression \Rightarrow more than one independent variable.

Eg: Restaurant Delivery System:

y : Total travel time (hr) — DEP. V

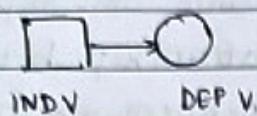
x_1 : Travel distance in miles } IND Vs

x_2 : No of deliveries

$$\therefore y = f(x_1, x_2)$$

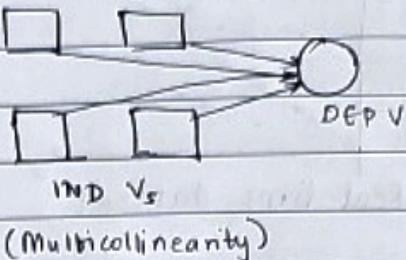
- Simple

One → One



Multivariable

Many to one



- Adding more no of variables (independent) in regression does not mean ~~go~~ better performance, infact it can make the things worse, i.e. overfitting.
- Adding more variables creates more relationships among them. IND Vs are not only potentially related to y, but also related to each other. ~~so~~ This is called multicollinearity
- Ideally, all IND Vs should be correlated with DEP V and NOT with each other.

- For multivariable regression:

(1) Simple Regression

(2) Scatter plots

(3) Correlation matrix.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n + \epsilon$$

\uparrow
Error

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

• Interpretation of coeffs in multivariable Regression:

$$\text{Eg: } y = 27 + 9x_1 + 12x_2$$

→ Estimated change in y corresponding to one unit change in variable when all other variables remain constant.

(Q) Predict the sale of the product for which using IND V x_1 product 1 and IND V x_2 product 2.

Sale/week			bias due to overfitting		
P ₁	P ₂		x =	y:	
x ₁	x ₂	y	1 1 4 1 2 5 1 3 8 1 4 2	1 6 8 12	
1	4	1	4x3	4x1	
2	5	6			
3	8	8			
4	2	12			

$$\hat{a} = ((x^T x)^{-1} x^T) y \text{ where}$$

$$a = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} \leftarrow \text{coeffs.}$$

$$x^T x = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 & 4 \\ 1 & 2 & 5 \\ 1 & 3 & 8 \\ 1 & 4 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 10 & 19 \\ 10 & 30 & 46 \\ 19 & 46 & 109 \end{bmatrix}$$

$$(x^T x)^{-1} x^T = \begin{bmatrix} 3.153 & -0.59 & -0.3 \\ -0.59 & 0.205 & 0.016 \\ -0.3 & 0.016 & 0.055 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} 1.36 & 0.47 & -1.02 & 0.19 \\ -0.319 & -0.098 & 0.1557 & 0.2622 \\ -0.065 & 0.0056 & 0.1857 & -0.125 \end{bmatrix}$$

$$\hat{a} = \begin{bmatrix} -1.619 \\ 3.4751 \\ -0.053 \end{bmatrix}$$

$$\hat{y} = a_0 + a_1 x_1 + a_2 x_2$$

$$\hat{y} = -1.619 + 3.475 x_1 - 0.053 x_2$$

> 2nd method of MVR

$$y = a_0 + a_1 x_1 + a_2 x_2$$

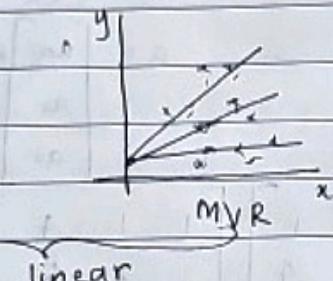
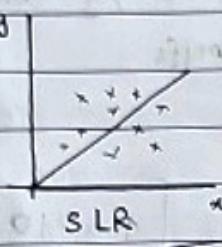
$$a_0 = y - a_1 x_1 - a_2 x_2$$

$$a_1 = \sum (x_1)^2 * \sum (x_1 y) - \sum (x_1 x_2) * \sum (x_2 y)$$

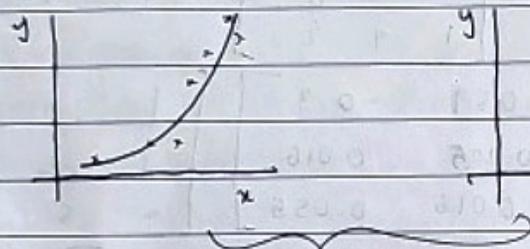
$$a_2 = \sum (x_2)^2 * \sum (x_2 y) - \sum (x_1 x_2) * \sum (x_1 y)$$

$$\sum (x_1^2) * \sum (x_2^2) - (\sum (x_1 x_2))^2$$

>



linear



$$y = a_0 + \sum_{i=1}^n x_i y_i + FP$$

non-linear

> Polynomial

$$y = a_0 + \sum_{i=1}^n x_i y_i + FP$$

Polynomial function

linear

> Non-linear to linear:

$$y = a e^{bx}$$

$$\ln y = \ln(a) + \ln(e^{bx})$$

$$= \ln(a) + bx \ln(e)$$

$$\boxed{\ln y = \ln(a) + bx}$$

$$y = a_0 + a_1 x + a_2 x^2$$

$$a = x^{-1} B$$

$$X = \begin{bmatrix} 1 & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix}$$

$$B = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \\ \sum x_i^2 y_i \end{bmatrix} \quad n = \text{no of samples}$$

(Q)	X	Y	x_i^2	x_i^3	x_i^4	$n y_i$	$x_i^2 y_i$
	1	1	1	1	1	1	1
	2	4	4	8	16	8	16
	3	9	9	27	81	27	81
	4	15	16	64	256	60	240
	$\sum x_i = 10$	$\sum y_i = 29$	$\sum x_i^2 = 30$	$\sum x_i^3 = 100$	$\sum x_i^4 = 354$	$\sum n y_i = 96$	$\sum x_i^2 y_i = 338$

$$X = \begin{bmatrix} 1 & 10 & 30 \\ 10 & 30 & 100 \\ 30 & 100 & 354 \end{bmatrix} \quad B = \begin{bmatrix} 29 \\ 96 \\ 338 \end{bmatrix}$$

$$X^{-1} B = \begin{bmatrix} 7.75 & -6.75 & 1.25 \\ -6.75 & 6.45 & -1.25 \\ 1.25 & -1.25 & 0.25 \end{bmatrix} \begin{bmatrix} 29 \\ 96 \\ 338 \end{bmatrix} = \begin{bmatrix} 11.216 \\ 3.6470 \\ 13.0122 \end{bmatrix}$$

$$= \begin{bmatrix} -0.75 \\ 0.95 \\ 0.75 \end{bmatrix}$$

$$y = -0.75 + 0.95x + 0.75x^2$$

> Model Selection: Depends on type of dataset and type of task. These are based on logical reasoning and comparison between ML models.

- Type of data:

(1) Supervised:

- Continuous, linear \rightarrow SLR, MVR, SVM
- Continuous, non-linear \rightarrow Polyⁿ R, ~~SVM~~ SVM
- Categorical \rightarrow logistic R, decision tree, RF

(2) Unsupervised: K-means, H-Cluster.

- Type of Task:

- (1) Regression: Linear OR Non-linear
- (2) Classification: Binary OR Multi-class
- (3) Clustering

- Bias-Variance Trade off(?)

- (1) Bias error ↴ ~~reducible~~ reducible
- (2) Variance error ↴ i.e., can minimize
- (3) Noise - unreliable, cannot minimize

- Bias is dif between ~~two~~ predicted and actual data points

$$y = f(x) + e$$

↳ irreducible error

$$f'(x) = p \quad A = f(x)$$

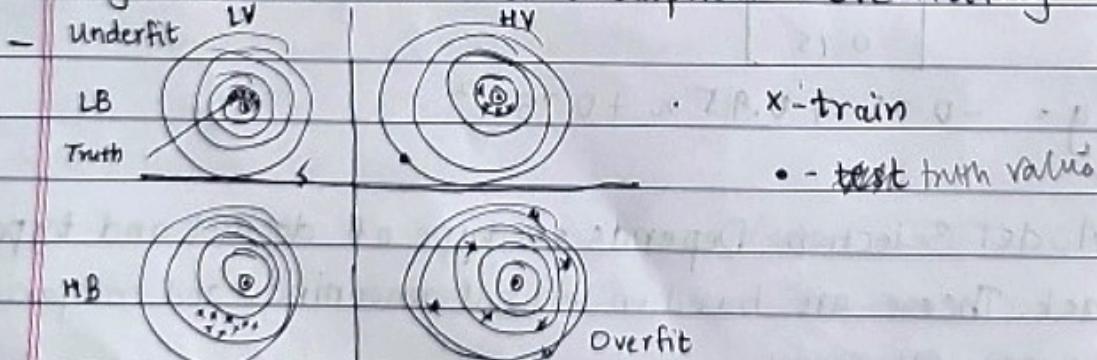
$$\text{Bias } f'(x) = E [f'(x) - f(x)]$$

- Variance is - when model takes into account the fluctuations in the data, i.e., noise as well.

$$\text{Variance } f(x) = E [x^2] - E(x)^2$$

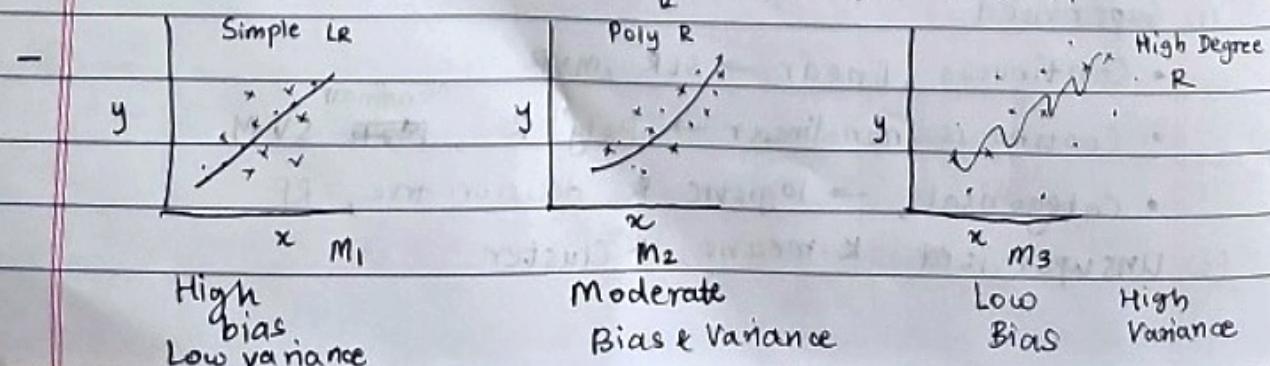
- High bias \rightarrow simple assumptions \rightarrow Underfitting

