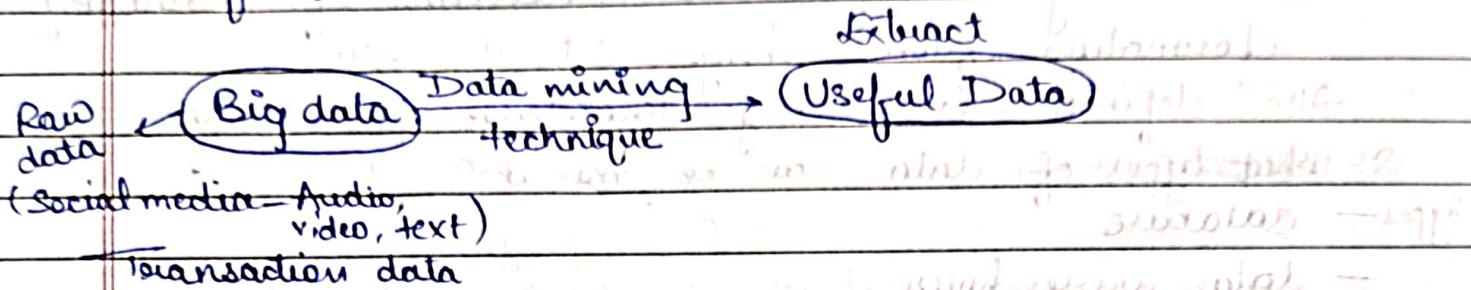


Data Mining CSE312A



(Valuable) Search)

→ It is the process of mining knowledge from large amount of data.



why we do data mining?

1. Companies and organisations get huge amount of data from different sources and platforms.
2. As size of database increasing & it is very difficult to manually search for useful information in it.
3. They use data mining technique which include AI & mathematical complex algo. for getting specific & useful data.
4. This specific data helps them for taking data-driven decisions.
5. We also get interesting, general direction in which data changes over time (+, -, flat)
 - (i) pattern → recognizable structure or trend in data that can be useful for making decision or predictions, insights of collected data.

→ Data Mining is also called as 'KDD' (Knowledge Discovery in Database).

Technique:-

Statistics	AI	ML
Cluster techniques :-	Diff.	→ KNN
	AI	→ Apriori
	Argo	→ K mean
	ANN	→ Naive bayes
	Back propagate	

Ex:- Price Selection in Amazon (0 - 1000)
Weather forecasting- (put data)

Data Mining is defined as procedure of extracting information from huge sets of data.

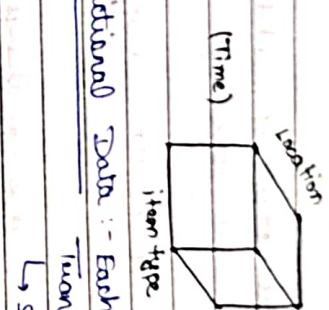
Also defined as mining knowledge from data.

Q. what types of data can be mined?

- data warehouse

- transactional

- time-series



3. Transactional Data :- Each record or attribute is called as Transaction.

↳ Sales, flight booking, user clicks on webpage)

Transaction has transaction ID, list of other items making transaction from transaction database, we can mine frequent patterns.

o Data Cleaning :-

↳ to remove noise and inconsistent data ex:- parsing the data

↳ cleaning is performed for detection of syntax errors

↳ Power divides whether the given string of data is acceptable within data specification

o Data Integration :-

where multiple data sources are combined

o Data Selection :-

where data relevant to the analysis task are retrieved

from the database.

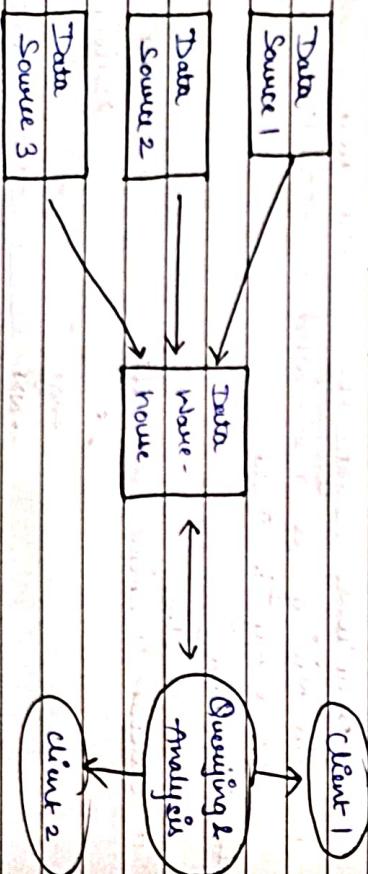
o Data Transformation :-

where data are transformed or consolidated into forms appropriate for mining by performing summarization or aggregation operation for instance.

o Data mining :- process where intelligent methods are applied in order to extract data pattern

o Pattern Evaluation:-

10 identify the truly interesting patterns representing knowledge base on some interesting measures.





knowledge representation :-

where visualization & knowledge representation technique are used to present the mind knowledge to the user.

