

Machine Learning & Data Science

→ Basics of ML, different from traditional programming, types of challenges

→ Visualisn, plotting
|
Numpy
Pandas
Matplotlib
sklearn

→ Statistics/statistical train, test

|
Validation
K-cross validn

Performance evaln for classificn

→ Classificn KNN

| Naive Bayes

Decision Tree → Ensemble

→ Regression → Mathematics learning.

Confidence &
Interval ANOVA

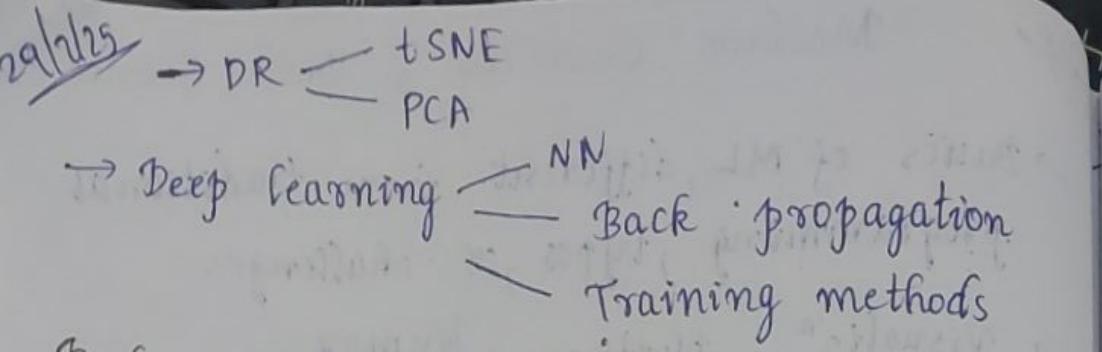
→ Regularization & Kernel method:

→ Clustering EM
Mixture models

|
K Means
DBSCAN

Hierarchical Clustering

Optics



- Book:
- ① Principles of DS
 - ② Hands on ML with scikit-learn, Keras, TensorFlow
 - ③ Data mining - Concepts & Techniques

Evaluation:

- Mid Term 25 M
- End Term 40 M
- Assignment 15 M
10 M
- 1 or 2 surprise tests 10 M

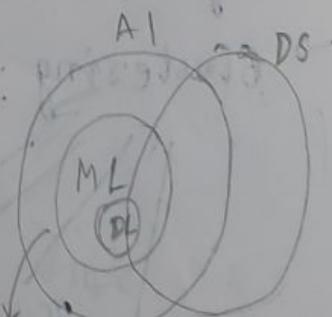
80% attendance is compulsory.

AI ML DS

AI - A mechanism that can do everything that a human does.

Less data - ML

Huge data, hidden characteristics - DL



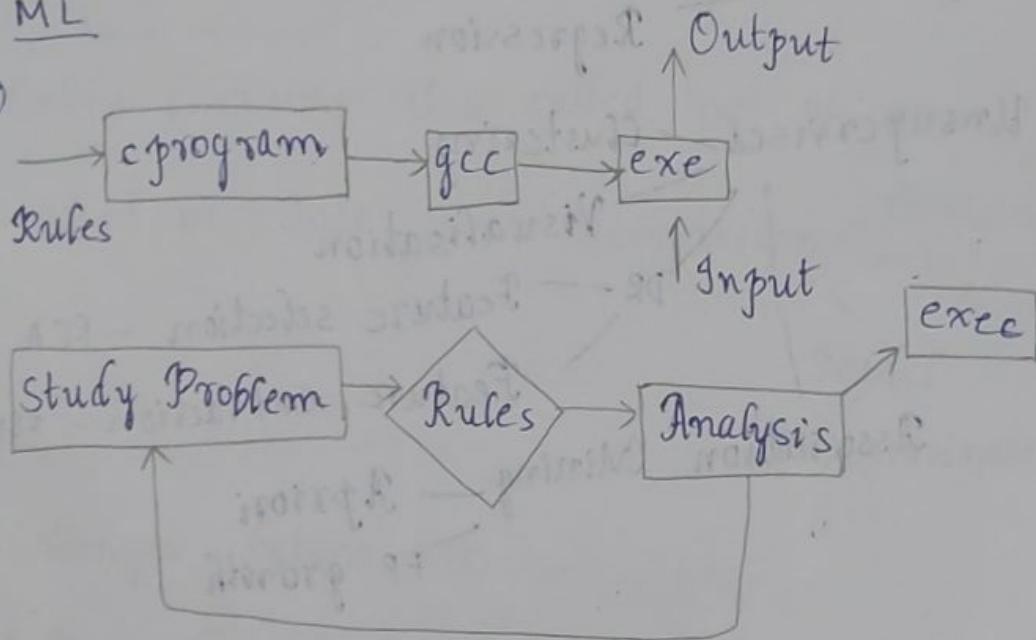
used for generating end product. (IoT sensors, etc)

29/7/25

DS → Visualisation
 DS → Missing → Data prep → Datasets
 DS → ML

ML

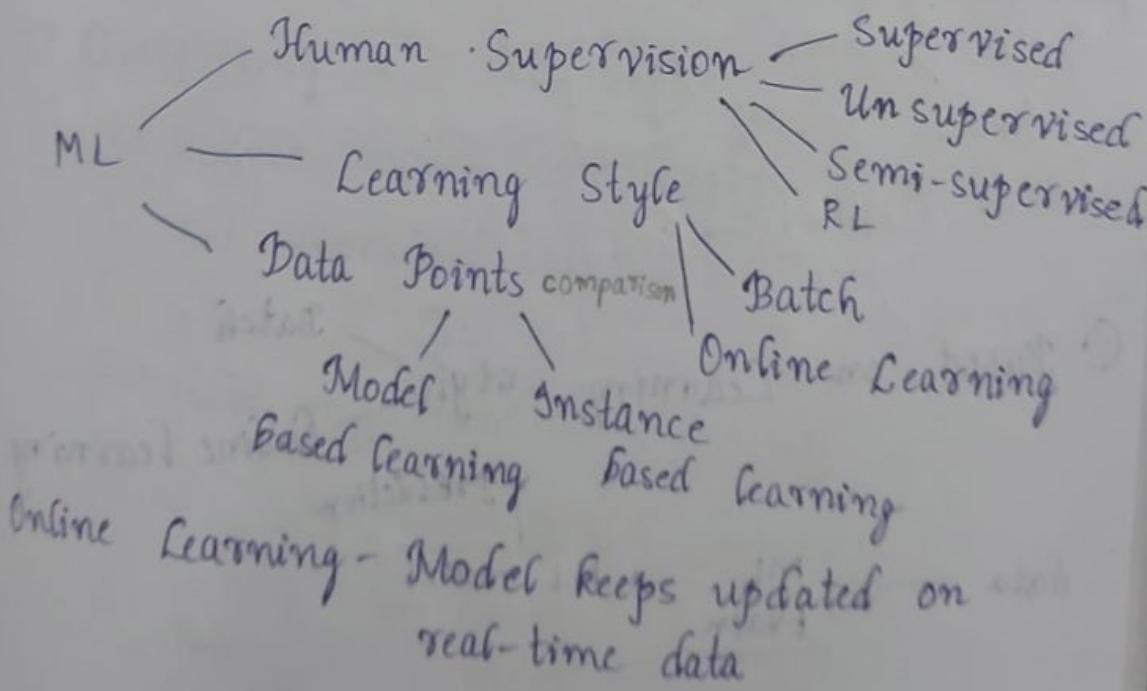
①



Does this involve any learning? No

→ Here, in C we will set our rules

→ In ML, machine sets its own rules.



Batch Learning - We need to give the data in 1 go

30/7/25 Once the training is completed, we cannot give any data again.

Supervised

Classification

Regression

Unsupervised

Clustering

Visualisation

DR Feature selection - PCA

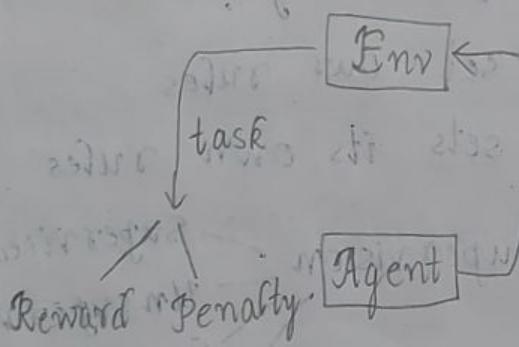
Feature extraction - t-SNE

Association Mining

Apriori

FP growth

Reinforcement : (No Human)



② Based on

Learning

style

Batch

Online Learning

↑ Prediction

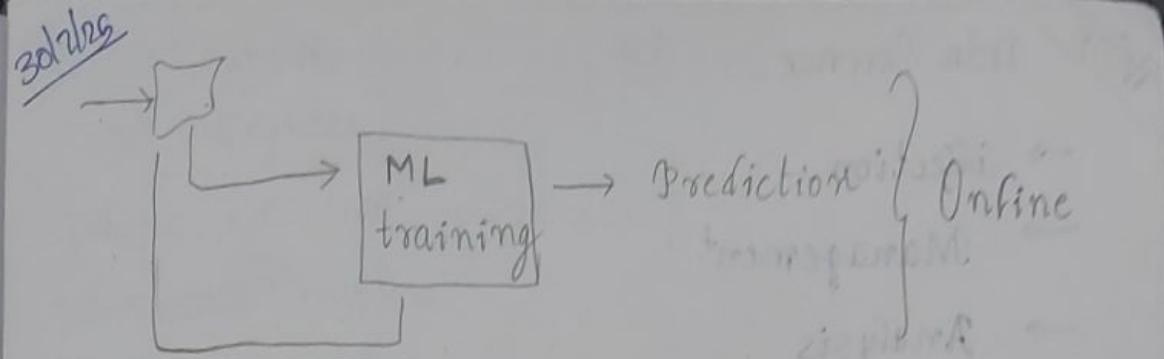
Data →

ML train

→ ML model

↑ Input

Batch



Online Learning also called "out of core L"

- ③ Based on Data points comparison
 - Instance (small dataset)
 - Model (Discriminative)

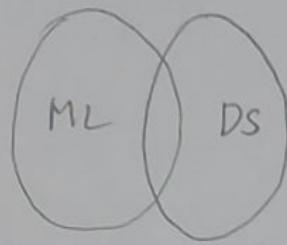
Challenges behind ML

- Data biased → Not represent
- Complex data → Low training accuracy
- Underfitting → Low (training + testing) accuracy
- Overfitting
 - ↳ Reduce model complexity
 - ↳ Set threshold value
 - ↳ Work on data set
 - ↳ Give important data/ feature

* Revise Pandas, Matplotlib

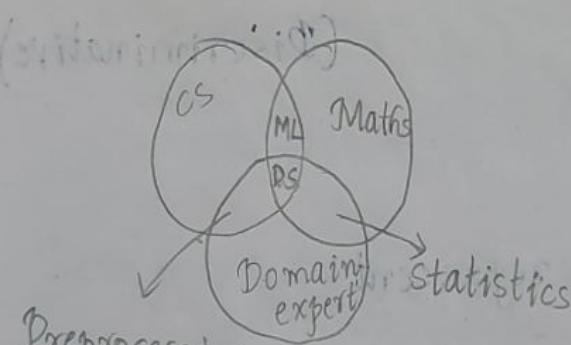
11/8/25 Data Science

- Collection
- Management
- Analysis



Domain experts collect the data.

Management - Store the data - cs experts
Analysis - Maths & statistics

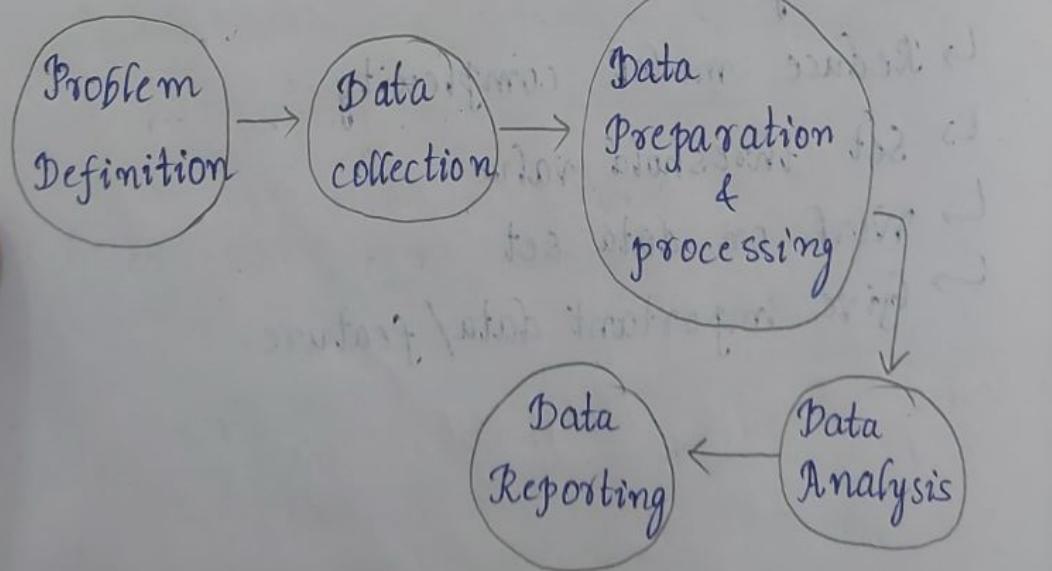


Preprocessing

Visualisation

Data Science Cycle

Descriptive Method - Mean, Median, Mode
Probabilistic method - Gaussian distribution



① We should have clear-cut objective } ①

18/28 → Experts } ①
→ Colleagues

Where we can collect data?