- High variance \rightarrow Too much assumption \rightarrow Overfitting

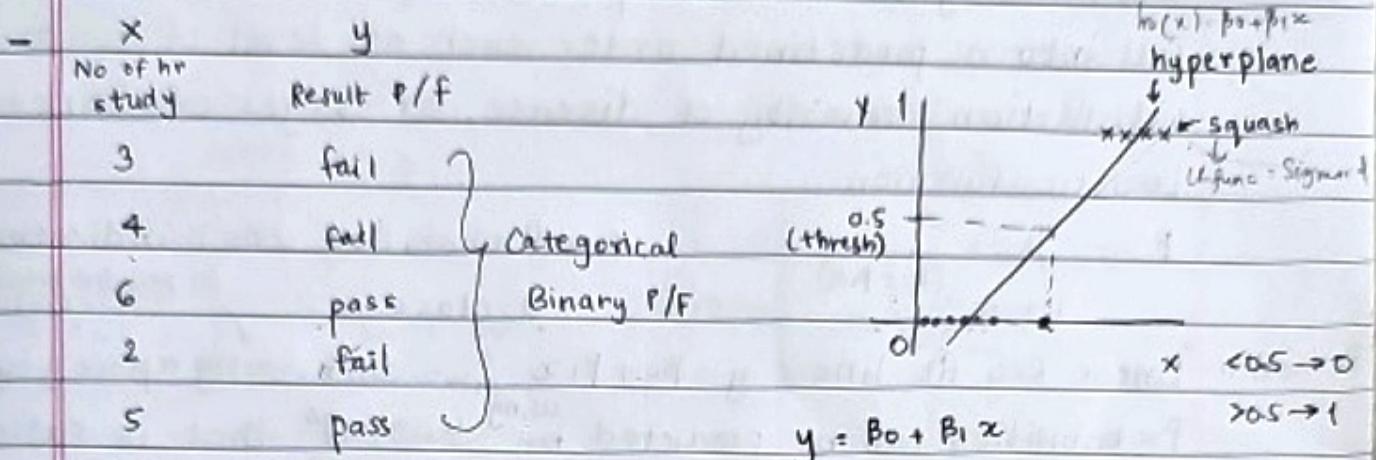


Bullseye diagram

↳ Bias-Variance Trade-off



Total error: Bias² + Variance + IR error



Adding data points may lead to outliers leading to misinterpretation of data.

- When you apply Linear Regression on the categorical output of the data points, it suffers with two problems.
 - (i) When you add more data pts, it faces the problem of outliers. Misclassification may occur.
 - (ii) Possibility of predicted data pts maybe greater than 1 or less than 0.

In order to overcome these, we switch to Logistic Regression.

- hyperplane $h_0(x) = B_0 + B_1x$
squash \Rightarrow Logistic Regression function = $\frac{1}{1+e^{-B_0-B_1x}}$ (Sigmoid)
- Logistic Regression is ^{the} appropriate regression to conduct when dependent variable is categorical / binary. It is predictive analysis. It is used to describe data & explain relationship between one DEP binary 1 and one or more nominal / ordinal / interval / ratio level IND V.
- Binary $\rightarrow 0/1$
- Multinomial logistic : Used to predict probability of one of three or more possible outcomes. Eg: Type of product a customer will buy, rating a customer will give a product.
Eg: Political party a person will vote for.

> Logistic Regression is linear in nature as we use $y = \beta_0 + \beta_1 x$ to get the hyperplane. This hyperplane is in 'S' shape due to use of sigmoid function for calculating probability.

- Ordinal Logistic: Predicts probability of an outcome that fall into a predefined order such as level of customer satisfaction, severity of disease, or stages of cancer.
- Logistic function:

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$P \rightarrow$ Probability which indicates a class. (as prob lies between 0,1)

Derivation: Linear Reg fit line: $y = \beta_0 + \beta_1 x$, we can't directly replace y by P . Probabilities can be predicted using "ODDS" that is ratio of prob of success and prob of failure.

$$\frac{P}{1-P} = \beta_0 + \beta_1 x \quad \text{--- (1)} \quad (\text{reduces no of data points} \Rightarrow \text{decrease in computation})$$

1-P always +ve $[0, \infty]$, but this is restricted range

Take log to get high correlation.

$$\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x \quad \text{--- (2)} \quad \text{Range: } (-\infty, \infty)$$

But we are not interested in log, but exponent value to predict probability.

$$e^{\ln(P/(1-P))} = e^{\beta_0 + \beta_1 x}$$

$$\therefore \frac{P}{1-P} = e^{(\beta_0 + \beta_1 x)}$$

$$\therefore P = \frac{e^{(\beta_0 + \beta_1 x)}}{1 + e^{(\beta_0 + \beta_1 x)}}$$

$$P = P \left(\frac{e^{(\beta_0 + \beta_1 x)}}{P} - e^{(\beta_0 + \beta_1 x)} \right)$$

$$\therefore 1 = \frac{e^{(\beta_0 + \beta_1 x)}}{1 + e^{(\beta_0 + \beta_1 x)}} - e^{(\beta_0 + \beta_1 x)}$$

$$\therefore 1 - \frac{P}{1 + e^{(\beta_0 + \beta_1 x)}} = -e^{(\beta_0 + \beta_1 x)}$$

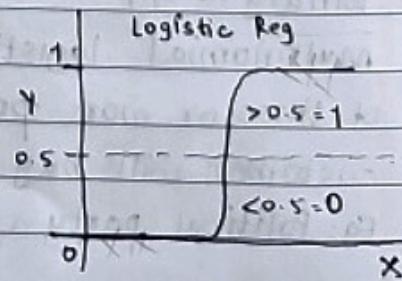
$$\therefore P \left[1 + e^{(\beta_0 + \beta_1 x)} \right] = e^{(\beta_0 + \beta_1 x)}$$

$$P = \frac{e^{(\beta_0 + \beta_1 x)}}{1 + e^{(\beta_0 + \beta_1 x)}} \quad \text{--- (3)}$$

Divide by $e^{(\beta_0 + \beta_1 x)}$

$$\therefore P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

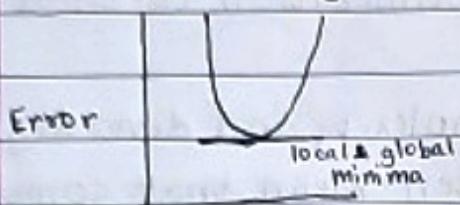
Logistic Reg



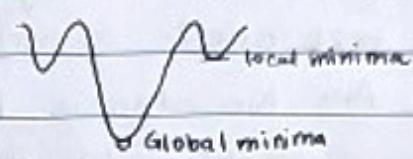
- You cannot classify points which lie exactly on threshold, but it is a very rare case so it is negligible.
- Cost function \Rightarrow Diff between predicted and actual

Linear Reg

logistic Reg



Error



Cost function

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

Cost function

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

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

maximum likelihood estimation

- On the given dataset apply logistic reg.

x (hrs ~~studying~~) y (~~w/l~~)

29 0

15 0

33 1

28 1

40 1

(1) calculate Probability of winning the game for player who practices 28 hrs.

(2) Atleast how many hrs, player should practice that makes him win the game with $P > 85\%$.

Assume log ODDS $\log(\text{odds}) = -64 + 2 \times \text{hr}$.

$$(1) \quad \log(\text{odds}) = -64 + 2 \times \text{hr} = p_0 + p_1 x$$

$$p_0 = -64, \quad p_1 = 2$$

$$P = \frac{1}{1 + e^{(-64 + 2x)}} \quad \text{Put } x = 28$$

$$= \frac{1}{1 + e^{(-64 - 56)}} = \frac{1}{1 + e^{-60}} = \frac{0.000033}{0.000033} = 0.000335 \Rightarrow 0.033\%$$

$$(2) \quad \frac{85}{100} < \frac{1}{1 + e^{(64 - 2x)}} \Rightarrow 1 + e^{(64 - 2x)} < \frac{100}{85}$$

$$\therefore e^{(64 - 2x)} < 0.17647 \Rightarrow 64 - 2x < -1.735 \\ \therefore x = 33 \text{ hrs} \Rightarrow 2x = 65.735$$

Player should practice for atleast 32.86 ~ 33 hrs so that he will win game with P of 85%.

- No Freeloader Rule:

There are so many algos for optimization so what is the best one?

Ans: No algo is best as complexity and diversity of real world problems often mean that some problems are easier to solve whereas others can be extremely difficult to solve. There is no one model, that works for every problem. Assumptions of great model for one problem may not hold for others. Thus, NFL states that no 2 algos are equivalent when their performance is average across all possible problems.

Mathematically NFL theorem: If an algo performs well on certain ~~a~~ class of problems, then it necessarily ~~not~~ pays for that with degraded performance on the set of all ~~other~~ remaining problems. In summary, there is no single best learning method.

(Q) Take NN with many layers optimized by Back Propagation i.e deep learning. Choose correct option:

- 1) Deep learning performs better than most other algos for real-world problems.
- 2) Deep learning can fit everything.
- 3) Deep learning performs better than other algos for all problems.