ATTRIBUTE TYPES:- Data Mining (NOIR classification)

- Nominal :-
- Related to name

• The values are symbols or name of things.

- Represent category - code or state, also called categorical attribute

• Eg:- hair color.

→ Ordinal :-

- Represent a meaningful order or ranking (magnitude can't calculate)
- Eg:- first, 2nd, 3rd, ... → associate rank → rank

→ Binary :-

- Only 2 states or categories are:- '0' or '1'

• Boolean (True or False)

- 1 (present) / 0 (absent)

• Types:-

- (i) Symmetric - (equal weightage)

→ Gender (M/F)

- (ii) Transitive - (Uniqueness)

→ HIV Test (+/-)

→ Numeric :-

- Quantitative

• Measurable

- Represented by Integer or Real values.

• Types:-

- (i) Nominal scaled :-

→ measured on equal sized unit (linear)

→ have order & can be '+'ve or '-'ve

→ allow us to compare & quantify the difference b/w values

→ Eg:- Temperature

(iii) Ratio Scale :-

→ Continuous & '0' is allowed

→ measurement on a non-linear scale

NOIR Classification :-

Scale is used to represent usually quantitatively, an item's, person's or an event's place in the spectrum

OR

Spectrum

Scale is defined as the device on an object used to measure or quantity any event or an object

Types of Scales

1. Nominal Scale :-

* Nominal scale is used to categorize data into mutually exclusive categories or groups.

* A nominal scale is the lowest level of measurement

* A scale in which the numbers, serve as "tags" or "labels" to classify or identify the objects.

* The data can be placed into categories but can't be multiplied, divided, added or subtracted from one another.

* Ex:- gender, country, player's T-shirt number etc.

• Types:-

- (i) Nominal scaled :-

→ measured on equal sized unit (linear)

Qualitative (Nominal / ordinal)

~~scale carry mark on order the event (only Toleriful)~~
 Labels are assigned commt by ordered
 [Armen Number, Phone Number]

2. Ordinal Data :-

- Data is ordered from least to most
- Un-determined numerical distance b/w each cat category (First place, 2nd place, 3rd place)
- Ordinal scale is used to measure variables in a natural order, such as rating or ranking totally satisfied, satisfied, Neutral, dissatisfied.

3. Interval Scaled Data :- (Interval Identifier)

- A scale with values
- Same numerical distance b/w each value.
 Eg. Rating 1, 2, 3
- Has an arbitrary zero point.
 0°C , we have temp.

- Interval scale contains properties of nominal & ordered data, but the difference b/w data points can be quantified.
- This type of data shows both the order of the variables & the exact difference b/w the variables.

- For eg. - time.

-2 -1 0 1 2 3

Ratio Scale:- (No Negative Value)

- A scale with zeroes
- Ratio scales of measurement include properties from all four scales of measurement.

o The data is nominal and identity, can be classified in order, contains intervals & can be broken down into exact value.

- o Data in the ratio scale can be added, subtracted, divided & multiplied.
- o Because of the existence of true zero value, the ratio scale doesn't have negative values
- o Example:- weight, height, distance are all ex. of ratio variables
- o True meaningful zero point (eg:- income, city distance)