- Search engine
- Web scraping
- Surveys
- Open ended interviews

Missing values

Outliers

} ③

→ Statistical techniques

→ Probability

→ ANOVA

→ t-
p-test

→ Sample

} ④

Data Science Market

① CRISP DM process

(CRISP-DM Guide)

Business
understanding

Data
understanding

Data Prepⁿ

Data

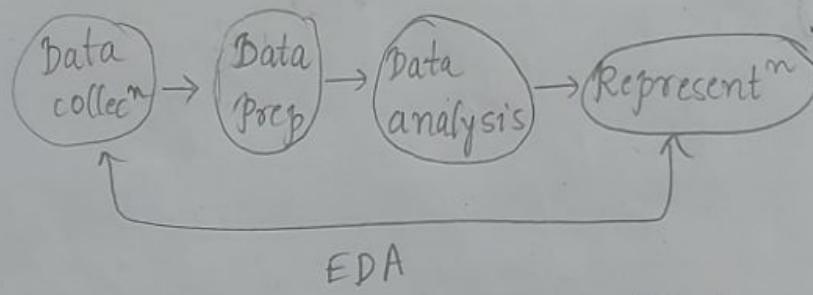
Modelling

Evalⁿ

- 18/25
- ① Business Understanding
 - Determine business objectives
 - Assess situation
 - Determine data mining goals
 - Produce project plan

② Data Understanding

- Collect initial data
- Describe data
- Explore data
- Verify data quality



③ Data Prepⁿ

- Select data
- Clean data
- Construct data
- Integrate data
- Format data

④ Modeling

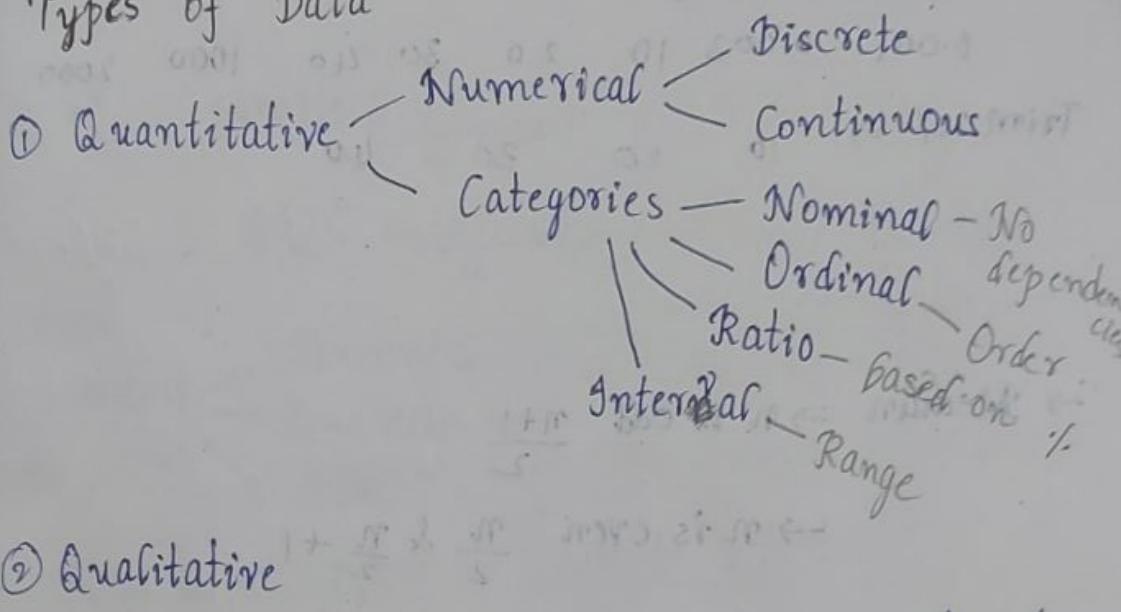
- Select modeling technique
- Generate test design
- Build model
- Assess model

Hyperparameter - We check & decide, which would give the best performance

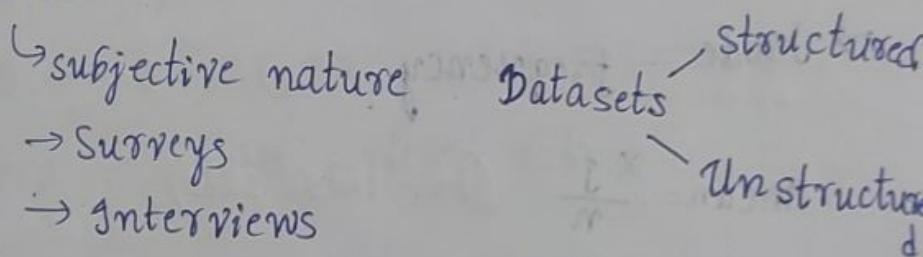
⑤ Deployment

- Plan deployment
- Plan monitoring & maintenance
- Produce Final Report
- Review project

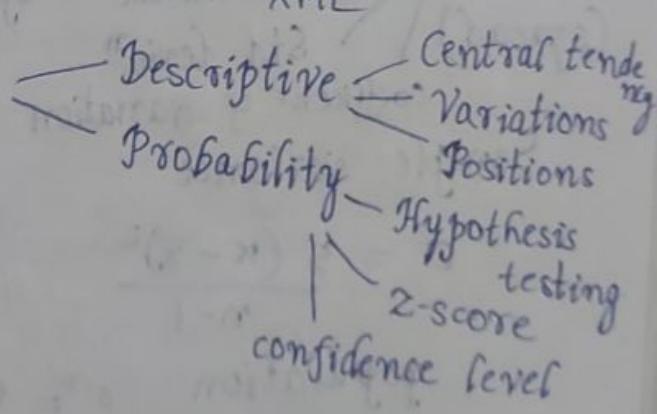
Types of Data



② Qualitative



Statistical Analysis



Descriptive

① Central tendency

If we are taking whole dataset, it is

5/8/25

represented by N. If we take a small sample, n.

N → population size μ mean $\mu = \frac{\sum x}{N}$

n → sample size \bar{x} $\bar{x} = \frac{\sum x}{n}$

① Central Tendency → Mean → Trimmed Mean

0.001 0.003 10 20 30 40 1000 2000

Trim 20% → 10 20 30 40

(Will be removed from lower value & higher value)

→ Median → n is odd $\frac{n+1}{2}$

→ n is even $\frac{n}{2}$ & $\frac{n}{2} + 1$

→ Mode - frequency

mean = $\frac{x \cdot f}{n}$ (in this case)

variation → Range → Max-Min (If outliers are present, it will get affected & we are not counting every value)

(spread) Variance → spread on mean
Std deviaⁿ

Sample variance

$$s^2 = \frac{\sum (x - \bar{x})^2}{n-1} \rightarrow \text{HW}$$

population variance $\sigma^2 = \frac{\sum (x - \mu)^2}{N}$

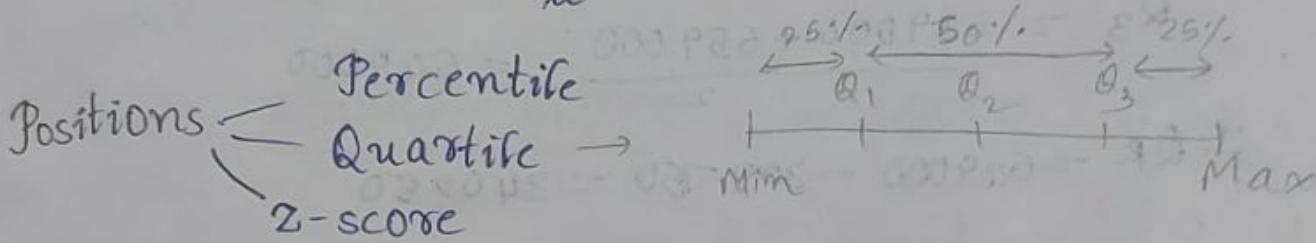
$$SD = s = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

If mean are different: (when we take from 2 datasets)

CV: sample = $\frac{s}{\bar{x}} \times 100\%$

popⁿ = $\frac{\sigma}{\mu} \times 100\%$



Find Q_1, Q_2, Q_3 .

5.4, 6.0, 6.3, 6.8, 7.1, 7.2, 7.4, 7.5, 7.9, 8.2, 8.7

→ We need to have ordered dataset

~~$25 = \frac{x}{11} \Rightarrow x = \frac{25 \times 11}{100} = \frac{11}{4}$~~

→ Take middle value → Q_2

$$Q_1 = 6.3$$

$$Q_3 = 7.9$$

Take mean → $Q_1 + Q_3$

IQR - Inter Quartile Range

Where my 50% of data is there?

$$IQR = Q_3 - Q_1$$

$$\text{Lower Bound} = Q_1 - (1.5 \cdot IQR)$$

$$\text{Upper Bound} = Q_3 + (1.5 \cdot IQR)$$

~~Ans~~
389950, 230500, 158000, 479000, 639000,
114950, 5500000, 387000, 659000, 529000,
575000, 488800, 1095000

114950, 58000, 230500, 387000, 389950, 479000,
488800, 529000, 575000, 639000, 1095000
659000
6500000

$$Q_2 = 488800$$

$$Q_1 = \frac{230500 + 387000}{2} = 308750.$$

$$Q_3 = \frac{639000 + 659000}{2} = 649000$$

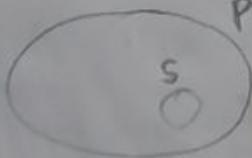
$$IQR = 649000 - 308750 = 340250$$

$$(1.5)(IQR) = (1.5)(340250) = 510375$$

Lower bound = $Q_1 - (1.5)(IQR) = -$

Written by me
Why $n-1$? (in sample variance)

In the case of skewed data, experiments have shown that, having $n-1$ gives better results. That is approximately equal to population variance σ^2 .



$$s^2 \approx \sigma^2$$

With the help of sample variance, we should estimate the population variance

Main: Degree of freedom concept

6/8/23 Z-score

$$Z\text{-score} = \frac{x-\mu}{\sigma}$$

Z can be
 $-ve \rightarrow$ below mean
left side
 $0 \rightarrow$ eq
 $+ve \rightarrow$ greater than μ
 $(\mu + \sigma) \quad (\mu + 2\sigma)$
means!

Suppose the x -value of house price is 270000, Mean is 350000, std is 40000, then calculate the z -score

$$z = \frac{270000 - 350000}{40000} = -\frac{80000}{40000} = -2$$

Probability Distribution

Binomial Distribution

↳ two success, failure

→ independent trials,

prob should not be changed

→ $x \rightarrow$ Discrete value

$$P(x) = \binom{n}{x} p^x (1-p)^{n-x}$$

Discrete $\begin{cases} \text{Binomial} \\ \text{Poisson} \end{cases}$

Continuous $\begin{cases} \text{Normal} \\ \text{Gaussian} \end{cases}$

A sample of 20 patients is selected for

6 (8) surgery, Calculate the prob that 18 out of 20 patients will be having the successful surgery. From the past data, prob of success is 92%.

$$\text{Sol: } p = 0.92, 1-p = 0.08$$

$$\begin{aligned} P(x) &= {}^{20}C_{18} \left(\frac{92}{100}\right)^{18} \left(\frac{8}{100}\right)^2 \\ &= \frac{20 \times 19 \times 18}{1 \times 2} \left(\frac{92}{100}\right)^{18} \left(\frac{8}{100}\right)^2 \\ &= 190 \left(\frac{92}{100}\right)^{18} \left(\frac{8}{100}\right)^2 \\ &= 0.271 \end{aligned}$$

Poisson Distribution

Each n every thing will be dependent on whole dataset

μ, σ

$$P(x) = \frac{\mu^x e^{-\mu}}{x!}$$

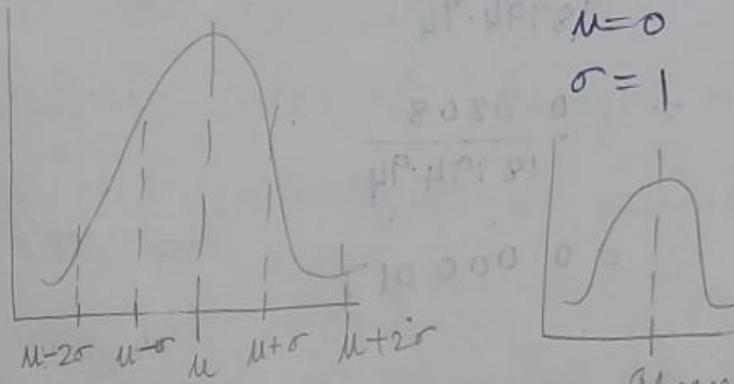
- (Q) From the past data, a traffic engineer demands the avg. no. of vehicles entering a biking garage, a 7 per 10 min period. Calculate the prob that the no. of the vehicles entering in a parking is 9 in certain 10 min period.