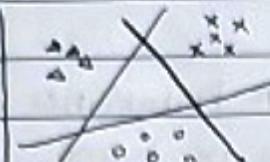
- Logistic Regression (Multiclass)

x_2	$\begin{matrix} \Delta & \Delta \\ \Delta & \Delta \\ \cdot & \cdot \end{matrix}$
	$\begin{matrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{matrix}$

x_1
(Binary)

It is simple. Variables are categorical. Confusion matrix \rightarrow derivation is simple
Eg: Spam / not spam

- Multiclass \rightarrow No of instances in dataset predicts no of classifiers.



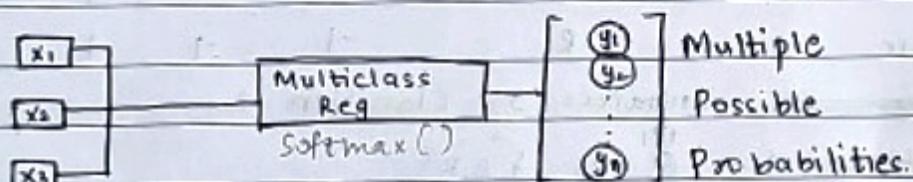
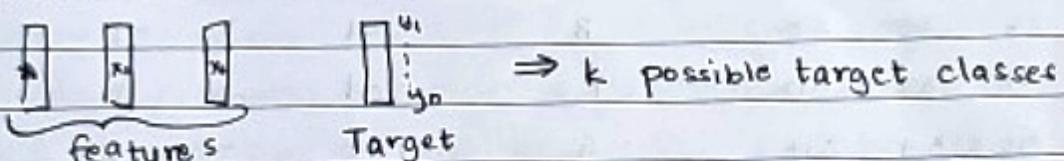
3 classes

This is complex. Confusion matrix is complex. Variables are categorical.

Eg: classification of fruits

For one to all (rest) \Rightarrow no of classifiers $k = N$ } Both can be
one to one $\Rightarrow \frac{N(N-1)}{2}$ classifiers } used
 $N = \text{No of instances}$

- Multiclass

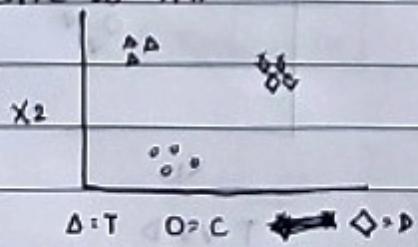


$$LR^n : y = mx + c$$

$$z_i(k+1) = \underbrace{w_{(k+1)m}}_{\text{wt. associated with feature}} \cdot x_{(m,i)} + b_{(k+1)} \leftarrow \text{Bias}$$

$$\begin{bmatrix} z_0 \\ z_1 \\ \vdots \\ z_n \end{bmatrix} = \begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & & \vdots \\ w_{n,0} & \dots & \dots & w_{n,n} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix}$$

- One to All



Instances = 3 \Rightarrow Classifiers = 3

Classifier 1: T vs {C, D} = y_1

2: C vs {T, D} = y_2

3: D vs {T, C} = y_3

- Then we apply softmax function to classify new test data pts into T, G, D classes.

$$y_1 \Rightarrow P_1$$

$$y_2 \Rightarrow P_2$$

$$y_3 \Rightarrow P_3$$

We will focus on the Probabilities and assign class which has majority vote count.

- $h_{\theta}^{(i)}(x) = P(y=i | x; \theta)$

$$(i=1, 2, 3)$$

(S)	Features	Classes	y_1	y_2	y_3
	$x_1 \quad x_2 \quad x_3$	G	1	-1	-1
	$x_4 \quad x_5 \quad x_8$	B	-1	1	-1
	$x_7 \quad x_6 \quad x_{12}$	R	-1	-1	1
	$x_{10} \quad x_9 \quad x_{14}$	G	1	-1	-1
	$x_{13} \quad x_{11} \quad x_{16}$	B	-1	1	-1
	$x_{17} \quad x_{15} \quad x_{21}$	R	-1	-1	1

~~$x_9 \quad x_{10} \quad x_{11}$~~ Instances = 3 Classifiers = 3.

Classifier 1 : G vs {B, R}

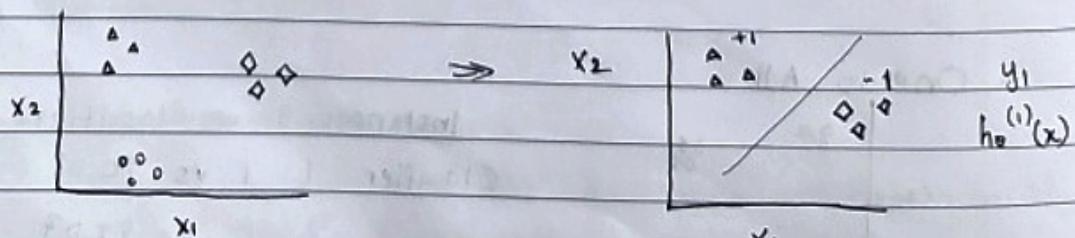
2 : B vs {G, R}

3 : R vs {G, B}

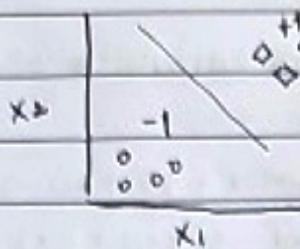
Softmax function will give probs P_1, P_2, P_3 .

$P_1 = 0.8, P_2 = 0.1, P_3 = 0.1$ lets say this is the case for one data pt, then ~~that~~ we take highest +ve $\Rightarrow P_1$ implying it belongs to G class.

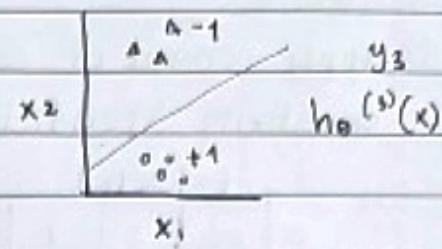
- One to One



T vs D



D vs C



C vs T

- (Q) Derive the cost function of multiclass Logistic Reg & compare with Binary logistic Reg cost function.

- (Q) Compare:

- Binary Logistic
- Sigmoid /logistic function is used to get the probabilities back from response variable y .
- Prob of +ve class is given as: $P = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}$
- $y = B_0 + B_1 * \text{feature_1} + B_2 * \text{feature_2} + \dots$
- Eg: Pass/Fail, spam/not spam
- Multiclass Logistic
- Softmax function is used to get probabilities back from response variable y .
- Prob of class i , given k classes, is: $\sigma(x)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad z = (z_1 \dots z_k) \in \mathbb{R}^k$
- $y_i = B_0 + B_1 * f_1 - 1 + B_2 * f_2 - 2 + \dots$
- One B is learnt per feature per class. Every class has its own hyperplane.
- Eg: Fruits, colors.

(Q)

Expected		Actual	
predicted	45 TP	8 FP	40 samples \rightarrow +ve
	15 FN	32 TN	40 samples \rightarrow -ve
			Total 100 samples.

(Q) Apply concept of confusion matrix and calculate accuracy, sensitivity, specificity and F1 score.

Expected					
		+ve		-ve	
Predicted	C1	52	3	7	2
	C2	2	28	2	0
C3	5	2	25	12	
C4	1	1	9	40	
	C1	C2	C3	C4	

+ve = 52 + 28 + 25 + 40 = 145
 Missclassified = 3 + 7 + 2 + 2 + 5 + 2 + 12 + 1 + 1 + 9 = 46
 Total = 191 samples
 Accuracy = $\frac{145}{191} \times 100 = 75.9\%$

$$\text{Accuracy} = \frac{TP + TN}{All}$$

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

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F1 \text{ score} = 2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

~~$$\text{Specificity} = \frac{\text{Precision C1}}{C1 + C2 + C3 + C4} = \frac{52}{52 + 3 + 7 + 2} = 0.8125$$~~

~~$$C2 = \frac{28}{28 + 2 + 2} = 0.875$$~~

~~$$C3 = \frac{25}{25 + 5 + 2 + 12} = 0.5556$$~~

~~$$C4 = \frac{40}{40 + 9 + 1 + 1} = 0.784$$~~

~~$$\text{Recall C1} = \frac{52}{60} = 0.867 \quad C2 = \frac{28}{34} = 0.823$$~~

~~$$C3 = \frac{25}{43} = 0.582 \quad C4 = \frac{40}{54} = 0.74$$~~

$$C1_{F1} = \frac{2(0.8125 \times 0.867)}{(0.8125 + 0.867)} = 0.8388$$

$$C_2 F_1 = \frac{2}{(0.875 + 0.823)} = 0.848$$

$$C_3 F_1 = \frac{2}{(0.568 + 0.582)} = 0.575$$

$$C_4 F_1 = \frac{2}{(0.784 + 0.74)} = 0.7614$$

$$\text{Precision} = \frac{TP_1 + TP_2 + TP_3 + TP_4}{TP_1 + FP_1 + TP_2 + FP_2 + TP_3 + FP_3 + TP_4 + FP_4}$$

$$= \frac{52 + 28 + 25 + 40}{52 + 12 + 28 + 4 + 25 + 19 + 40 + 11} = 0.759$$

$$\text{Recall} = \frac{TP_1 + TP_2 + TP_3 + TP_4}{TP_1 + FN_1 + TP_2 + FN_2 + \dots}$$

$$= \frac{52 + 28 + 25 + 40}{52 + 8 + 28 + 6 + 25 + 18 + 40 + 14} = 0.759$$

$$F_1 \text{ score} = \frac{2}{2} \frac{(0.759 \times 0.759)}{2 (0.759)} = 0.759$$

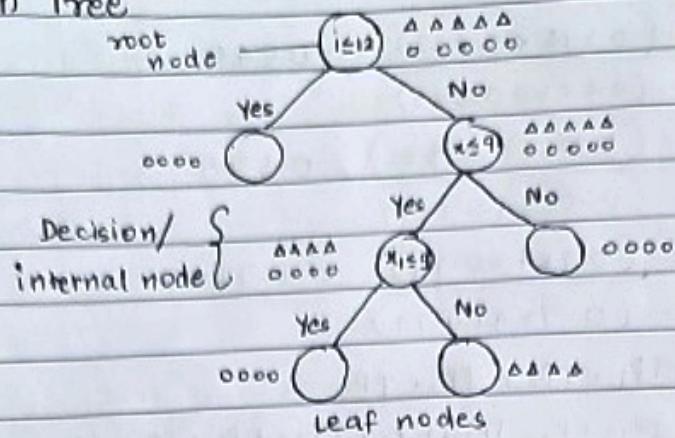
(Q) What factors can contribute to the popularity of Logistic Reg?