Graph Theory

Application :- Communication, Computation, Bio, chemistry

$$G = (V, E)$$

$V = \text{vertex (Node)} \{v_1, v_2, v_3, \dots, v_n\}$

$E = \text{edges (line)} \{e_1, e_2, e_3, \dots, e_m\}$

$$\text{O} \quad a \quad \text{O}$$

$$\text{O} \quad b \quad \text{O}$$

$$\text{O} \quad c \quad \text{O}$$

$$\text{O} \quad d \quad \text{O}$$

$$\text{O} \quad e \quad \text{O}$$

$$\text{O} \quad f \quad \text{O}$$

$$\text{O} \quad g \quad \text{O}$$

$$\text{O} \quad h \quad \text{O}$$

$$\text{O} \quad i \quad \text{O}$$

$$\text{O} \quad j \quad \text{O}$$

$$\text{O} \quad k \quad \text{O}$$

$$\text{O} \quad l \quad \text{O}$$

$$\text{O} \quad m \quad \text{O}$$

$$\text{O} \quad n \quad \text{O}$$

$$\text{O} \quad o \quad \text{O}$$

$$\text{O} \quad p \quad \text{O}$$

$$\text{O} \quad q \quad \text{O}$$

$$\text{O} \quad r \quad \text{O}$$

$$\text{O} \quad s \quad \text{O}$$

$$\text{O} \quad t \quad \text{O}$$

$$\text{O} \quad u \quad \text{O}$$

$$\text{O} \quad v \quad \text{O}$$

$$\text{O} \quad w \quad \text{O}$$

$$\text{O} \quad x \quad \text{O}$$

$$\text{O} \quad y \quad \text{O}$$

$$\text{O} \quad z \quad \text{O}$$

$$\text{O} \quad \alpha \quad \text{O}$$

$$\text{O} \quad \beta \quad \text{O}$$

$$\text{O} \quad \gamma \quad \text{O}$$

$$\text{O} \quad \delta \quad \text{O}$$

$$\text{O} \quad \epsilon \quad \text{O}$$

$$\text{O} \quad \zeta \quad \text{O}$$

$$\text{O} \quad \eta \quad \text{O}$$

$$\text{O} \quad \theta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

$$\text{O} \quad \lambda \quad \text{O}$$

$$\text{O} \quad \kappa \quad \text{O}$$

$$\text{O} \quad \rho \quad \text{O}$$

$$\text{O} \quad \sigma \quad \text{O}$$

$$\text{O} \quad \tau \quad \text{O}$$

$$\text{O} \quad \pi \quad \text{O}$$

$$\text{O} \quad \vartheta \quad \text{O}$$

$$\text{O} \quad \varphi \quad \text{O}$$

$$\text{O} \quad \psi \quad \text{O}$$

$$\text{O} \quad \chi \quad \text{O}$$

$$\text{O} \quad \omega \quad \text{O}$$

$$\text{O} \quad \nu \quad \text{O}$$

$$\text{O} \quad \mu \quad \text{O}$$

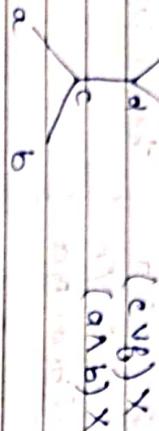


$(a,c) \cap a$

$(a,c) \rightarrow c \text{ or } a$

$(b,d) \rightarrow c \text{ or } a$

c Max in lattice



a

b

Mean :- the mean or average is the sum of all values

in a dataset divided by the no. of values.

- It represents the central point around which the data is distributed.

Ex:- Salary Package in Lakhs

$$2, 3, 4, 5, 8, 9, 32$$

$$\text{Mean} \Rightarrow 2+3+4+5+8+9+32 = 7$$

Median :- It is the middle value of a dataset when it is ordered.

$$\text{if odd no. of values} \Rightarrow \frac{\text{No.} + 1}{2}$$

$$\text{if even no. of values} = \frac{N_i + N_j}{2}$$

Node :- The mode is the value that appears most frequently in a dataset

$$\text{Ex:- } 2, 3, 4, 3, 5, 7, 3, 8, 30$$

$$\text{Mode} = 3$$

$$\text{Ex:- } 2, 3, 4, 5, 7, 8, 30$$

$$\text{Mode} = 0$$

Ex Real time :- E-commerce

Range :- It is a measure of the spread or dispersion of a dataset.

- It is defined as the difference b/w the maximum & minimum values in a set of observation.

$$\text{Ex:- } 2, 3, 4, 5, 8, 9, 32$$

Min

Max

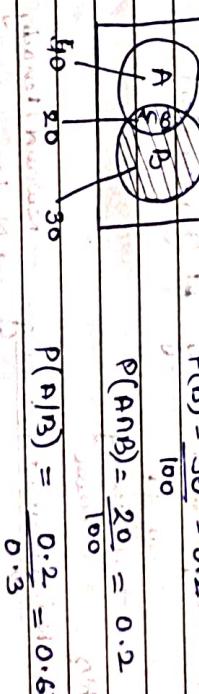
$$\text{Max - Min} = 32 - 2 = 30$$

Conditional Probability Event A & B.

$$P\left(\frac{A}{B}\right) = \frac{P(A \cap B)}{P(B)}$$

↳ Probability of event A given that event B is already occur.

Ex:-



Q:- In a group of 100 computer buyers, 40 bought CPU, 30 purchased monitor, and 20 purchased CPU and monitor. If a computer buyer chose at random and bought a CPU, what is the probability they also bought a Monitor?