6/8/2023 Ans $e = 2.718$ Ans 0.101

$$P(x) = \frac{(-2.718)^7}{9!}$$

PMF - Probability Mass fn

Normal Distribution



Mean
Median } overlap
Mode }

$$x \sim N(\mu, \sigma^2)$$

Mean \rightarrow Std

$$P(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

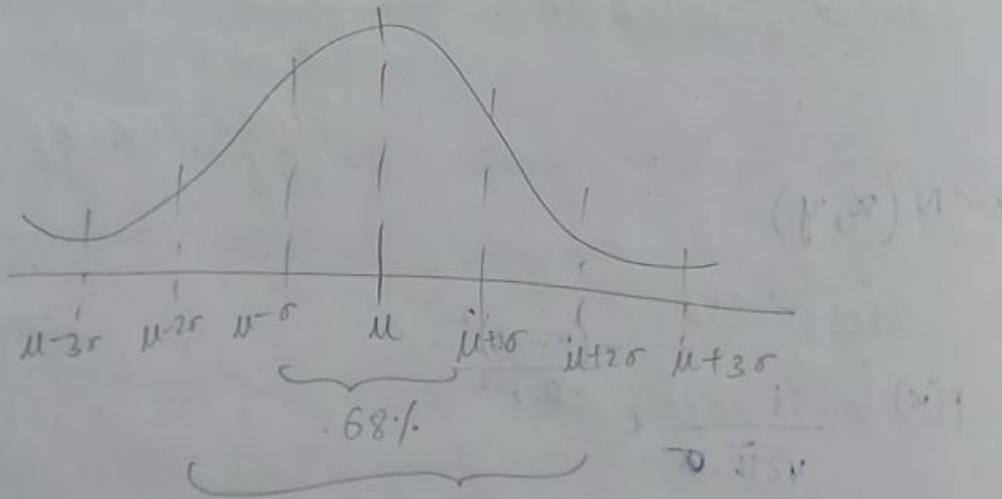
Each n every normal distribution will be standardized by z-score.

- b) At a software company, the mean of employee's salary is 60000 with a std dev 7500, calculate the prob that a random employee earns more than 68000.

6th Ques

$$\begin{aligned}
 P(68000) &= \frac{1}{\sqrt{2\pi} \cdot 7500} e^{-\frac{(68000 - 60000)^2}{2 \cdot (7500)^2}} \\
 &= \frac{1}{18794.94} (2.718) \frac{64 \times 10^6}{5625 \times 10^4} \\
 &= \frac{1}{18794.94} (2.718) \frac{6400}{5625} \\
 &= \frac{1}{18794.94} (2.718)^{-1.137} \\
 &= \frac{0.3208}{18794.94} \\
 &= 0.000017
 \end{aligned}$$

Ans: 14%



- Q) Suppose an automated designer is interested in designing automatic seats to accomodate different heights of 95%. The heights of customers follow normal distribution with the mean value 68 ft.

6/8/25 Std 3. For wt range of heights, should the designer model the car seats
NEXT CLASS: CONFIDENCE LEVEL, ANOVA, HYPOTHESIS TESTING

Inferential Statistics

We take some observⁿ from sample & use it for populⁿ.

→ We find out error also

Confidence interval → Margin of error

Any statistic → Point of estimate $\hat{p}_{oe} \pm e$
Confidence Interval

Lower bound of CI = $\hat{p}_{oe} - ME$ ↑ confidence level

Upper bound of CI = $\hat{p}_{oe} + ME$ ↓ sample mean is

A sample of 100 patients is taken. Determined at 64 years old. Assume that the corresponding MOE for a 95% CI is calculated to be 4 yrs.

$$= 64 - 4 \rightarrow 60 \text{ to } 68$$

$$= 64 + 4$$

Mean is going to lie btw CI 60-68

Mean → Assumption, $\mu \rightarrow$ populⁿ mean

Proportion $\mu_x \rightarrow$ sample mean

$$\mu_x = \mu$$

$$S_x = \frac{\sigma}{\sqrt{n}} \quad n \rightarrow \text{sample size}$$

Central Limit Theorem → When we are

~~1382~~ going to aggregate the samples, it will form normal distribution.

We might get skewed graph.

2 cond^{ns}: ① We don't know dataset distribution

We know mean, std

If sample size > 30. | for σ in shew

② Std is not known

CL
→ Mean $\begin{cases} \sigma \text{ known} \rightarrow z\text{-test} \\ \sigma \text{ Not known} \rightarrow t\text{-test} \end{cases}$

σ is known

Confidence level for mean = $z_c \frac{\sigma}{\sqrt{n}}$

Lower B of CI \rightarrow Mean - E_c

Upper B of CI \rightarrow Mean + E_c

Q) Sample size = 50, POE = Populⁿ mean = 86000
 σ = 9000. Distribⁿ isn't known. Calculate the mean & std for sample & comment the shape of distribⁿ for sample.

sol)
50
85000 } Populⁿ
90000
Sample

$$\begin{aligned} \text{Sample std dev } s_x &= \frac{\sigma}{\sqrt{n}} \\ &= \frac{9000}{\sqrt{50}} \\ &= \frac{9000}{5\sqrt{2}} \end{aligned}$$

~~13 shs~~ * According to CLT, if sample size > 30, then it follows normal distribution.

Find CI for 95% confidence, $z = 1.96$

~~error of margin~~ $CL = 1.96 \times \frac{9000}{\sqrt{50}} = \frac{1.96}{\sqrt{100}} \times \frac{9000}{5\sqrt{2}} = 98\sqrt{2} \times 18 = 1764\sqrt{2} \approx 730.2 = 2494.66$

Q) Suppose sample of 50 students is taken, in that 50 students, the mean of time spent on HW is 12.5 hrs. The popn std devⁿ is 6.3 hrs. Create a forecasted CI using CL 90% & 95% and provide a conclusion regarding these confidence level that which interval is wider.

$z\text{-score for } 90\% = 1.645, 95\% = 1.960$

Sol):- $ss = 50 \Rightarrow$ Normal distribⁿ

Mean = 12.5, $\sigma = 6.3$

CI Mean $= (1.645) \times \frac{6.3}{\sqrt{50}} = 1.4656$

CI Mean $(95\%) = (1.960) \times \frac{6.3}{\sqrt{50}} = 1.7462$

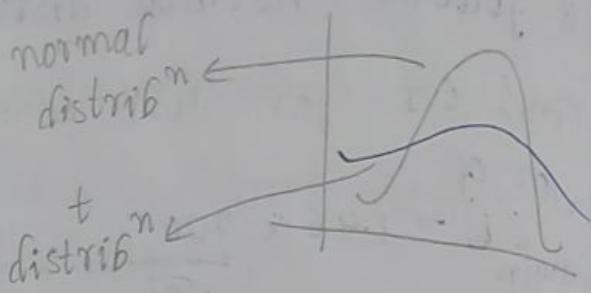
~~(90%) Lower B = 12.5 ± 1.465~~

~~(95%) Upper B = 12.5 ± 1.746 → Wider~~

→ Random Sample p → If σ is not known,
→ 30 or ND find out s_x

13/8/29

$$E_C = t_c \frac{S_x}{\sqrt{n}}$$



degree of freedom = $n - 1$

→ Check from table

Dof ... 98% 99%

$\frac{1}{2}$	2	3	4
1			

Q) Suppose in a company 5000 employees & there Administr^n decided to take sample of 16 employee to find mean salary. The sample data indicate, sample mean of salary is 15.8 lakhs & the std is 3.2 lakh. Calculate the 99% CL, CI for salary.

$$t_c = 2.947$$

~~Boole's approx~~

Sol): $E_C = (2.947) \frac{(3.2)}{\sqrt{16}} = 2.3576$

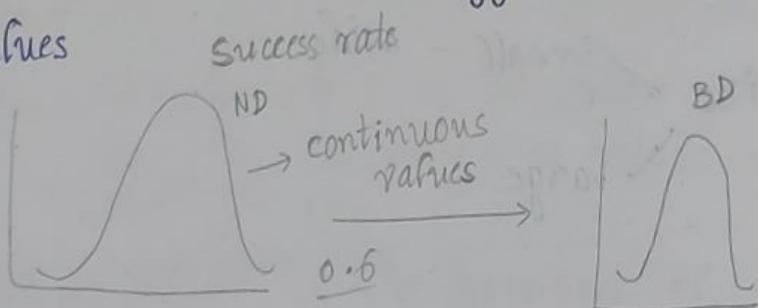
$$CI = 15.8 \pm 2.357$$

$$= [13.44, 18.16]$$

13/14 Confidence Interval for Proportions

$$\hat{p} = \frac{x}{n}$$

Normal distribⁿ is applicable for continuous values



$$\left[\begin{array}{l} np \\ n(1-p) \end{array} \right] \rightarrow \text{at least } 5$$

If it is greater than 5, Normal distribⁿ can represent binomial distribⁿ also.

$$E_c = z_c \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Q) A medical researcher wants to see the change in the proporⁿ of smokers from 5 yrs ago & the % of adult % was 28%. A random sample of 1500 adults is taken & 360 were found smokers. Calculate the 95% confidence interval for the true popⁿ portion of adults who smoke. $z_c = 1.960$

$$\hat{p} = \frac{360}{1500}$$

$$\hat{p} \pm E_c$$

$$E_c = (1.960) \sqrt{\frac{0.24(1-0.24)}{1500}} = 0.0005$$

$np \cdot n(1-p)$ should be atleast 5

~~13.8 hrs~~

Q) How can we determine the sample size if we take ~~some~~ small sample?

Sample size

- Small - wide
- Large -

Mean

$$E_c = z_c \frac{\sigma}{\sqrt{n}}$$

$$\begin{aligned} n &= \left(\frac{z_c \sigma}{E_c} \right)^2 \\ n &\geq 2 \end{aligned}$$

$$n = \left(\frac{z_c \sigma}{E_c} \right)^2$$

Proportion

$$E_c = z_c \sqrt{\frac{(1-p)p}{n}}$$

$$n = \left(\frac{z_c}{E_c} \right)^2 (1-\hat{p})\hat{p}$$

Q) A Benefits analyst ~~is~~ interested in 95% CL for _____ for mean salary of employee. What should be the ~~max~~ sample size is fixed. ~~population~~ size if margin of error 1000 is desired. $\sigma(\text{populn}) = 8000$

$$- (z_c = 1.645)$$

$$\Rightarrow n = \left(\frac{1.645 \times 8000}{1000} \right)^2 = \left(\frac{13160}{1000} \right)^2 = (13.16)^2$$

$$= 173.18$$

≈ 174 / whole number will be taken

~~13/8/25~~ 20) One candidate is Smith is

- (i) 1 candidate is Smith is planning a survey to determine 95% C.I. for the proportion of voters who planned to vote for Smith. How many should be surveyed. $E_c = 3\%$, $Z = 1.96$

Case-1: Assume there is no prior estimate for the 0.5 proportion of voters who vote for Smith default

Case-2: Based on the previous elections, 42% people voted for Smith

(i) $E_c = 3\%$, $Z = 1.96$

95% CI

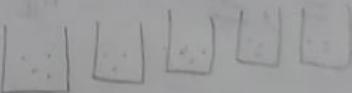
$$n = \left(\frac{1.96}{0.03} \right)^2 (0.95)(0.05)$$
$$= 1043$$

Bootstrap sampling

Mean $\rightarrow 30$, ND

Proportion $\rightarrow np$, $n(1-p) \rightarrow 8$

Bootstrap \rightarrow sample with replacement



What is the advantage of bootstrap sampling?