(Q) Is the decision boundary linear or non-linear in case of Logistic Reg model?

(Q) What is the impact of outliers on the Logistic Reg?

(Q) Which algorithm is better in case of outliers present in dataset from logistic, linear reg and SVM?

- Decision Tree



→ Steps:

Algo

CART
(Classification &
Regression Tree)
↓
Gini Index

ID3 ← requires high computational power
with unstable tree structure when
no of nodes increase.
↓
Entropy & info gain

• Steps:

- (1) Compute the entropy for dataset that is called as Entropy(s).
- (2) For every attribute/feature:
 - (i) Calculate entropy for all other values: Entropy(A)
 - (ii) Take avg information entropy for current attribute.
 - (iii) Calculate gain for current attribute.
- (3) Pick up the highest gain attribute.
- (4) Repeat until we get the tree we desire.

• Formula:

- (1) Calculate entropy (amt of uncertainty in dataset)

$$\text{Entropy} = -\frac{P}{P+n} \log_2 \left(\frac{P}{P+n} \right) - \frac{n}{P+n} \log_2 \left(\frac{n}{P+n} \right)$$

- (2) Calculate avg info

$$I(\text{attribute}) = \sum_{P+n} p_i n_i \text{Entropy}(A)$$

- (3) Calculate info gain

$$\text{Gain} = \text{Entropy}(s) - I(\text{attribute})$$

i.e., dif in entropy before and after splitting dataset into decision trees.

Eg:	Sr No	Outlook	Temp	Humidity	Wind	Play Tennis
	1	Sunny	hot	High	Weak	N
	2	Sunny	hot	High	Strong	N
	3	Overcast	hot	High	Weak	Y
	4	Rainy	Mild	High	Weak	Y
	5	Rainy	Cool	Normal	Weak	Y
	6	Rainy	Cool	Normal	Strong	N
	7	Overcast	Cool	Normal	Strong	Y
	8	Sunny	Mild	High	Weak	N
	9	Sunny	Cool	Normal	Weak	Y
	10	Rainy	Mild	Normal	Weak	Y
	11	Sunny	Mild	Normal	Strong	Y
	12	Overcast	Mild	High	Strong	Y
	13	Overcast	Hot	Normal	Weak	Y
	14	Rainy	Mild	High	Strong	N

Compute and draw decision tree. Decide if player plays tennis, based on weather conditions.

$$\rightarrow p=9 \quad n=14 \quad (+ve \text{ and } -ve)$$

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

Now split attributes : Eg: Outlook will be split into sunny, overcast, rainy.

$$\text{Entropy}(\text{sunny}) = -\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) = 0.971 \\ (p=2, n=3)$$

$$\text{Entropy}(\text{Rainy}) = -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.971 \\ (p=3, n=2)$$

$$\text{Entropy}(\text{Overcast}) = -\frac{4}{4} \log_2 \left(\frac{4}{4} \right) = 0 \\ (p=4, n=0)$$

$$I(\text{Outlook}) = \frac{5}{14} (0.971) + \frac{5}{14} (0.971) + 0 = 0.693$$

$$\text{Gain} = 0.94 - 0.693 = 0.247$$

For Temp,

$$\text{Entropy (Hot)} = -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right) = 1$$

(p=2, n=2)

$$\text{Entropy (Mild)} = -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2 \left(\frac{2}{6}\right) = 0.918$$

(p=4, n=2)

$$\text{Entropy (cool)} = -\frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{4} \log_2 \left(\frac{1}{4}\right) = 0.811$$

(p=3, n=1)

$$I(\text{Temp}) = \frac{4}{14}(1) + \frac{6}{14}(0.918) + \frac{4}{14}(0.811) = 0.911$$

$$\text{Gain: } 0.94 - 0.911 = 0.029$$

For Humidity:

$$\text{Entropy (High)} = 0.985$$

$$\text{Entropy (Normal)} = 0.59$$

$$I(\text{Humidity}) = 0.788$$

$$\text{Gain} = 0.152$$

For Wind:

$$\text{Entropy (Strong)} = 1$$

$$\text{Entropy (Weak)} = 0.811$$

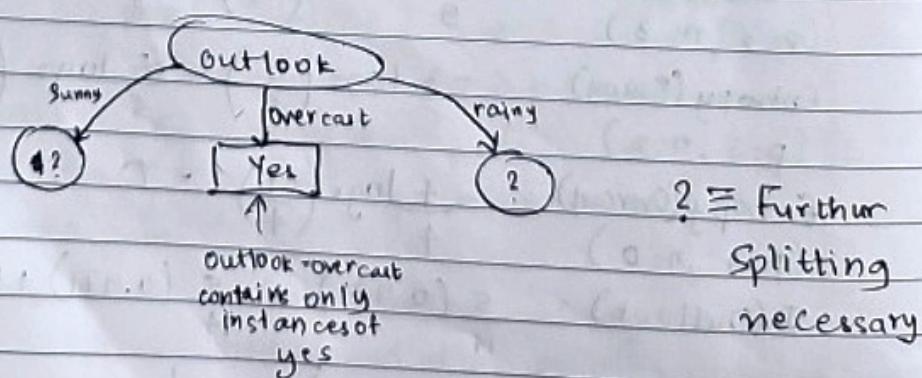
$$I(\text{Wind}) = 0.892$$

$$\text{Gain} = 0.048$$

Gains in Ascending order:

$$\begin{array}{cccc} \text{Gain(Temp)} & \leftarrow & \text{Gain(Wind)} & \leftarrow \\ 0.029 & & 0.048 & \\ \text{Gain(Humidity)} & \leftarrow & \text{Gain(Outlook)} & \\ 0.152 & & 0.247 & \end{array}$$

Choose Outlook as root node.



- Now split datasets according to Sunny and Rainy
- For Sunny:

Sr No	Temp	Humid	Wind	Play
1	hot	High	Weak	N
2	hot	High	Strong	N
3	Mild	High	Weak	N
4	Cool	Normal	Weak	Y
5	Mild	Normal	Strong	Y

$$\text{Entropy}(S) = 0.971 \text{ (Only Sunny)}$$

$$\text{Entropy}(\text{hot}) = 0$$

$$\text{Entropy}(\text{cool}) = 0$$

$$\text{Entropy}(\text{mild}) = 1$$

$$I(\text{Temp}) = 0.4 \quad \text{Gain} = 0.571 \quad (0.971 - 0.4)$$

$$\text{Entropy}(\text{high}) = 0$$

$$\text{Entropy}(\text{normal}) = 0$$

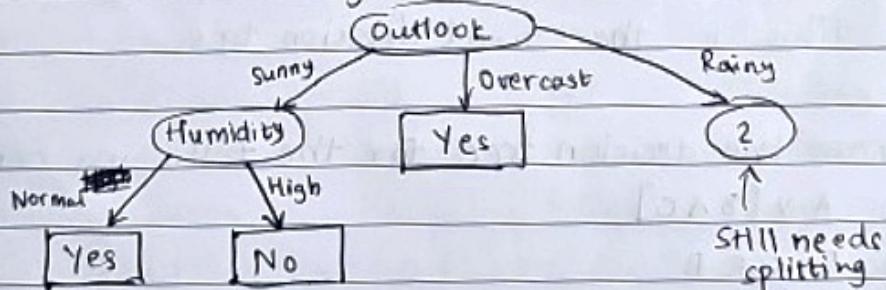
$$I(\text{Humid}) = 0 \quad \text{Gain} = 0.971$$

$$\text{Entropy}(\text{Weak}) =$$

$$\text{Entropy}(\text{Strong}) =$$

$$I(\text{Wind}) = 0.951 \quad \text{Gain} = 0.02$$

So we take humidity as next root



For rainy:

Sr No	Temp	Humid	Wind	Play	Entropy(S)
1	Mild	High	Weak	Y	0.971
2	Cool	Normal	Weak	Y	
3	Cool	Normal	Strong	N	
4	Mild	Normal	Weak	Y	
5	Mild	High	Strong	N	

$$\text{Entropy (Mild)} = -2 \times \frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) = 0.918$$

($p=2, n=1$)

$$\text{Entropy (Wet)} = 2 \times \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) = 1$$

($p=1, n=1$)

$$I(\text{Temp}) = \frac{2}{5} (1) + \frac{3}{5} (0.918) = 0.9508$$

$$\text{Gain} = [0.02]$$

$$\text{Entropy (High)} = 1$$

($p=1, n=1$)

$$I(\text{Humid}) = 0.9508$$

$$\text{Gain} = [0.02]$$

$$\text{Entropy (Normal)} = 0.918$$

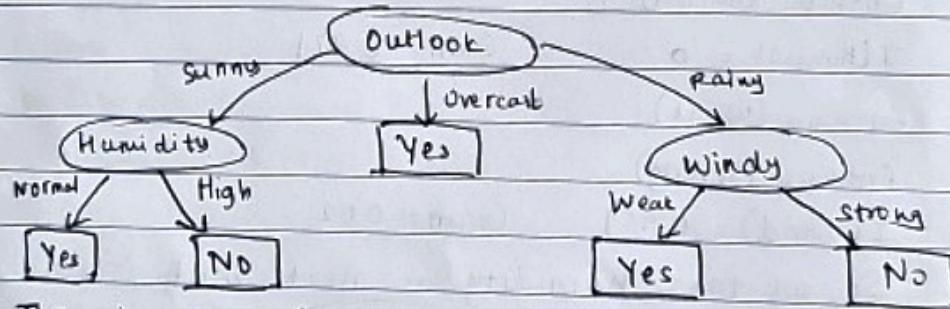
($p=2, n=1$)

$$\text{Entropy (Weak)} = 0 \quad I(\text{Wind}) = 0$$

$$\text{Entropy (Strong)} = 0$$

$$\text{Gain} = [0.971]$$

So take Windy as outlook



This is the final decision tree.

(Q) Draw the decision tree for the following condition:

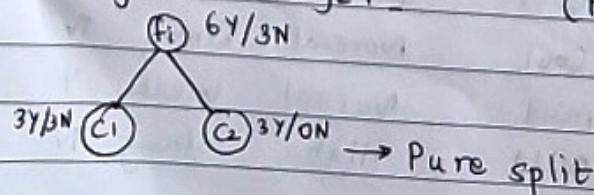
$$(1) A \vee [B \wedge C]$$

$$(2) A \times D \wedge B$$

$$(3) [A \wedge B] \vee [C \wedge D]$$

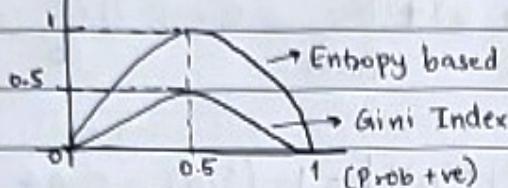
- Entropy based Decision tree:

$$H(s) = -P_+ \log_2 P_+ - P_- \log_2 P_- \quad (P_+ \Rightarrow +ve, P_- \Rightarrow -ve)$$



$$\text{for } C_2, H(C_2) = -\frac{3}{3} \log_2 (3/3) - 0 = 0$$

$$C_1, H(C_1) = -\frac{1}{2} \log_2 (1/2) - \frac{1}{2} \log_2 (1/2) = 1$$



$$\text{- Gini Index : } GI = 1 - \sum_{i=1}^n (P_i)^2 \quad \xrightarrow{\text{+ve & -ve}} \quad = 1 - [(P_+)^2 + (P_-)^2]$$

$$\text{For } C_1, GI(C_1) = 1 - [(1/2)^2 + (1/2)^2] = 1 - 1/2 = 1/2$$

Gini Index is computationally efficient as compared to the entropy based Decision tree.

(a) Make a decision tree according to Gini Index.

	Weekend	Weather	Parent	Money	Decision
1	w ₁	Sunny	Yes	Rich	Cinema
2	w ₂	Sunny	No	Rich	Tennis
3	w ₃	Windy	Yes	Rich	Cinema
4	w ₄	Rainy	Yes	Poor	Cinema
5	w ₅	Rainy	No	Rich	StayIn
6	w ₆	Rainy	Yes	Poor	Cinema
7	w ₇	Windy	No	Poor	Cinema
8	w ₈	Windy	No	Rich	Shopping
9	w ₉	Windy	Yes	Rich	Cinema
10	w ₁₀	Sunny	No	Rich	Tennis

$$GI = 1 - \sum_{i=1}^n (P_i)^2 = 1 - (0.6^2 + 0.2^2 + 0.1^2 + 0.1^2) = 0.58$$

For Money, Poor = 0.3, Rich = 0.7

$$GI(\text{Poor}) = 1 - (3/3)^2 = 0$$

$$GI(\text{Rich}) = 1 - [(3/7)^2 + (2/7)^2 + (1/7)^2 + (1/7)^2] = 0.694$$

$$\text{Weighted avg} = GI * \frac{\text{Poor}}{\text{Total}} + GI * \frac{\text{Rich}}{\text{Total}} = 0 \times \frac{3}{10} + 0.694 \times \frac{7}{10} = 0.486$$

For parent,

$$GI(\text{Yes}) = 1 - (5/5)^2 = 0$$

$$GI(\text{No}) = 1 - [(1/5)^2 + (2/5)^2 + (1/5)^2 + (1/5)^2] = 1 - 0.28 = 0.72$$

$$\text{Weighted avg} = GI \times \frac{\text{Yes}}{\text{Total}} + GI \times \frac{\text{No}}{\text{Total}} = 0 + 0.72 \times \frac{5}{10} = 0.36$$

For weather,

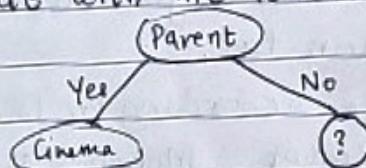
$$GI(\text{sunny}) = 1 - [(1/3)^2 + (2/3)^2] = 1 - (5/9) = 0.44$$

$$GI(\text{windy}) = 1 - (1^2 + 3^2 / 4^2) = 1 - 10/16 = 0.375$$

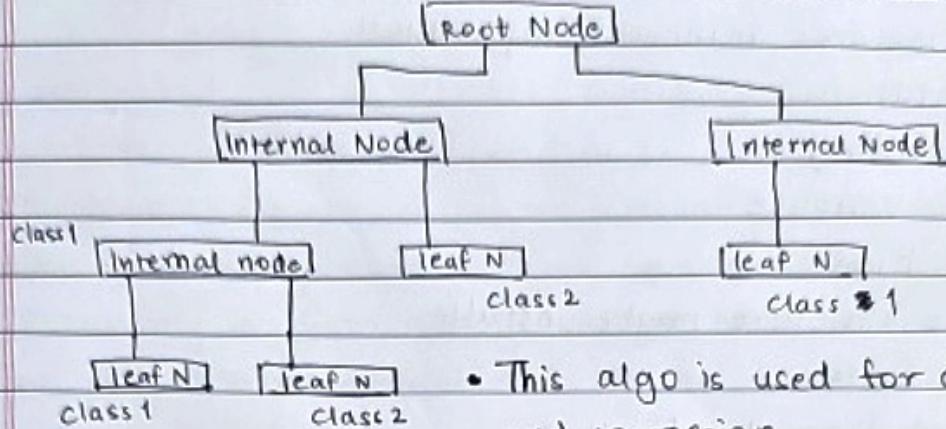
$$GI(\text{rainy}) = 0.44$$

$$\text{Weighted avg} = 0.44 \times 2 \times 3/10 + 0.375 \times 4/10 = 0.416$$

The attribute with the lowest Gini index will be considered as root



- CART Tree (Classification And Regression Tree)



- This algo is used for classification and regression.

- The best split point of each input is taken.
- Based on the best split point of each input in step 1, the new "Best" split point is identified.
- Split the chosen input according to the best split point.
- Continue the splitting until a stopping rule is satisfied or no further desirable splitting is available.

- CART for classification:

- It works by recursively splitting the training data into smaller and smaller subsets based on certain criteria. The goal is to split data in such a way that minimizes the impurity within each subset. For classification, CART uses Gini impurity.

- CART for Regression:

- Regression is an algo where target variable is continuous and tree is used to predict its value. CART for regression is a decision tree that creates a tree like structure to predict continuous target variable. It works by splitting training data recursively into smaller datasets based on specific criteria. The objective is to split the data in a way that minimizes residual reduction in each subset.

- Advantages:

- (1) Results ^{are} simplistic
- (2) C.ART can implicitly perform feature selection.

(3) Outliers have no meaningful effect on CART.

(4) It requires minimal supervision.

- Limitations

(1) Overfitting

(2) High Variance

(3) Low bias

(4) Tree structure maybe unstable.

- Use Cases

(1) Blood donor classification

(2) Environmental and Ecological data

(3) Financial applications.

CGPA	Interactive	Pract Knowledge	Comm skills	Job Offer
≥ 9	Yes	V Good	Good	Y
> 8	N	good	moderate	Y
≥ 9	N	avg	poor	N
≤ 8	N	avg	good	N
≥ 8	Y	good	mod	Y
≥ 9	Y	good	mod	Y
< 8	Y	good	poor	N
≥ 9	N	V Good	good	Y
≥ 8	Y	good	good	Y
≥ 8	Y	avg	good	Y

$Y=7$, $N=3$ (According to Job offer)

$$GI = 1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2 \Rightarrow GI(T) = 0.42 \text{ Imp}$$

CGPA Y N

$$\geq 9 \quad 3 \quad 1 \quad GI(T, CGPA \in (\geq 9, \geq 8)) = 1 - \left(\frac{7}{8}\right)^2 - \left(\frac{1}{8}\right)^2 = \frac{0.2187}{256} = 0.0861$$

$$\geq 8 \quad 4 \quad 0 \quad GI(T, CGPA \in (< 8)) = 1 - \left(\frac{2}{2}\right)^2 = 0$$

$$\leq 8 \quad 0 \quad 2 \quad \text{GII}(T, \text{CGPA}) = \frac{|S_1|}{|T|} GI(S_1) + \frac{|S_2|}{|T|} GI(S_2)$$

$$= \frac{8}{10} \times 0.2187 + 0 \times \frac{2}{10}$$

$$GII(T, CGPA) = 0.17496$$

$$\text{Now, } GI(T, \text{CGPA} \in \{>9, \leq 9\}) = 1 - (3/6)^2 - (3/6)^2 = 0.5$$

$$GI(T, \text{CGPA} \in \{>8\}) = 1 - (4/4)^2 = 0$$

$$GI(T, A - \text{CGPA}) = \frac{6}{10} \times 0.5 + \frac{4}{10} \times 0 = [0.3]$$

$$\text{Now, } GI(T, \text{CGPA} \in \{>8, \leq 8\}) = 0.44$$

$$\{>9\} = 0.375$$

$$\therefore GI(T, A - \text{CGPA}) = [0.414]$$

So, first GI is minimum. Subtract from GI(T).

$$\therefore \Delta Gini = GI(T) - GI(T, \text{CGPA}) = 0.42 - 0.1755 = [0.244] \quad (1)$$