$\hookrightarrow 6C_1 \times 4C_1$

Sample without replacement

- (i) To increase population size
- (ii) Helps us use valid
- (iii) Helps remove bias in sample

stratified sampling:

If dataset contains 60% class A & 40% class B,

then sample should also contain 60% class A & 40% class B.

① Population → repeatedly generate sample

② \bar{x}, \bar{s}_x

③ Sort

④ 95%, $P_{2.5} - P_{97.5}$

⑤ 90%, $P_5 - P_{95}$

Hypothesis Testing

We were defining statistical parameters & we were mapping it with popⁿ "

(Previously)

Null H_0

$\hat{=}, \leq, \geq$

Alternate H_a

$H_0 \mu = 25\%$

$H_a \mu \neq 25\% \quad \neq, >, <$

After claiming, we need to apply some statistical approach.

$\sigma \rightarrow z\text{-test}$

$\rightarrow t\text{-test}$

Then, we'll set some level of significance

$Z \leq p \quad (0\%) \quad \text{and} \quad Z > p$

19/8/26

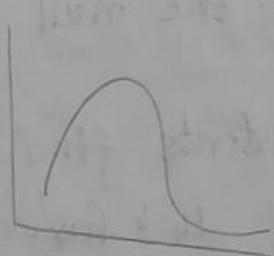
		Prediction	
		Yes	No
YES S A C T A N O	TP	FN	→ Type 2 error
	FP	TN	

Type 1 error

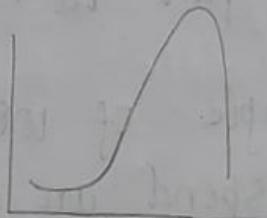
If my null hypothesis is true & I'm failing to reject that one - Type 1

bcoz actually it is true but we're failing.

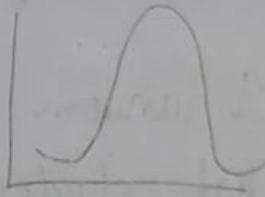
If my null hypothesis is false & we fail,
- Type 2



left-tailed



right-tailed



2-tailed



known

$$\text{mean } z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

proporⁿ In any sample, how many are there

$$\hat{P} = \frac{x}{n} \rightarrow np \text{ & } n(\hat{P}-1) \text{ at least 5.}$$

$$\text{not known } t = \frac{\bar{x} - \mu}{s / \sqrt{n}}$$

$$z = \frac{\hat{P} - p}{\sqrt{p(1-p)/n}} \quad p \rightarrow \text{claim proportion of hypothesis}$$

~~(algebra)~~

Q) An auto repair facility claims the avg time for oil-change is less than 20 mins.

atleast, atmost \rightarrow null

Null $\mu \geq 20$

less than, greater \rightarrow alterna

Alt. $\mu < 20$

Q) A medical researcher claims that the proportion of adults in the United States who r smokers is atmost 25%.

Null $\mu \leq 25\%$.

Alt $\mu > 25\%$.

* If \hat{P} value is less than^{or equal to} significance value, we reject

* If \hat{P} value is greater than significance value, we fail to reject the null

Q) A random sample of 1000 students find that, students spend an avg of 14.4 hrs/week on social media. Use this sample data to test the claim that Ctg student spend atleast 15 hrs/week on social media. Level of significance = 0.05, Popn of std devn = 3.25
soc):-

Null $\mu \geq 15$

$\mu < 15$

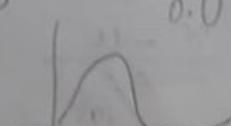
$$Z = \frac{14.4 - 15}{\sqrt{\frac{3.25^2}{1000}}}$$

- 0.05

- 1.83

0.0324 < 0.05

$\sqrt{1000}$



19/8/2023 We're going to reject the null hypothesis
→ The claim will be on sample, that is going to be reflected in populⁿ.

Based on t-test

Q) A smart phone manufacturer claims that the mean battery life of the phone is 25 hrs. Consumers & collect the sample of 50 SP. and determine the mean value 24.1 hr & std dev 4.1 hr. Use this sample data to test the claim made by smart phone manufacturer. Use the level of significance, 0.10.

Null $H_0 = 25$

Alt $H_1 \neq 25$

$$t = \frac{24.1 - 25}{\sqrt{\frac{4.1^2}{50}}} = \frac{-0.9}{0.127}$$

$$\frac{s}{\sqrt{n}} = \frac{4.1}{\sqrt{50}} = 0.562$$

$$\hookrightarrow 0.0635$$

2 tails \Rightarrow double value

$$0.127$$

→ we're failing

to reject the null hypothesis

$$0.127 > 0.1$$

i.e., we're accepting.

so, the manufacturer's claim is true

Based on Proporⁿ

Q) A Cdg professor claims that the proportion of students using the Chatgpt in their

19/8/29 assignment is less than 45%. Professor selects sample of 200 students & figure it out that 74 out of 200 use gpt. Use $\alpha = 0.05$

SOL:-

$$\text{Null } \mu \geq 0.45 \quad p = 0.45$$

$$\text{Alt } \mu < 0.45 \quad \hat{p} = 74$$

$$Z = \frac{74 - 0.45}{\sqrt{0.45(1-0.45)/200}} =$$

$$\text{first, } n\hat{p} = 200 \times 0.37 = 74 \geq 5$$

$$n(\hat{p}-1) = 200(0.37-1) = 5$$

$$Z = -2.274$$

$$Z \rightarrow p\text{-value} = 0.012 < 0.05$$

→ We will reject the null hypothesis, accept alternate hypothesis.
∴ The claim is true.

Hypothesis Testing for two samples

μ_1, μ_2 - mean of sample 1 & sample 2

$$\text{Null } \mu_1 = \mu_2 \quad \mu_1 \geq \mu_2 \quad \mu_1 \leq \mu_2$$

$$\text{Alt } \mu_1 \neq \mu_2 \quad \mu_1 < \mu_2 \quad \mu_1 > \mu_2$$

Samples can be dependent or independent
Independent

→ Random & Indep

~~20/8/25~~ → Sample should be of size 30 or follow normal distribution

→ sample variance or SD

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Q) Suppose a SE is claiming that the avg salary of women engineer is less than male. Two samples are taken, sample 1 = 45, sample 2 = 50) LOS = 0.05

Stm	WE	ME
Sample mean	72700	76900
SD	9800	10700
Size	45	50

WE \leq ME \Rightarrow Null. $\mu_1 \geq \mu_2$
 $\mu_1 \quad \mu_2$ At $\mu_1 < \mu_2$

If there is no difference, then there is no change in claim. so, there is always in null hypothesis.

$$\therefore \mu_1 = \mu_2 \Rightarrow (\mu_1 - \mu_2) = 0$$

$$t = \left(\dots \right) = -1.997 \xrightarrow{0.025 < 0.05}$$

2018/29

∴ The claim is true

Dependent samples:

difference (d) - before data - after data

$$\bar{d} = \frac{\sum d}{n}$$

$$t = \frac{\bar{d} - \mu_d}{s_d / \sqrt{n}}$$

→ When drug is taken, what is the bp &
when drug is not taken, what is bp?

μ_d , K constant

Null $\mu_d = 0$ or K , $\mu_d \geq K$

Alt $\mu_d \neq 0$ or K , $\mu_d < K$

$$\mu_d \leq K$$

$$\mu_d > K$$

(Q) Cholesterol lowering medicine.

Claim - Medic'm helps to lower cholesterol level., LOS = 0.05

Null $\mu_d \leq 0$

Alt $\mu_d > 0$

Claim - cholesterol is lower

∴ before - after > 0

But if claim is cholesterol is higher
then before - after < 0

20/8/05

Patient	BM	AM
1	218	210
2	232	241
3	269	223
4	265	244
5	248	227
6	298	223
7	263	252
8	281	276
9	290	281
10	271	259

$$\bar{d} = \frac{\sum d}{n} = 13.9 \rightarrow \text{On avg}$$

$$S_d = 12.414$$

$$t = \frac{13.9 - 0}{\frac{12.414}{\sqrt{10}}} = 3.541 \quad \hookrightarrow 0.003 \quad 0.003 < 0.05$$

Two proportion

$$\begin{array}{ll} x_1 & x_2 \\ n_1 & n_2 \end{array}$$

$$\hat{p}_1 = \frac{x_1}{n_1}, \quad \hat{p}_2 = \frac{x_2}{n_2}$$

$$\bar{P} = \frac{x_1 + x_2}{n_1 + n_2}$$

$n_1 \bar{P}, n_2 \bar{P}, n_1(1-\bar{P}), n_2(1-\bar{P})$ should be atleast 5

$$Z = \frac{(\hat{P}_1 - \hat{P}_2) - (P_1 - P_2)}{\sqrt{\hat{P}(1-\hat{P})} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

Based on our hypothesis

Q) Investigate to address lung cancer.
 People who took drug in 211 was cancer among 250, 172 among 230 free

$$LOS = 0.01$$

Claim: Proportion of subjects candidates $P_1 > P_2$ taking cancer-free drug are more cancerfree than the candidates who r not taking drug.

$$\text{Sol: } n_1 \bar{P} =$$

$$\hat{P}_1 = \frac{211}{250}, \hat{P}_2 = \frac{172}{230}$$

$$\bar{P} = \frac{211+172}{250+230} = 0.798$$

$$n_1 \bar{P} = 250 \times 0.798 \\ = 199.5$$

$$Z = 2.6$$

$$0.004 < 0.01$$

ANOVA test \leftarrow random independent

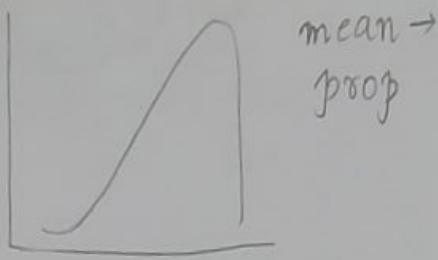
Analysis btw the variance

$\rightarrow Z, t \rightarrow$ one sample

$\rightarrow t \rightarrow$ two sample \leftarrow

22/8/25

F → distribution



$$\textcircled{1} \quad \bar{x}_1, \bar{x}_2, \bar{x}_3 \quad \& \quad s_1^2, s_2^2, s_3^2$$

$$\textcircled{2} \quad \text{Grand mean } \bar{\bar{x}}$$

$$\textcircled{3} \quad SSB = n_1 (\bar{x}_1 - \bar{\bar{x}})^2 + n_2 (\bar{x}_2 - \bar{\bar{x}})^2 + n_3 (\bar{x}_3 - \bar{\bar{x}})^2$$

$$\textcircled{4} \quad SSW = (n_1 - 1) s_1^2 + (n_2 - 1) s_2^2 + (n_3 - 1) s_3^2$$

$$\textcircled{5} \quad df_n = \text{no of groups} - 1$$

$$df_d = \text{total ss} - \text{no of group}$$

$$\textcircled{6} \quad MSB = \frac{SSB}{df_n} \quad \textcircled{7} \quad MSW = \frac{SSW}{df_d}$$

$$\textcircled{8} \quad F = \frac{MSB}{MSW}$$

two degree

→ no. of group

→ total no samples - no. of
g ~

A	B	C
s	s	s
\bar{x}_1	\bar{x}_2	\bar{x}_3

$$s_1^2 \quad s_2^2 \quad s_3^2$$

F < level of

P-value < level of sn

From Utube:

If $F > LOS$

Reject H_0

$F \leq LOS$

Fail to Reject
 H_0

S.No	A	B	C	D
1	8	12	18	13
2	10	11	12	9
3	12	9	16	12
4	8	14	6	16
5	7	4	8	15

Claim: Height of popⁿ > Mean value of popⁿ

$$\text{Mean}(A) = \frac{45}{5} = 9 \quad \text{Var}(A) = \frac{16}{4} = 4 \quad \text{Std}(A) = 2$$

$$\text{Mean}(B) = \frac{50}{5} = 10 \quad \text{Var}(B) = \frac{58}{4} = 14.5 \quad \text{Std}(B) = 3.8$$

$$\text{Mean}(C) = \frac{60}{5} = 12 \quad \text{Var}(C) = \frac{104}{4} = 26 \quad \text{Std}(C) = 5.09$$

$$\text{Mean}(D) = \frac{65}{5} = 13 \quad \text{Var}(D) = \frac{30}{4} = 7.5 \quad \text{Std}(D) = 2.73$$

② Grand mean $\bar{x} = \frac{9+10+12+13}{4} = 11$

③ SSB = $5(9-11)^2 + 5(10-11)^2 + 5(12-11)^2 + 5(13-11)^2$
 $= 20 + 5 + 5 + 20$
 $= 50$

④ SSW = $(5-1)4 + (5-1)14.5 + (5-1)26 + (5-1)7.56$
 $= 208$

⑤ df_m = 3
df_d = 20 - 3 = 17

22/8/22

$$\textcircled{6} \quad \text{MSB} = \frac{50}{3} \\ = 16.66$$

$$\textcircled{7} \quad \text{MSW} = \frac{208}{17} \\ = 12.23$$

Correlation Analysis \rightarrow Regression

Representation

Whether the features are going to have impact on output or not.