Now wrt interactive

Interactive	Y	N	$GI(T, \text{Interactive} \in \{\text{Yes}, \text{No}\})$
Yes	5	1	$= 1 - (5/6)^2 - (1/6)^2 = 0.27$
No	2	2	$GI(T, \{\text{No}\}) = 1 - (2/4)^2 - (2/4)^2 = 0.5$

$$GI(T, A - \text{interactive}) = \frac{6}{10} \times 0.27 + \frac{4}{10} \times 0.5 = [0.362]$$

$$\Delta Gini(\text{Inter}) = 0.42 - 0.362 = [0.058] \quad (2)$$

Now wrt Pract knowledge

Prac	Y	N	$GI(T, \{\text{Vgood}, \text{good}\}) = 1 - (6/7)^2 - (1/7)^2 = 0.245$
Vgood	2	0	$GI(T, \{\text{avg}\}) = 1 - (1/3)^2 - (2/3)^2 = 0.44$
good	4	1	$GI(T, A - \text{Prac}) = \frac{7}{10} \times 0.245 + \frac{3}{10} \times 0.44 = [0.3035]$
avg	1	2	

$$GI(T, \{\text{A}, \text{avg}\}) = 0.4688 \quad \left. \begin{aligned} GI(T, A - \text{Prac}) &= \frac{8}{10} \times 0.4688 = [0.375] \\ GI(T, \{\text{Vgood}\}) &= 0 \end{aligned} \right\}$$

$$GI(T, \{\text{Vg}, \text{avg}\}) = 1 - (3/5)^2 - (2/5)^2 = 0.48 \quad \left. \begin{aligned} GI(T, \{\text{A}\}) &= [0.4] \\ GI(T, \{\text{A}\}) &= 0.32 \end{aligned} \right\}$$

$$\Delta Gini(\text{Prac}) = 0.42 - 0.3035 = [0.116] \quad (3)$$

Now wrt Comm skills,

Comm	Y	N	$GI(T, \{\text{M}, \text{P}\}) = 1 - (3/5)^2 - (2/5)^2 = 0.48$
G	4	1	$GI(T, \{\text{G}\}) = 1 - (4/5)^2 - (1/5)^2 = 0.32$
M	3	0	$GI(T) = [0.4]$
P	0	2	

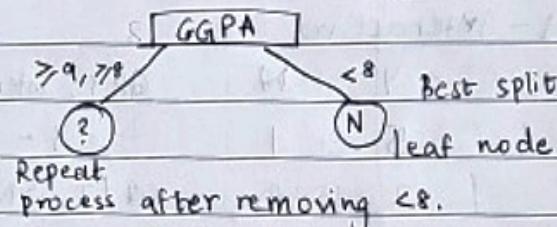
$$GI(T, \{G, M\}) = 1 - (7/8)^2 - (1/8)^2 = 0.218 \quad \left. \begin{array}{l} GI(1) = 0.1744 \\ GI(T, \{P\}) = 0 \end{array} \right\}$$

$$GI(T, \{G, P\}) = 1 - (4/7)^2 - (3/7)^2 = 0.489 \quad \left. \begin{array}{l} GI(T) = 0.342 \\ GI(T, \{M\}) = 0 \end{array} \right\}$$

Lowest is 0.1744

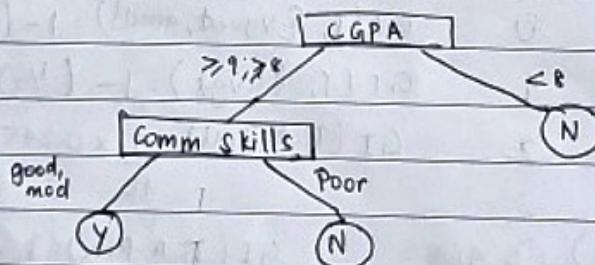
$$\Delta Gini(\text{Comm}) = 0.42 - 0.1744 = 0.2456 \quad \text{--- (4)}$$

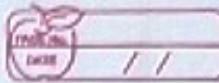
Now taking ~~the~~ highest $\Delta Gini$ as root. Note: Comm and CGPA are almost same so take any.



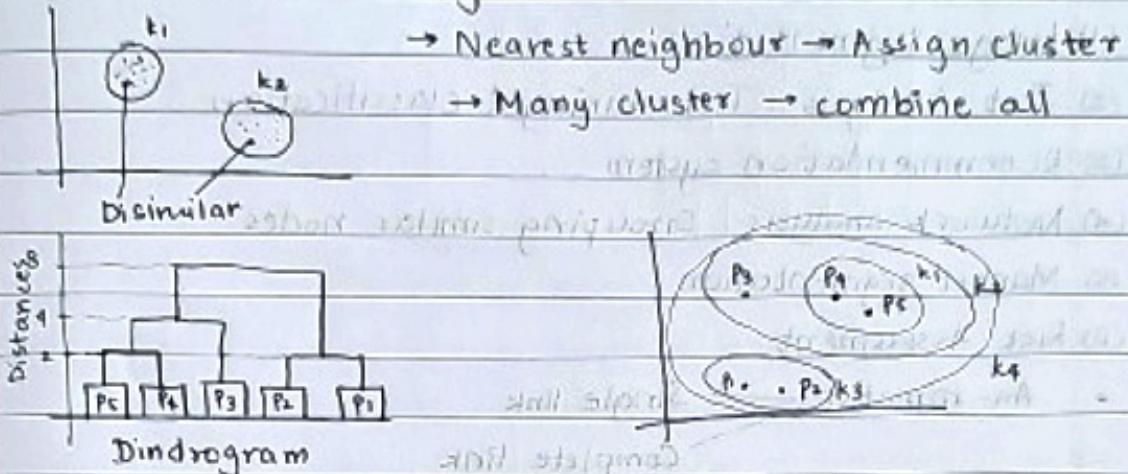
Attribute	GI	Gini	
Interactive	0.056	0.1624	direct answers calculated by mnm
Pract.	0.125	0.0934	
Comm	0	0.2184	

So take comm as root,





- Hierarchical Clustering:



Above is an example of Bottom to Top (One data pt to entire cluster)

Hierarchical
Agglomerative Dendroclustering (not in syllabus)

Bottom Top approach

Top Bottom Approach

(One cluster broken into data pts)

• Merits:

- (1) It is simple to implement & easy to interpret
- (2) It is robust, means no predefined cluster k , so given hierarchical structure is reliable & accurate
- (3) Algo is flexible. It is suitable to any type of data.
- (4) This algo is versatile, it can be used for both supervised & unsupervised learning.

• Demerits:

- (1) It is very sensitive to outliers, means if such outliers are present, hierarchical structure may not be accurate.
- (2) Computationally expensive, no of datapoints are more then it is time consuming
- (3) Not suitable for multidimensional data.
- (4) Cannot handle categorical variables as it cannot convert categorical data to numerical.

- Use cases:

- (1) Image Segmentation
- (2) Text Analysis : Text mining & classification
- (3) Recommendation system
- (4) Network analysis: Grouping similar nodes
- (5) Market segmentation
- (6) Risk Assessment

- Agglomerative
 - Single link
 - Complete link

(Q) Show the hierarchical clustering with single link on given dataset.

→	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0					
P ₂	3	0				
P ₃	8	9	0			
P ₄	9	8	1	0		
P ₅	5	10	6	7	0	
P ₆	4	9	7	8	2	0

If these pts are given, just ignore. Focus on lower triangle.

Bcoz of bottom up approach, see last row, i.e P₆. It has minimum distance from P₅, so make a cluster P_{5,6}.

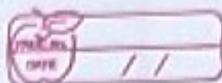
	P ₁	P ₂	P ₃	P ₄	P _{5,6}
P ₁	0				
P ₂	3	0			
P ₃	8	9	0		
P ₄	9	8	1	0	
P _{5,6}	4	9	6	7	0

Taking min distance of points from cluster.

$$\text{Eg: } P_1, P_{5,6} = 5, \quad P_1, P_6 = 4$$

$$\text{So, } P_1 (P_{5,6}) = 4.$$

Now see row P₄. Minimum is P₃ & P₅ so merge.

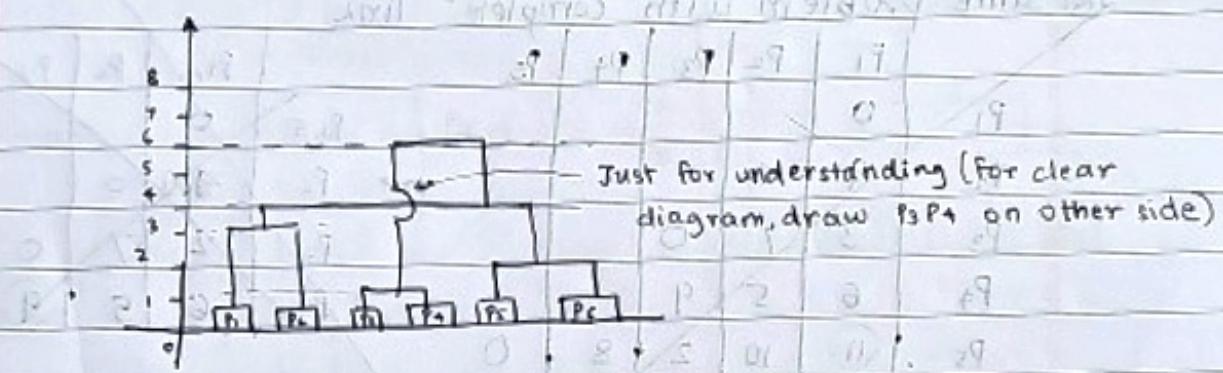


	P_1	P_2	$P_{3,4}$	$P_{5,6}$		$P_1 P_2$	$P_3 P_4$	$P_5 P_6$
P_1	0						0	
P_2	(3)	0					8	0
$P_{3,4}$	8	8	0				4	6
$P_{5,6}$	4	9	6	0				

S6 now merge $P_1 P_2 P_3 P_4$

	$P_{1,2,5,6}$	$P_{3,4}$
$P_{1,2,5,6}$	0	minimum values and unique val -
$P_{3,4}$	6	0

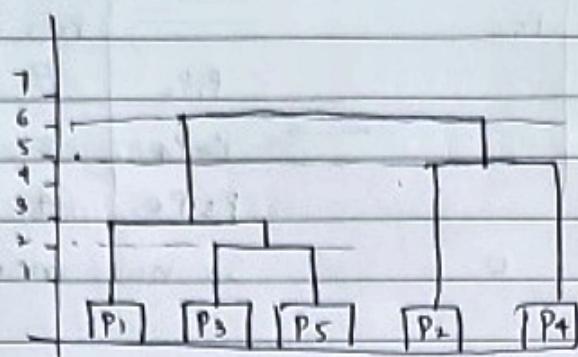
will diagrams dti u making min? in



(Q) Apply the agglomerative algo and represent the dendrogram on given dataset with single link.

	P_1	P_2	P_3	P_4	P_5		P_1	P_2	$P_3,5$	$P_4,6$
P_1	0						0	0		
P_2	9	0						9	0	
P_3	3	7	0						3	7
P_4	6	5	9	0					6	5
P_5	11	10	(2)	8	0				11	8

	P_1	$P_{2,4}$	$P_{3,5}$		$P_{1,3,5}$	$P_{2,4}$
P_1	0					
$P_{2,4}$	6	0				
$P_{3,5}$	(3)	7	0			



- In complete link, consider maximums instead of minimums in same process as above.

(a) Same problem with complete link

	P_1	P_2	P_3	P_4	P_5		$P_{1,5}$	P_2	P_3	P_4
P_1	0						0			
P_2	9	0					9	0		
P_3	3	7	0				2	7	0	
P_4	6	5	9	0			6	5	9	0
P_5	11	10	2	8	0					

	$P_{1,5}$	P_2	$P_{3,4}$		$P_{1,5}$	$P_{2,3,4}$	
$P_{1,5}$	0				0		
P_2	9	0			9		
$P_{3,4}$	2	5	0		5		

	$P_{1,5}$	P_2	P_3	P_4		$P_{1,5}$	P_2	$P_{3,4}$	
$P_{1,5}$	0					0			
P_2	10	0				10	0		
P_3	3	7	0			3	7	0	
P_4	8	5	9	0		8	7	0	

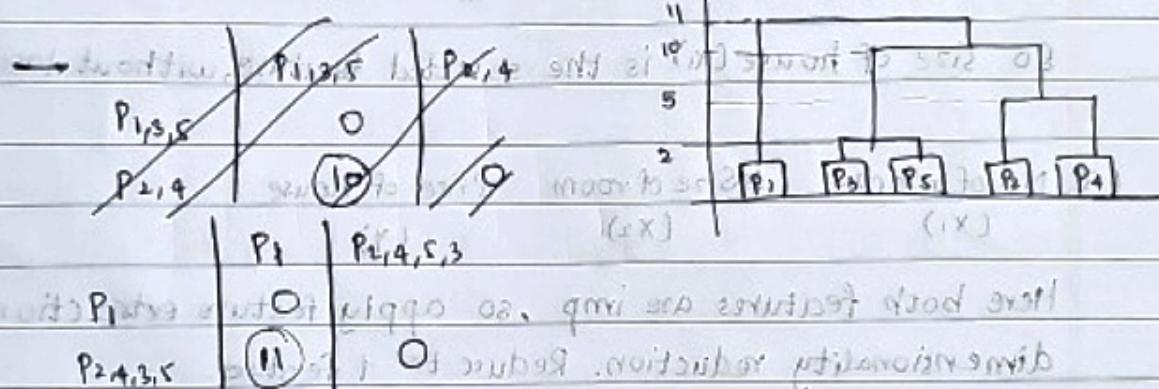
	$P_{1,3,4,5}$	P_2
	0	10

- For complete link, consider same clusters, i.e., you have to merge pts with min. distance between them. But, while merging consider max distances.

(a) Same with complete.

	P_1	P_2	$P_{3,5}$	P_4	Result for 3/2		P_1	$P_{2,4}$	$P_{3,5}$	Result for 3/2
P_1	0						P_1	0		
P_2	9	0					P_2	9	0	
$P_{3,5}$	11	10	0		Result for 3/2		$P_{3,5}$	11	10	Result for 3/2
P_4	6	5	9	0	Result for 3/2					

We take min = 5 and merge.



- Reducing Curse of Dimensionality can be done in 2 ways:

- Feature selection: Take only imp features highly correlated with target
- Feature extraction: Transform the given feature set to a new, improved feature set.

- Advantages of Dimension Reduction:

- Reduces COD
- Improves performance of model
- Visualize the data, i.e., understanding the data.

- $\text{cov}(x,y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$ is used to quantify relationship.

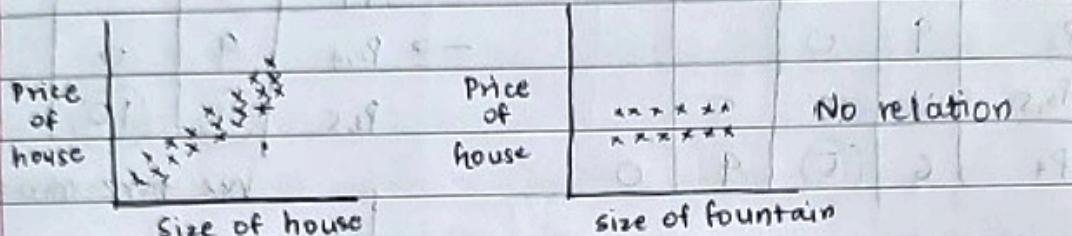
$$\text{cov}(x,y) = 1 \text{ linear } \begin{cases} \nearrow & \text{if } x \approx 0 \\ \downarrow & \end{cases} \Rightarrow \text{No relationship}$$

$$= -1 \text{ linear } \begin{cases} \nearrow & \\ \downarrow & \end{cases}$$

- Correlation \rightarrow Pearson correlation

$$= \frac{\text{cov}(x, y)}{\sqrt{6x^* 6y}} = -1 \text{ to } 1$$

(Q) Size of house (x_1) size of fountain (x_2) Price of house (y)



So size of house (x_1) is the selected feature, without loss of info.

(Q) No of rooms Size of room Price of house
 (x_1) (x_2) (y)

Here both features are imp, so apply feature extraction, i.e., dimensionality reduction. Reduce to 1 feature.

$2f \rightarrow \text{old f set}$



Transform



$D_f \rightarrow \text{New f set}$



Size of house



Predict y without loss

Principal Component Analysis ($x - \bar{x}$) \cdot $(x - \bar{x})^T$

(i) Calculate mean for each

$$\bar{x}_i = \frac{1}{N} \sum_{j=1}^N (x_{i1} + x_{i2} + \dots + x_{iN})$$

(2) Calculate Cov

$$\text{Cov}(x_i, x_j) = \frac{1}{N-1} \sum_{k=1}^N (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)$$

(3) Make Cov-matrix

(4) Calculate eigenvalues and eigen vectors of cov matrix.

$\det(S - \lambda I) = 0$ Polynomial eqn of degree n has n real roots which are eigenvalues of S.

$$(x_{11} - \lambda)(x_{12} - \lambda) \cdots (x_{1N} - \lambda) = 0$$

$$[(x_{11} - \lambda)(x_{12} - \lambda) + (x_{12} - \lambda)(x_{22} - \lambda) + \cdots + (x_{1N} - \lambda)(x_{NN} - \lambda)] = 0$$

$$[(x_{11} - \lambda)^2 + (x_{12} - \lambda)(x_{21} - \lambda) + (x_{13} - \lambda)(x_{31} - \lambda) + \cdots + (x_{1N} - \lambda)(x_{N1} - \lambda)] = 0$$

$$[(x_{11} - \lambda)^2 + (x_{12} - \lambda)(x_{21} - \lambda) + (x_{13} - \lambda)(x_{31} - \lambda) + \cdots + (x_{1N} - \lambda)(x_{N1} - \lambda)] = 0$$

$$\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} = \begin{bmatrix} (x_{11}, x) \text{ v.o} & (x_{12}, x) \text{ v.o} \\ (x_{21}, x) \text{ v.o} & (x_{22}, x) \text{ v.o} \end{bmatrix}$$

$$\begin{bmatrix} 0 & \lambda \\ \lambda & 0 \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} - (\lambda x_{11}) \text{ v.o}$$

$$\begin{bmatrix} x_{11} - \lambda x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} =$$

$$(x_{11} - \lambda x_{11})(x_{22} - \lambda x_{22}) - (\lambda x_{11})(\lambda x_{22}) = 0$$

$$0 = 100 + 115\lambda - 5\lambda^2$$

$$\frac{210.00}{5\lambda} < \frac{5 + 8\lambda + 0.05}{5\lambda} = 1$$

Teacher's Signature:

(Q) F Ex 1 Ex 2 Ex 3 Ex 4

$$x_1 \quad (4, 8, 13, 7)$$

$$x_2 \quad 11 \quad 4 \quad 5 \quad 14$$

variance

(x, x) var

Derive new feature set from given data using PCA.