+ve \rightarrow \uparrow indep \uparrow dep

0 \rightarrow x

-ve \rightarrow \uparrow indep \downarrow dep

Pearson correlation = $\frac{n \sum xy - (\sum x)(\sum y)}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{\sum y^2 - (\sum y)^2}}$

↑
value \rightarrow strength
sign \rightarrow direction

$r \rightarrow$ Populn P.C Coefficient

Linear relationship $\Rightarrow |r| > \frac{2}{\sqrt{n}}$

Q) Suppose the chief financial officer is investigating the correln between rest of price

In a sample data size, 10 points are taken.
Now, based on this analysis can chief financial

~~22/8/2023~~ officer conclude that the correlation is significant b/w the stock price & unemployment rate with the LOS = 0.06.

Soln) $r = -0.68$

$$n = 10$$

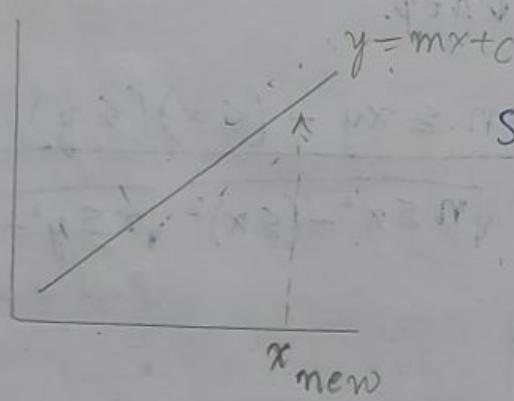
$$|-0.68| > \frac{2}{\sqrt{10}}$$

$$0.68 > 0.63$$

Yes, the relationship is there.

They're going to be very correlated

Linear Regression



Classification
Supervised ML → Regression
(Discrete)
(Continuous)
Loss / cost

$$\hat{y} = w_0 f_0 + w_1 f_1 + w_2 f_2 + w_3 f_3 + w_4 f_4 + c$$

$$c = w_0 f_0$$

$$f_0 = 1$$

$$\hat{y} = w_0 f_0 + w_1 f_1 + \dots + w_n f_n$$

w_0	f_0	w_1	f_1	w_2	f_2	w_3	f_3	w_4	f_4	y

$$\hat{y} = WF - \hat{y} = WX$$

$$\rightarrow y = ax + b$$

→ Closed Normal form $\frac{\partial \text{Loss}}{\partial w}$

→ Gradient descent $\rightarrow n$

26/8/23 Linear Regression

$y = mx + c$	$w_1 w_2 w_3 w_4$	$x_1 x_2 x_3 x_4$	y



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

Loss / Cost / Error / Bias

We will minimize

$$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + c$$

$$\text{where } x_0 = 1 \quad \text{MAE} = \frac{1}{n} |y - \hat{y}|$$

$$\hat{y} = wX \quad \text{MSE} = \frac{1}{n} (y - \hat{y})^2$$

$$w = [w_0 \dots w_n]$$

$$\text{RMSE} = \sqrt{\frac{1}{n} (y - \hat{y})^2}$$

$$x = [x_0 \dots x_n]$$

Cross entropy

$$\text{where } x_0 = 1$$

$$\hat{y} = y \cdot \ln(\hat{y}) + (1 - \hat{y}) \cdot \ln(1 - \hat{y})$$

LR → Closed normal form

→ Gradient descent

$$\hat{y} = wX$$

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

$$\frac{\partial \text{Loss}}{\partial w} = \frac{1}{n} [2(y - \hat{y})] \times \frac{\partial}{\partial w} (y - \hat{y})$$

$$\therefore 2X(y - \hat{y}) * \frac{\partial}{\partial w} (y - Xw)$$

→ Equate it to 0.

$$\cancel{26/8/2023} \quad 2(y - \hat{y}) * \frac{\partial}{\partial w}(y - x^T w) = 0$$

$$2(y - x^T w)^* (-x) = 0$$

$$2 \cdot x^T(y - x^T w) = 0$$

$$x^T(y - x^T w) = 0$$

$$x^T y - x^T x^T w = 0$$

$$x^T x^T w = x^T y$$

$$w = (x^T x)^{-1} x^T y$$

Code

$$X = 2^* np.random.rand(m, 1)$$

$$y = 4 + 3^* X + np.random.randn(m, 1)$$

from sklearn.preprocessing

normal distributed values

import add_dummy_feature

$$X_b =$$

$$a.T \rightarrow a^T$$

$$\theta_{\text{best}} = np.linalg.inv(X_b.T \cdot \dot{dot}(X_b))$$

$$\theta_{\text{best}}$$

$$dot(X_b.T).dot$$

$$X_{\text{new}} = np.array([0, 1])$$

$$\text{print}(X_{\text{new}})$$

$$X_{\text{new}}_b = \text{add_dummy_feature}(X_{\text{new}})$$

$$y_{\text{predict}} = X_{\text{new}}_b \cdot \dot{dot}(\theta_{\text{best}})$$

$$y_{\text{predict}}$$

26/8/18 Using sklearn

```
from linear_model import Lasso
```

Lasso-reg = Lasso()

Lasso-reg.fit(X, y)

Lasso-reg.intercept_, Lasso-reg.coef_

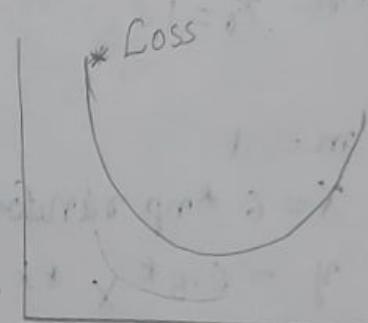
$$\downarrow \omega_0$$

$$\downarrow \omega_1$$

Gradient descent

$$\omega_{i+1} = \omega_i - \eta \frac{\partial}{\partial \omega} \text{Loss}$$

Learning rate



In the starting, take some curve should be random values for ω . convex curve

Batch G-D - Whole dataset is given in Ig

G-D Stochastic G-D - Any 1 random instance

Mini-batch

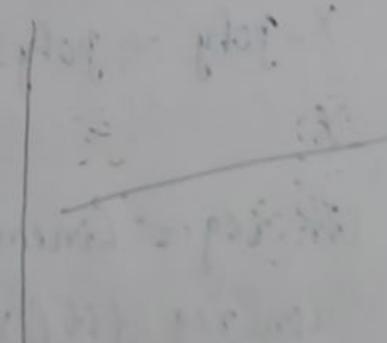
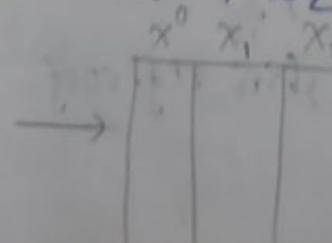
Underfitting ←

Polynomial Regression

$x, d=2$

x^0, x^1, x^2

w_0, w_1, w_2



$(x, y), \text{degree}=2$

x, y, xy, x^2, y^2
 w_0, w_1, w_2, w_3

~~26/8/23~~ degree - d , num of features - n
→ How many new features will come?

$${}_{n+d}^{n+d} C_d = \frac{(n+d)!}{n! \times d!}$$

$$y = wX$$

$$= w_0 + w_1 x$$

$$w_0 x_0 + w_1 x$$

$$\text{where } x_0 = 1$$



Code

$$m = 1$$

$$X = 6 * np.random.rand(m, 1) - 3$$

$$y = 0.5 * X ** 2 + X + 2 + np.random.randn(m)$$

for finding polynomial features:

from

import PolynomialFeatures

poly_features = PolynomialFeatures

(degree=2, include_bias=False)

X_poly = poly_features.fit_transform(X)

X

lin_reg = LinearRegression()

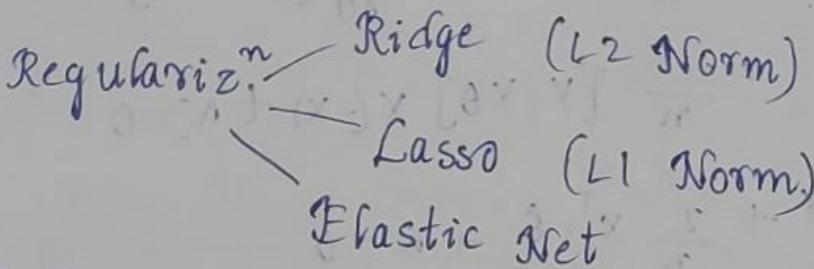
lin_reg.fit(X_poly, y)

lin_reg.intercept_, lin_reg.coef_

Regularization

$$\text{Cost} = \frac{1}{n} (y - \hat{y})^2 + \underline{\alpha \sum_{i=1}^m \theta_i^2}$$

Minimize w.r.t. θ



1) Ridge $J(\theta) = \text{MSE} + \alpha \sum_{i=1}^m \theta_i^2$

$$\text{Cost} = \frac{1}{n} (y - \hat{y})^2 + \alpha \sum_{i=1}^m \theta_i^2$$

$$= \frac{1}{n} (y - \hat{y})^2 + \alpha \cdot \theta^2$$

$$\frac{\partial \text{Cost}}{\partial \theta} = \frac{1}{n} \times 2 \times (y - \hat{y}) \times \frac{\partial}{\partial \theta} (y - \hat{y}) + \alpha \times 2\theta$$

$$= \frac{2}{n} (y - X\theta) \times \frac{\partial}{\partial \theta} (y - X\theta) + 2\alpha\theta$$

$$= \frac{2}{n} (y - X\theta) \times (-X) + 2\alpha\theta = 0$$

$$= -XY + X^T X \theta + 2\alpha\theta = 0$$

$$\therefore \theta = \frac{XY}{X^T X + 2\alpha}$$

$\alpha \rightarrow \text{large} \Rightarrow \theta \rightarrow \text{very small} \Rightarrow \text{Underfitting}$

2) Lasso Regulariz.

$$J(\theta) = \text{MSE} + \alpha \sum_{i=1}^n |\theta_i|$$

$$\text{cost} = \frac{1}{n} [y - \hat{y}]^2 + \alpha \theta$$

$$\begin{aligned}\frac{\partial}{\partial \theta} \text{cost} &= \frac{1}{n} X_2 [y - \hat{y}] \times \frac{\partial}{\partial \theta} [y - \hat{y}] + \alpha \\ &= \frac{2}{n} [y - \theta x] \times \frac{\partial}{\partial \theta} [y - \theta x] + \alpha \\ &= \frac{2}{n} [y - \theta x] \times (-x) + \alpha = 0 \\ -\frac{2}{n} x^T y + \frac{2}{n} x^T x \theta + \alpha &= 0\end{aligned}$$

$$\theta = \frac{\frac{2}{n} x^T y - \alpha}{\frac{2}{n} x^T x}$$

3) Elastic Net

$$J(\theta) = \text{MSE} + \gamma \alpha \sum_{i=1}^m \theta_i + \frac{1-\gamma}{2} \alpha \sum_{i=1}^n \theta_i^2$$

L1
L2

We're adding something to MSE, so that our model will work in real-time scenario

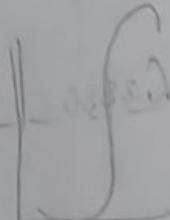
Logistic Regression

→ Classification

Apply sigmoid fn 0-1 $\Rightarrow \frac{1}{1+e^{-t}}$

$0.5 \leq N$ Negative Class

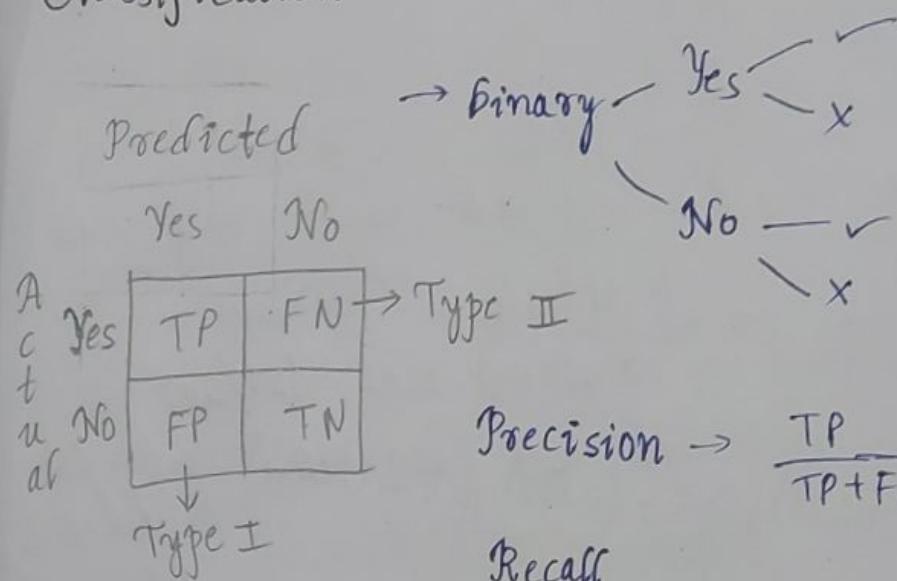
$> P$ Positive Class



Softmax - multi-class classification

$$\frac{e^{z_i x}}{\sum_{j=1}^n e^{z_j x}}$$

Classification



Recall

(Completeness) → $\frac{TP}{TP+FN}$

We need a
mid-way of

Precision & Recall → F1-score →

Accuracy → $\frac{TP+TN}{TP+TN+FP+FN}$

$$\frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

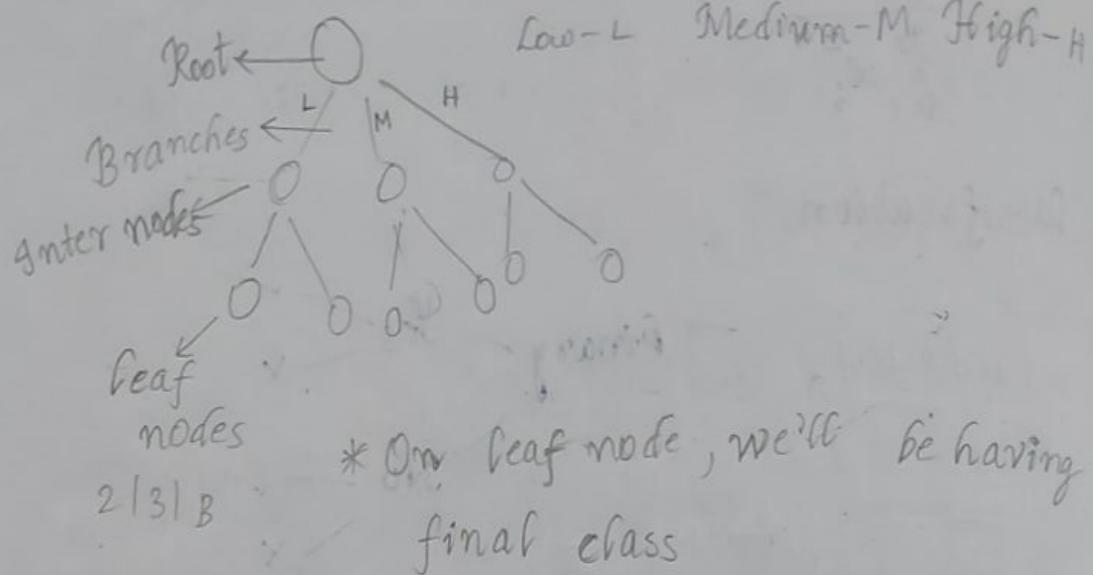
Underfitting → Training \times Testing \times

Overfitting → Training \checkmark Testing \times

Loss_{error} = Bias² + Variance + irreducible Noise

2/9/23

Decision tree



2/31/B

- It can work on n-dimensional dataset
- Prior info is not required
- Domain-expert
- Overfitting - Because it's a rule-based algorithm

Advantage

No Learning
Disadvantage

How are we going to split & Attribute selection?

Discrete \Rightarrow L/M/H or Yes/No

Continuous \Rightarrow Mean
 \leq / \geq [value \leq Mean]

Attribute Selection \rightarrow Entropy or Infoⁿ. Gain

Gain Ratio \rightarrow C4.5
ID3

Index Root \rightarrow CART

Entropy (or) IG → Measure that quantifies the avg. amount of inform" or uncertainty.

→ Randomness

$$IG_E(\text{Class}) = -\frac{1}{N} \sum_i P_i$$

$$= - \sum_i P_i \log_2 P_i$$

→ how the class is going to be distributed

$$\begin{matrix} 0 & 0 & 1 & 1 & 2 & 0 & 3 & 4 \\ \frac{3}{8} & \frac{2}{8} & \frac{1}{8} \end{matrix}$$

$$IG_A(\text{Class}) = \sum_{i=1}^N \frac{|D_i|}{|D|} \times IG(D_i)$$

$$\Delta \text{Gain}_A(\text{Class}) = IG(\text{Class}) - IG_A(\text{Class})$$

Root

D₁ D₂

→ When all instances are saying, this is the class

→ If all instances & rows are fully utilized

Basic criteria for stopping DT:

- ① All instances are agreeing for 1 class
- ② All attributes are utilized
- ③ Subset that we are generating does not have any instance.

8/9/25

S.N	Age	Income	Student	Credit Rating	Buy - camp
1.	Y	H	N	F	N
2.	Y	H	N	E	N
3.	M	H	N	F	Y
4.	S	M	N	F	Y
5.	S	L	Y	F	Y
6.	S	L	Y	E	N
7.	M	L	Y	E	Y
8.	Y	M	N	F	N
9.	Y	L	Y	F	Y
10.	S	M	Y	F	Y
11.	Y	M	Y	E	Y
12.	M	M	N	E	Y
13.	M	H	Y	F	Y
14.	S	M	N	E	N

Y- Youth, M- Middle Age, S- Senior Citizen

$$IG(\text{Class}) = - \sum_{i=1}^n p_i \log_2 p_i$$

$$IG_{\text{Att}}(\text{Class}) = \sum_{i=1}^N \frac{|D_i|}{|D|} \cdot IG(D_i)$$

$$\Delta G_{(A)} = IG_{\text{class}} - IG_{\text{Att(class)}}$$

$$Y = \frac{9}{14} \quad IG(\text{Buy-Comp}) = -\frac{9}{14} \cdot \log_2 \frac{9}{14}$$

$$N = \frac{5}{14} \quad -\frac{5}{14} \log_2 \frac{5}{14}$$

$$= -\frac{9}{14}(-0.637) - \frac{5}{14}(-1.486)$$

$$= 0.4095 + 1.0398 + 0.5303$$

$$\neq 1.449$$

$$= 0.939$$

$$\approx 0.94$$

$$\log_a b = \frac{\log b}{\log a}$$

$$\log_2 \frac{9}{14} = \frac{\log \frac{9}{14}}{\log 2}$$

$$\begin{matrix} IG_{\text{Age}} & = \frac{5}{14} IG(D_1) + \frac{4}{14} IG(D_2) + \frac{5}{14} IG(D_3) \end{matrix}$$

$$Y_S = \frac{2}{5} \quad M_4 = \frac{4}{5} \quad S_5 = \frac{3}{5}$$

$$= \frac{5}{14} \left[\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right]$$

$$M_4 = \frac{4}{5} \quad S_5 = \frac{3}{5}$$

$$+ \frac{4}{14} \left[-\frac{4}{4} \log_2 \frac{4}{4} - 0 \right]$$

$$+ \frac{5}{14} \left[-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right]$$

$$= \frac{5}{14} \left[-\frac{2}{5} (-1.321) - \frac{3}{5} (-0.736) \right]$$

$$+ \frac{4}{14} \left[-1(0) - 0 \right] + \frac{5}{14} \left[-\frac{3}{5} (-0.736) - \frac{2}{5} (-1.321) \right]$$

$$= \frac{5}{14} [0.5284 + 0.441] + \frac{5}{14} [0.4416 + 0.5284]$$

$$= \frac{5}{14} \times 0.9694 + \frac{5}{14} \times 0.9694 = 0.692$$

$$\approx 0.694$$

$$\Delta G = 0.94 - 0.694$$

$$= 0.246$$

$$IG(\text{Student}) = IG_{\text{student}} (\text{Class}) = \frac{7}{14} IG(D_1)$$

~~N - 17~~
~~Y - 17~~

~~Y - 1~~
~~M - 3~~

~~B - 4~~

$$\begin{aligned} &= \frac{7}{14} \left[-\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \right] + \frac{7}{14} \left[-\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \right] \\ &= \frac{7}{14} [0.191 + 0.401] + \frac{7}{14} [0.524 + 0.461] \\ &= \frac{7}{14} \times 0.592 + \frac{7}{14} \times 0.985 \\ &= 0.296 + 0.493 = 0.789 \end{aligned}$$

$$\Delta G_{(\text{student})} = 0.94 - 0.789 = 0.151$$

$$IG_{\text{Income}} (\text{Class}) = \frac{4}{14} IG(D_1) + \frac{6}{14} IG(D_2) + \frac{4}{14} IG(D_3)$$

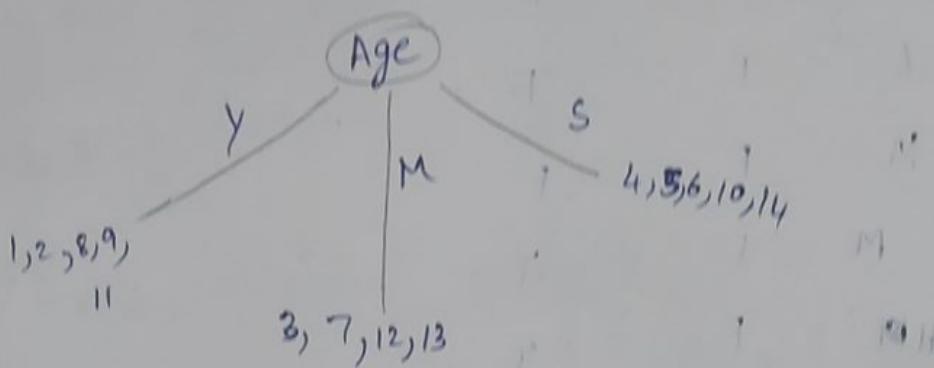
~~L - 1~~
~~M - 4~~
~~H - 2~~
~~V - 2~~

$$\begin{aligned} &= \frac{4}{14} \left[-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right] + \frac{6}{14} \left[-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right] \\ &= \frac{4}{14} [0.311 + 0.5] + \frac{6}{14} [0.390 + 0.528] + \frac{4}{14} [1] \\ &= \frac{4}{14} (0.811) + \frac{6}{14} (0.918) + 0.286 \\ &= 0.232 + 0.393 + 0.286 = 0.911 \end{aligned}$$

$$\Delta G_{\text{Income}} = 0.94 - 0.911 = 0.029$$

Similarly, $\Delta G_{\text{Credit Rating}} = 0.048$

3/9/20 * Age has highest - 0.246
So, it will be root node



$$IG(D_1) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$IG_{\text{income}}(D_1) = \frac{2}{5} IG(D_1) + \frac{2}{5} IG(D_1) + \frac{1}{5} IG(D_1)$$

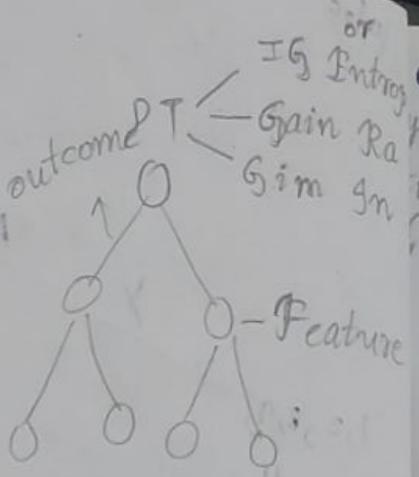
$$IG_{\text{student}}(D_1) = \frac{2}{5} \left[-\frac{2}{5} \log_2 \frac{2}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right]$$

$$IG_{\text{CR}}(D_1)$$

$H \leq 2$
 $M \leq 1$
 $L \leq 1$

9/9/28

	GPA	Studied	Passed
1	L	F	F
2	L	T	T
3	M	F	F
4	M	T	T
5	H	F	T
6	H	T	T



Added by me, not feature

$$IG(\text{Class}) = - \sum_{i=1}^n P_i \log_2 P_i$$

$$IG_A(\text{Class}) = \sum_{i=1}^n \frac{|D_i|}{|D|} IG\left(\frac{D_i}{|D|}\right)$$

$$\Delta G = IG(\text{Class}) - IG_A(\text{Class})$$

Stopping Criteria for DT: \rightarrow All instances belong to
the same class

- \rightarrow No features are left
- \rightarrow No instances are left

$$\begin{aligned}
 IG(\text{Class}) &= -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \\
 &= 0.390 + 0.528 \\
 &= 0.918 \approx 0.92
 \end{aligned}$$

$$\begin{aligned}
 IG_{GPA}(\text{Class}) &= \frac{2}{6} [IG(D_1)] + \frac{2}{6} [IG(D_2)] + \frac{2}{6} [IG(D_3)]
 \end{aligned}$$