(1) Calculate mean

$$\bar{x}_1 = 8 \quad \text{mean in x1 dimension}$$

$$\bar{x}_2 = 8.5 \quad \text{mean in x2 dimension}$$

(2) Calculate cov

$$\begin{aligned} \text{cov}(x_1, x_2) &= \frac{1}{N-1} \sum_{k=1}^n (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j) \\ &= \frac{1}{3} [(4-8)(11-8.5) + (8-8)(4-8.5) + (13-8)(5-8.5) \\ &\quad + (7-8)(14-8.5)] \\ &= -11 \end{aligned}$$

$$\begin{aligned} \text{cov}(x_1, x_1) &= \frac{1}{3} [(4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2] \\ &= 14 \end{aligned}$$

$$\begin{aligned} \text{cov}(x_2, x_2) &= \frac{1}{3} [(11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2] \\ &= 23 \end{aligned}$$

$$(3) S = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) \end{bmatrix} = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

(4) Eigenvalues

$$\cancel{(S - \lambda I)} = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$\cancel{\begin{bmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{bmatrix}} = \begin{bmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{bmatrix}$$

$$\det(S - \lambda I) = (14-\lambda)(23-\lambda) - (11 \times 11)$$

$$= \lambda^2 - 37\lambda + 201 = 0$$

$$\lambda_1 = 30.3848, \lambda_2 = 6.615$$

(5) Compute eigenvectors of this A(3) to get two unit eigenvectors.

$$U = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \end{bmatrix} = (I - \lambda I)(U)$$

$$\therefore \begin{bmatrix} (14-\lambda)u_1 - 11u_2 \\ -11u_1 + (23-\lambda)u_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$(14-\lambda)u_1 - 11u_2 = 0$$

$$-11u_1 + (23-\lambda)u_2 = 0$$

$$\frac{u_1}{11} = \frac{u_2}{14-\lambda} = t$$

$$\therefore u_1 = 11t, u_2 = (14-\lambda)t$$

$$u_1 = \begin{bmatrix} 11 \\ 14-\lambda \end{bmatrix}$$

$$\|u_1\| = \sqrt{11^2 + (14-\lambda)^2} = 19.7348$$

$$e_1 = \begin{bmatrix} 11/\|u_1\| \\ (14-\lambda)/\|u_1\| \end{bmatrix} = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$$

$$e_2 = \begin{bmatrix} 0.83 \\ 0.5574 \end{bmatrix}$$

(5) Compute first PC1

$$e_1^T \mathbf{x} = [x_{1k} - \bar{x}_1 \quad x_{2k} - \bar{x}_2] \cdot [0.5574 \quad -0.8303] = \frac{(x_{1k} - \bar{x}_1)}{(81.4) \cdot 9} \cdot \frac{(x_{2k} - \bar{x}_2)}{(81.4) \cdot 9} = 0.5574(4-8) + 0.8303(11-8.5)$$

$$= 0.5574(-4) + 0.8303(2.5) = -2.230535/72.9 = -0.0308535$$

F	Fx1	Fx2	Fx3	Fx4	Fx5	Fx6	Fx7	Fx8	Fx9
PC1	-4.30535	3.73	5.69	-5.12387					

$$\alpha_{PC1} = 90^\circ - 86^\circ = 4^\circ$$

$$\cos(\alpha) = \frac{\cos(\theta) \cos(\phi) - \sin(\theta) \sin(\phi) \cos(\psi)}{\sqrt{1 - \sin^2(\theta) \sin^2(\phi) \cos^2(\psi)}}$$

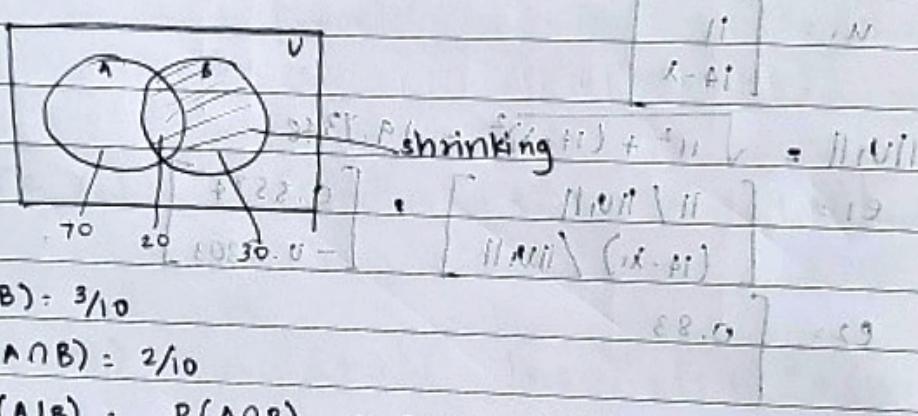
(Q) Apply the concept of PCA and calculate eigenvalues.

$$\begin{array}{cc}
 X & Y \\
 2.5 & 2.4 \\
 0.5 & 2.2 \\
 2.2 & 1.9 \\
 1.9 & 2.2 \\
 3.1 & 3.0
 \end{array}
 \quad \bar{x} = 2.04 \quad \bar{y} = 2.34$$

$$\text{cov}(X, X) = \frac{1}{4} [3.1752] = 0.938$$

$$\text{cov}(Y, Y) = \frac{1}{4} [10.158] = 0.2895$$

- Naive Baye's Theorem



$$P(B) = 3/10$$

$$P(A \cap B) = 2/10$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.2}{0.3} = 0.67$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{0.2}{0.67} = 0.299$$

$$(1 - P(A)) \cdot P(\bar{A}) = 1 - 0.67 = 0.33$$

$$P(A \cap B) = P(B) P(A|B) = P(A) P(B|A)$$

$$P(A|B) = \frac{P(A) P(B|A)}{P(B)} = \frac{0.67 \cdot 0.299}{0.3} = 0.47$$

$x_1, x_2, \dots, x_n \quad \left\{ \begin{matrix} i/p \\ f_1, f_2, \dots, f_n \end{matrix} \right. \quad \left\{ \begin{matrix} i/p \\ f_1, f_2, \dots, f_n \end{matrix} \right. \Rightarrow \text{Target}$

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y) P(x_2|y) \dots P(x_n|y)}{P(x_1) P(x_2) \dots P(x_n)}$$

$$= P(y) \prod_{i=1}^n P(x_i|y)$$

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

$$\therefore P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

$$y = \underset{\text{highest probability}}{\arg\max} P(y) \prod_{i=1}^n P(x_i|y)$$

(Q)	No	Outlook	Temp	Humidity	Play Tennis
	1	Sunny	High	High	N
	2	Sunny	High	High	N
	3	Overcast	High	High	Y
	4	Rain	Mild	High	Y
	5	Rain	Cool	Normal	?
	6	Rain	Cool	Normal	N
	7	Overcast	Cool	Normal	Y
	8	Sunny	Mild	High	N
	9	Sunny	Cool	Normal	Y
	10	Rain	Mild	Normal	Y
	11	Sunny	Mild	Normal	Y
	12	Overcast	Mild	High	Y
	13	Overcast	Hot	Normal	Y
	14	Rain	Mild	High	N

From the given dataset, what is probability of person to play tennis if climate is sunny and hot?

$$\rightarrow P(\text{sunny}) = 2/14 \text{ Outlook}$$

	Y	N	P(Y)	P(N)	(N)
Sunny	2	3	2/14	3/14	
Overcast	4	0	4/14	0/14	
Rain	3	2	3/14	2/14	
Overcast					
Total	9	5	100%	100%	



Temp:

	Y	N	$P(Y)$	$P(N)$
Hot	2	2	$2/4$	$2/5$
Mild	4	2	$4/9$	$2/5$
Cool	3	1	$3/9$	$1/5$
Total	9	5	100%	100%

Play Tennis?

	$P(Y) \cdot P(N)$	Actual	Actual
Yes	$9/14$	Actual	Actual
No	$5/14$	Actual	Actual
Total	100%	Actual	Actual

$$P(Y | \text{sunny, hot}) = P(\text{sunny} | \text{yes}) P(\text{hot} | \text{yes}), P(Y)$$

$$= \frac{2}{9} \times \frac{2}{4} = \frac{1}{9}$$

Denominator is always constant (so neglect).

$$P(Y | \text{sunny, hot}) = 0.0317$$

$$P(Y | \text{sunny, hot}) = 0.0857$$

Normalize,

$$\frac{P(Y)}{P(Y) + P(N)} = \frac{0.0317}{0.0317 + 0.0857} = 0.2656 \approx [0.27]$$

$$P(N) = 1 - 0.27 = [0.73]$$

• Naive Bayes

- Naive Bayes Classifier: How to handle zero probability

(a)

TEXT

- 1 A great game
- 2 The election was over
- 3 Very clean match
- 4 A clean but forgettable game
- 5 It was a close election

Category

Sports

Non sports

Sports

Non sports

Non sports

Classify the given sentence ("A very close game") using Naive Bayes Classifier.

Total no of words in sports = 11

Total no of words in non-sports = 9

Total no of unique words = 14

So we need to consider probability of each word in given sentence.

Sports:

$$P(A) = 2/11 \quad P(\text{very}) = 1/11 \quad P(\text{close}) = 0/11 \quad P(\text{game}) = 2/11$$

So $P(\text{close}) = 0$, this zero probability is handled by Laplace's smoothing:

$$\Theta_i = \frac{x_i + \alpha}{N + \alpha d}$$

$\Theta_i = P(\text{word})$

$x_i = \text{Word count}$

α is ~~0~~ N: Total no of words

d: No of unique words

α is > 0 but always $\alpha = 1$ constant.

$$P(A) = \frac{2+1}{11+14} = \frac{3}{25}$$

$$\frac{2}{23}$$

$P(W/\text{Non sports})$

$$P(\text{very}) = \frac{1+1}{11+14} = \frac{2}{25}$$

$$\frac{1}{23}$$

$$= \frac{2 \times 1 \times 2 \times 1}{(23)^4}$$

$$P(\text{close}) = \frac{0+1}{11+14} = \frac{1}{25}$$

$$\frac{2}{23}$$

$$= \frac{1^{1+4}}{(23)^4} \times 10^{-5}$$

$$P(\text{game}) = \frac{2+1}{11+14} = \frac{3}{25}$$

$$\frac{1}{23}$$

$P(W/\text{Sports})$

$P(W/\text{Non sports})$

$$P(W/\text{Sports}) = \frac{3 \times 2 \times 1 \times 3}{25 \times 25 \times 25 \times 25}$$

$$= \frac{0.46}{10^{-5}}$$