$D_1 = \{L, M\}$
 $D_2 = \{L, H\}$
 $D_3 = \{M, H\}$

$$\text{alghs} = \frac{2}{6} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] + \frac{2}{6} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right]$$

$$+ \frac{2}{6} \left[-\frac{2}{2} \log_2 \frac{2}{2} - 0 \right]$$

$$= \frac{2}{6} [1] + \frac{2}{6} [1] + \frac{2}{6} [0] = \frac{4}{6} = 0.667$$

$$IG_{\text{studied}} (\text{class}) = \frac{3}{6} IG(D_1) + \frac{3}{6} IG(D_2)$$

T = 0
F = 1
- 2

$$= \frac{3}{6} \left[-\frac{3}{3} \log_2 \frac{3}{3} - 0 \right] + \frac{3}{6} \left[-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right]$$

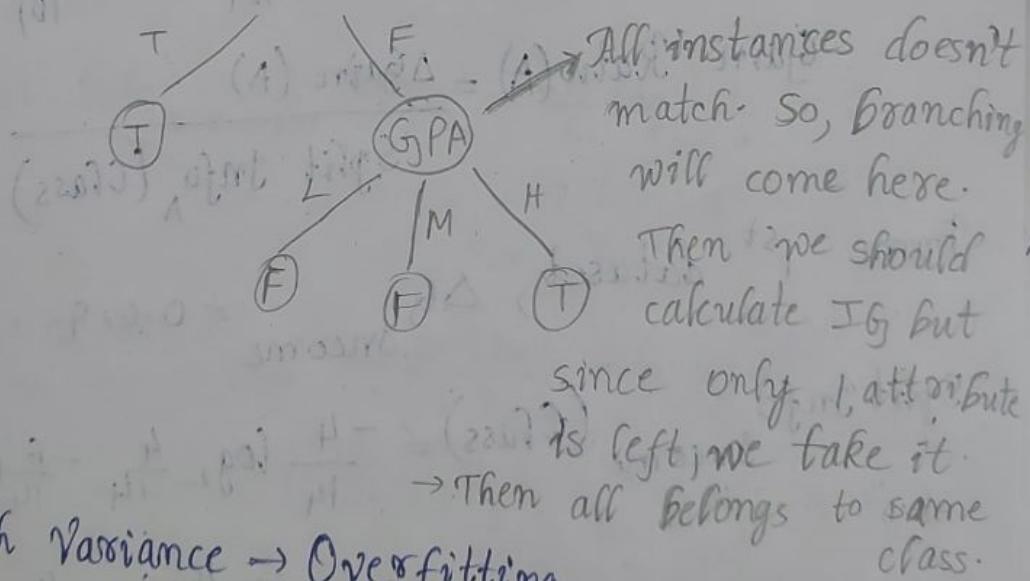
$$= \frac{3}{6} [0] + \frac{3}{6} [0.528 + 0.390]$$

$$= \frac{3}{6} [0.918] = 0.459$$

$$\Delta G_{GPA} = 0.92 - 0.67 = 0.25$$

$$\Delta G_{\text{studied}} = 0.92 - 0.46 = 0.46$$

Studied



High Variance \rightarrow Overfitting

(Training accuracy is high & testing accuracy is less)

Why?

~~9/9/25~~ Problems of Information Gain or Entropy method

$$IG_{S.N}(\text{Class}) = \frac{1}{6} \left[\frac{1}{2} \log_2 \frac{1}{2} \right] + \frac{1}{6} \left[\frac{1}{2} \log_2 \frac{1}{2} \right] \\ = 0$$

$$\Delta G_{S.N} = 0.92 - 0 \\ = 0.92$$

→ Uncertainty very low for these columns then ΔG high.

Gain Ratio

$$\text{Split Info}_A(\text{Class}) = - \sum_{i=1}^{|D_i|} \frac{|D_i|}{|D|} \log_2 \frac{|D_i|}{|D|}$$

$$\text{Gain Ratio}(A) = \frac{\Delta \text{Gain}(A)}{\text{Split Info}_A(\text{Class})}$$

From 1st dataset, $\Delta G_{\text{Income}} = 0.029$

$$\begin{aligned} \text{Split Info}_{\text{Income}}(\text{Class}) &= - \frac{4}{14} \log_2 \frac{4}{14} - \frac{6}{14} \log_2 \frac{6}{14} \\ &\quad - \frac{4}{14} \log_2 \frac{4}{14} \\ &= 0.516 + 0.524 + 0.516 \\ &= 1.556 \end{aligned}$$

$$\text{Gain Ratio}_{(\text{Income})} = \frac{0.029}{1.556} = 0.019$$

also
2nd dataset:

$$\Delta G_{GPA} = 0.26$$

$$\text{Split Info}_{GPA}(\text{Class}) = -\frac{2}{6} \log_2 \frac{2}{6} - \frac{2}{6} \log_2 \frac{2}{6}$$
$$= -\frac{2}{6} \log_2 \frac{2}{6}$$
$$= 0.528 \times 3 = 1.584$$

$$\text{Gain ratio}(GPA) = \frac{0.26}{1.584} = 0.164$$

$$\Delta G_{\text{studied}} = 0.46$$

$$\text{Split Info}_{\text{studied}}(\text{Class}) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$
$$= 0.5 + 0.5 = 1$$

$$\text{Gain ratio}_{(\text{studied})} = \frac{0.46}{1} = 0.46$$

$$\Delta G_{S\text{-No}} = 0.92$$

$$\text{Split Info}_{\cancel{S\text{-No}}}(\text{Class}) = -\frac{1}{6} \log_2 \frac{1}{6} \times 6$$
$$= 0.431 \times 6$$
$$= 2.586$$

$$\text{Gain ratio}_{S\text{-No}} = \frac{0.92}{2.586} = 0.358$$

Name of DT based on Gain Ratio = 0.358

Name of DT based on Information Gain = ID3

Data mining & Techniques - Chapter 8

Decision Tree

Gini Index \rightarrow CART

→ This will be used for classification & regression.

$$\text{Gini (class)} = 1 - \sum_{i=1}^n p_i^2$$

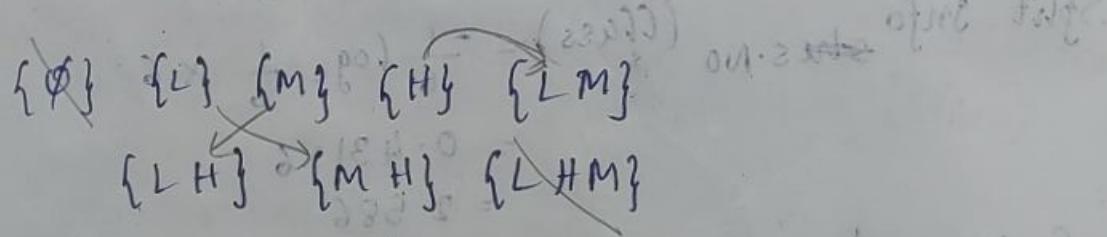
$$\text{Gini}_A(\text{class}) = \frac{|D_1|}{|D|} \text{Gini}(D_1) + \frac{|D_2|}{|D|} \text{Gini}(D_2)$$

$$\Delta \text{Gini}_A(\text{class}) = \text{Gini}(\text{class}) - \text{Gini}_A(\text{class})$$

→ Gini index is working on the binary split

→ We're going to form the sets.

3 attributes:



→ Whichever split has higher Δ value, that split will be considered.

* In the question, if it is mentioned 'multiway' split, then no need to do that process.

S.No	Body Temp	Birth	4-legged	Hibernates	Class (Mammal)
1	Warm	Y	Y	Y	Y
2	Warm	Y	Y	N	Y
3	Warm	Y	N	Y	N
4	Warm	Y	N	N	N
5	Cold	N	Y	Y	N
6	Cold	N	Y	N	N
7	Cold	N	N	Y	N
8	Cold	N	N	N	N
9	Warm	N	N	N	N
10	Cold	Y	N	N	N

$$\begin{aligned}
 \text{Gini Index (Class)} &= 1 - \sum_{i=1}^n p_i^2 \\
 &= 1 - \left(\frac{2}{10}\right)^2 - \left(\frac{8}{10}\right)^2 \\
 &= 0.32
 \end{aligned}$$

$$\begin{aligned}
 \text{Gini Body Temp (Class)} &= \frac{5}{10} \text{Gini (Warm)} + \frac{5}{10} \text{Gini (Cold)} \\
 &= \frac{5}{10} \left[1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 \right] + \frac{5}{10} \left[1 - 0 - \left(\frac{5}{5}\right)^2 \right] \\
 &= \frac{5}{10} [0.48] + \frac{5}{10} (0) = 0.24
 \end{aligned}$$

$$\text{Gini Body temp (Class)} = 0.32 - 0.24 = 0.08$$

$$\begin{aligned}
 \text{Gini Birth (Class)} &= \frac{5}{10} \text{Gini (Y)} + \frac{5}{10} \text{Gini (N)} \\
 &= 0.24 \\
 \therefore \Delta \text{Gini Birth} &= 0.08
 \end{aligned}$$

10/9/25

$$Gini_{4\text{-legged}} = \frac{4}{10} \left[1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \right] + \frac{6}{10} \left[1 - \left(\frac{5}{6}\right)^2 \right]$$

$$Y \leq 2 \quad = \frac{4}{10} \times 0.5 + 0$$

$$N \leq 6 \quad = 0.2$$

$$\Delta Gini_{4\text{-legged}} = 0.12$$

$$Gini_{Hibernate} = \frac{4}{10} \left[1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 \right] + \frac{6}{10} \left[1 - \left(\frac{1}{6}\right)^2 - \left(\frac{5}{6}\right)^2 \right]$$

$$Y \leq 3 \quad = \frac{4}{10} \times 0.375 + \frac{6}{10} \left[0.278 \right]$$

$$N \leq 5 \quad = 0.15 + 0.467$$

$$= 0.317$$

$$\Delta G_{Hibernate} = 0.32 - 0.317 = 0.003$$

```

graph TD
    Root["4-legged"] -- Y --> Leaf1["1, 2, 5, 6"]
    Root -- N --> Leaf2["N"]
  
```

$$Gini_{Class D_1} = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 0.5$$

$$Gini_{BT(D_1)} = \frac{2}{4} \left[1 - \left(\frac{2}{2}\right)^2 - 0 \right] + \frac{2}{4} \left[1 - \left(\frac{2}{2}\right)^2 - 0 \right]$$

$$= 0$$

$$\Delta Gini_{BT} = 0.5 - 0 = 0.5$$

10/9/25

$$\text{Gini}_{\text{Birth}} (D_1) = \frac{2}{4} \left[1 - \left(\frac{1}{2}\right)^2 - 0 \right] + \frac{2}{4} \left[1 - \left(\frac{1}{2}\right)^2 - 0 \right]$$

$$= 0$$

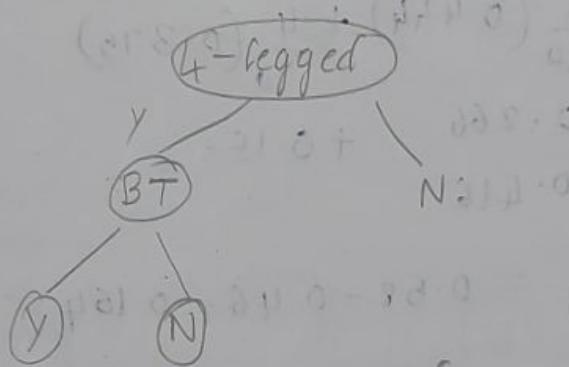
$$\Delta \text{Gini}_{\text{Birth}} = 0.5 - 0 = 0.5$$

$$\text{Gini}_{\text{Hibernate}} (D_1) = \frac{2}{4} \left[1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \right] + \frac{2}{4} \left[1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \right]$$

$$= \frac{2}{4} (0.5) + \frac{2}{4} (0.5)$$

$$= 0.5$$

$$\Delta \text{Gini}_{\text{Hibernate}} = 0.5 - 0.5 = 0$$



SNO	Weather	Parents	Money	Multiway Split	
				Rich	Poor
1	S	Y	R	C	T
2	S	N	R	T	C
3	W	Y	R	C	T
4	R	Y	P	C	S
5	R	N	R	S	T
6	R	Y	P	C	C
7	W	N	P	C	S
8	W	N	R	S	H
9	W	Y	R	C	C
10	S	N	R	T	T

Rich Cinema
Poor Tennis
Stay in Shopping

$$\text{Gini (Class)} = 1 - \left[\left(\frac{6}{10}\right)^2 + \left(\frac{2}{10}\right)^2 + \left(\frac{1}{10}\right)^2 + \left(\frac{1}{10}\right)^2 \right] \\ = 0.58$$

$$\begin{aligned} \text{Gini Weather (Class)} &= \frac{3}{10} \left[1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 \right] \\ &\quad + \frac{4}{10} \left[1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \right] \\ &\quad + \frac{3}{10} \left[1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \right] \\ &= \frac{6}{10} (0.444) + \frac{4}{10} (0.375) \\ &= 0.266 + 0.15 \\ &= 0.416 \end{aligned}$$

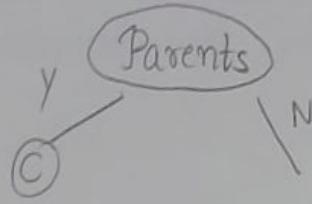
$$\Delta \text{Gini}_{\text{Weather}} = 0.58 - 0.416 = 0.164$$

$$\begin{aligned} \text{Gini Parents (Class)} &= \frac{5}{10} \left[1 - \left(\frac{5}{5}\right)^2 \right] + \frac{5}{10} \left[1 - \left(\frac{1}{5}\right)^2 - \left(\frac{2}{5}\right)^2 \right. \\ &\quad \left. - \left(\frac{1}{5}\right)^2 - \left(\frac{1}{5}\right)^2 \right] \\ &= 0 + \frac{5}{10} [0.72] = 0.36 \\ \Delta \text{Gini}_{\text{Parents}} &= 0.58 - 0.36 = 0.22 \end{aligned}$$

$$\begin{aligned} \text{Gini Money (Class)} &= \frac{7}{10} \left[1 - \left(\frac{3}{7}\right)^2 - \left(\frac{2}{7}\right)^2 - \left(\frac{1}{7}\right)^2 - \left(\frac{1}{7}\right)^2 \right] \\ &\quad + \frac{3}{10} \left[1 - \left(\frac{3}{3}\right)^2 \right] \\ &= \frac{7}{10} [0.694] + 0 = 0.486 \end{aligned}$$

$$\Delta \text{Gini}_{\text{Money}} = 0.58 - 0.486 = 0.094$$

10/10/25



$$\text{Gini}(\text{Class}_P) = 1 - \left[\left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{1}{5}\right)^2 \right]$$

$$= 1 - 0.28 = 0.72$$

$$\text{Gini}_{\text{Weather}}(\text{Class}_P) = \frac{2}{5} \left[1 - \left(\frac{2}{2}\right)^2 \right] + \frac{2}{5} \left[1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \right]$$

$$+ \frac{1}{5} \left[1 - \left(\frac{1}{1}\right)^2 \right]$$

$$= \frac{2}{5} (0.5) = 0.2$$

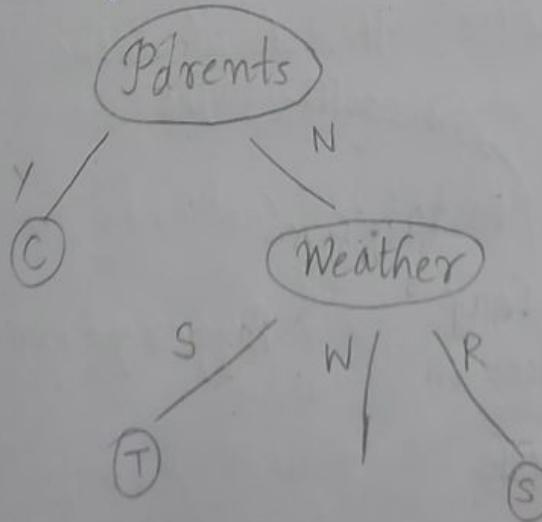
$$\Delta \text{Gini}_{\text{Weather}} = 0.72 - 0.2 = 0.52$$

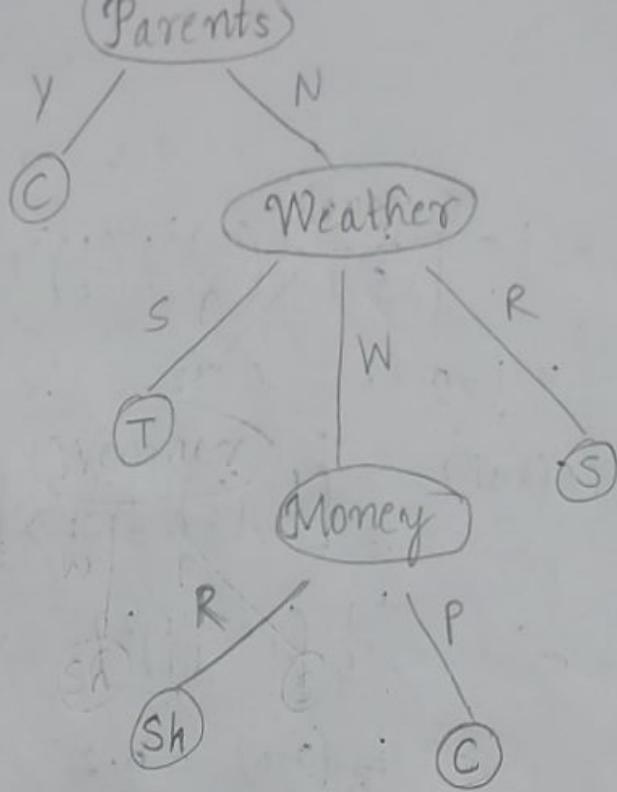
$$\text{Gini}_{\text{Money}}(\text{Class}_P) = \frac{4}{5} \left[1 - \left(\frac{2}{4}\right)^2 - \left(\frac{1}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \right]$$

$$+ \frac{1}{5} \left[1 - \left(\frac{1}{1}\right)^2 \right]$$

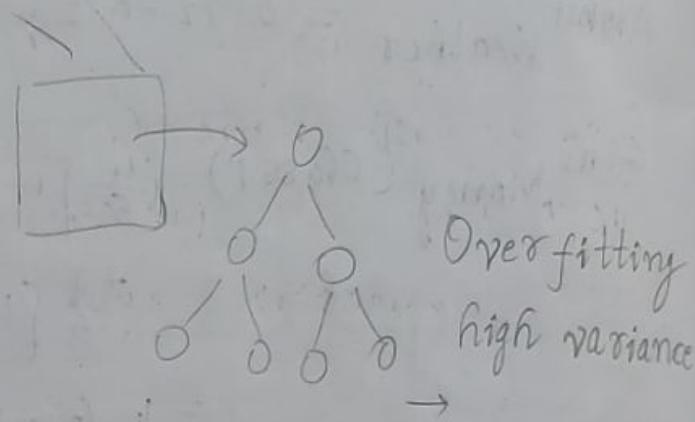
$$= \frac{4}{5} (0.625) = 0.5$$

$$\Delta \text{Gini}_{\text{Money}} = 0.72 - 0.5 = 0.22$$





Tree Pruning



Tree Pruning

- ← Pre-Pruning
- ← Post-Pruning

→ We will set some threshold value

→ Max-depth

→ Min-split-sample

→ Min-sample

instance

(50 or 60)

→ Min no. of leaf instances in

hyperparameters

→ Max_leaf node

→ Max_features

12/10/26 Assignment: Take any dataset, set hyperparameters, if value low-impact? } based on this
" " high-impact? } what's the analysis?

Take the features 1 by 1.

→ We need to check if x_1 & x_2 are required for y_2 or not. Test is Chi-square test

Pre-Pruning → Thresholds
→ Dependency
→ Improvement (difference in delta values)
for each split

Post → Whole decision tree should be constructed

→ Removal of a branch impacts the DT accuracy or not. For this, we can apply cross validation approach.

KNN

Predicⁿ based on Majority Voting Class

Brightness	Saturation	Class	
40	20	Red	
60	50	Blue	$k=5$
60	90	Blue	$\text{Brightness} = 20$
10	25	Red	$\text{Satur}^n = 35$
70	70	Blue	
60	10	Red..	$\text{Class} = ?$
25	80	Blue	

<u>Dist(new, X_i)</u>	
② 25	$\sqrt{(20-40)^2 + (35-20)^2}$
③ 33.54	$\sqrt{(20-50)^2 + (35-50)^2}$
68.01	$\sqrt{(20-60)^2 + (35-90)^2}$
① 14.1	$\sqrt{(20-10)^2 + (35-25)^2}$
61.03	$\sqrt{(20-70)^2 + (35-70)^2}$
⑤ 47.17	$\sqrt{(20-60)^2 + (35-10)^2}$
④ 45.27	$\sqrt{(20-25)^2 + (35-80)^2}$

KNN - Instance based

Limitation of KNN:

- ① Applicable only for small dataset.
 RAM Capacity cannot be that much
 to calculate many distances.
- ~~(1)~~ $x_1 \rightarrow 0.1, x_2 \rightarrow 100000 \rightarrow$ Varying ranges b/w
 features

Apply normalization first

- ② It is applicable only for numeric data
 Issue of missing values, even if we
 fill it, it will effect the performance

KNN Limitations

- Small
- Numeric

Normalization is required

Training time is 0, testing time is high.

$$P(C/x) = \frac{P(x/c) \cdot P(c)}{P(x)} \xrightarrow{\text{Prior}} P(c) \approx P(x/c) \quad P(\frac{A}{B}) = \frac{P(B/A) \cdot P(A)}{P(B)}$$

If 2 classes: $P\left(\frac{C_1}{x_{\text{new}}}\right) \quad P\left(\frac{C_2}{x_{\text{new}}}\right)$

Model-based learning.

It will tell the class properties.

$$P(x/c) = \prod_{i=1}^K P(x_i/c)$$

$\hookrightarrow x_1, \dots, x_n$

$$P\left(\frac{C_1}{x_{\text{new}}}\right) = \frac{P(x_{\text{new}}/c_1) \times P(c_1)}{P(x)}$$

$$P\left(\frac{C_2}{x_{\text{new}}}\right) = \frac{P(x_{\text{new}}/c_2) \times P(c_2)}{P(x)}$$

2 more terms:

- ① Generative classifier
- ② Discriminative classifier

→ Naive-Bayes is Generative classifier

Code

Initial

`X_iris = iris.data[["petal length (cm)", "petal width (cm)"]].values`

`y_iris = iris.target`

`print(iris)`

`tree_clf = Dec.Tre.Clf.`

`(max_depth=,`

Assignment

② Apply DT on.

) regression

16/9/25

Plotting the tree:

```
matplotlib import pyplot as plt
from sklearn import tree
```

* By default, it gives binary split & uses Gini Index

~~Ensemble Learning~~

Naive Bayes

$$P(C/x) = P(x|C) \times P(C)$$

* Buy Computer Dataset

Age = y

Income = Medium

Student = Yes

CR = fair,

Class = ?

$$P(Yes) = \frac{9}{14}$$

$$P(No) = \frac{5}{14}$$

$$P(Age - Y/Yes) = \frac{2}{9}$$

$$P(CR - F/Yes) = \frac{6}{9}$$

$$P(Income - M/Yes) = \frac{4}{9}$$

$$P(St - Yes/Yes) = \frac{6}{9}$$

$$P\left(\frac{Y}{C}\right) = \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{6}{9}$$

Syllabus for Mid

→ Basic ML Challenges

→ P(E), A, P, R, Gini, k-cross, variance, bias

→ Life cycle of DS

→ Statistical → Descriptive → Probability

→ Inferential statistics

→ Linear Regression

→ Classification - K-NN

→ Decision tree
Naive Bayes

Time: 1 & half hr

Marks: 25 M

16/9/29 $P(Yes/X) = P(X/Yes) \times P(Yes)$ $\left[\because P(C/X) = P(C) \cdot P_C \right]$

$$= \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{6}{9} \times \frac{9}{14} = 0.028$$

$$P(Age - Y/No) = \frac{3}{5}$$

$$P(Salary - M/No) = \frac{2}{5}$$

$$P(ST - Yes/No) = \frac{1}{5}$$

$$P(CR-F/No) = \frac{2}{5}$$

$$P\left(\frac{X}{No}\right) = \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5}$$

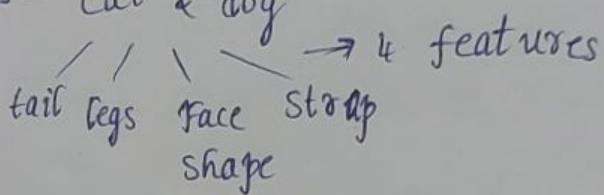
$$P\left(\frac{No}{X}\right) = P\left(\frac{X}{No}\right) \times P(No)$$

$$= \frac{3}{2} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{5}{14}$$

$$= 0.007$$

Generative Classifier: It respects each & every class.

Ex: 2 classes - Cat & dog



It will find the probability for each & every feature of cat & dog. They will be considered as indept

Discriminative Classifier: For new instance, it will check face & strap on neck. They are going to set some margin (or) decision boundary.