Ans:- A = bought computer = 40
B = purchased monitor = 30
 $P(A) = \frac{40}{100} = 0.4$
 $P(B) = \frac{30}{100} = 0.3$



Steps of Supervised learning Algorithms:-

- Find determine the type of training dataset.
- Collect/Gather the labelled training data.
- Split the training dataset into training dataset, test dataset, validation dataset.

- Determine the sfp features of the training dataset, which should have enough knowledge so that the model can accurately predict the sfp.
- Support vector machine, decision tree etc.
- Execute the algo. on training dataset.
- Evaluate the accuracy of model by providing the test set.
- If the model predicts the correct sfp, which means our model is accurate.

Types of supervised learning algo.

Ex:-
1. Classification
2. Regression

- Regression
Ex:-
What is the temp going to be tomorrow?
Will it be cold or hot tomorrow?
- Classification
Ex:-
Is it going to rain tomorrow?

Regression: now used if there is a relationship b/w sfp variable and the sfp variable.

used for prediction of continuous & real variable. or to predict a numerical value based on sfp features.

Ex:- Weather Forecasting, Market Trends

Temp, Age, Salary, Price

Regression Analysis Algorithm:-

- Linear Regression
- Logistic Regression
- Support Vector Regression
- Decision Tree Regression
- Random Forest Regression

Classification:- are used when sfp variable is categorical, which means there are 2 classes such as Yes-No,

No- female, True - False etc.

The goal is to assign sfp data to predefined categories or classes.

Ex:- 1. Determining an email is spam or not.

2. The outcome can be a binary classification

3. Logistic regression

4. Rainfall tomorrow or not.

5. Something is harmful or not.

6. Support Vector Machines.

7. Decision Trees

8. Random Forest

9. Logistic Regression

10. SVM

11. DT

12. RF

13. LR

14. KNN

15. NB

16. QDA

17. GNB

18. SGD

19. K-Means

20. PCA

21. NN

22. RF

23. DT

24. LR

25. SVC

26. KNN

27. QDA

28. GNB

29. SGD

30. PCA

31. NN

32. RF

33. DT

34. LR

35. SVC

36. KNN

37. QDA

38. GNB

39. SGD

40. PCA

41. NN

42. RF

43. DT

44. LR

45. SVC

46. KNN

47. QDA

48. GNB

49. SGD

50. PCA

51. NN

52. RF

53. DT

54. LR

55. SVC

56. KNN

57. QDA

58. GNB

59. SGD

60. PCA

61. NN

62. RF

63. DT

64. LR

65. SVC

66. KNN

67. QDA

68. GNB

69. SGD

70. PCA

71. NN

72. RF

73. DT

74. LR

75. SVC

76. KNN

77. QDA

78. GNB

79. SGD

80. PCA

81. NN

82. RF

83. DT

84. LR

85. SVC

86. KNN

87. QDA

88. GNB

89. SGD

90. PCA

91. NN

92. RF

93. DT

94. LR

95. SVC

96. KNN

97. QDA

98. GNB

99. SGD

100. PCA

101. NN

102. RF

103. DT

104. LR

105. SVC

106. KNN

107. QDA

108. GNB

109. SGD

110. PCA

111. NN

112. RF

113. DT

114. LR

115. SVC

116. KNN

117. QDA

118. GNB

119. SGD

120. PCA

121. NN

122. RF

123. DT

124. LR

125. SVC

126. KNN

127. QDA

128. GNB

129. SGD

130. PCA

131. NN

132. RF

133. DT

134. LR

135. SVC

136. KNN

137. QDA

138. GNB

139. SGD

140. PCA

141. NN

142. RF

143. DT

144. LR

145. SVC

146. KNN

147. QDA

148. GNB

149. SGD

150. PCA

151. NN

152. RF

153. DT

154. LR

155. SVC

156. KNN

157. QDA

158. GNB

159. SGD

160. PCA

161. NN

162. RF

163. DT

164. LR

165. SVC

166. KNN

167. QDA

168. GNB

169. SGD

170. PCA

171. NN

172. RF

173. DT

174. LR

175. SVC

176. KNN

177. QDA

178. GNB

179. SGD

180. PCA

181. NN

182. RF

183. DT

184. LR

185. SVC

186. KNN

187. QDA

188. GNB

189. SGD

190. PCA

191. NN

192. RF

193. DT

194. LR

195. SVC

196. KNN

197. QDA

198. GNB

199. SGD

200. PCA

201. NN

202. RF

203. DT

204. LR

205. SVC

206. KNN

207. QDA

208. GNB

209. SGD

210. PCA

211. NN

212. RF

213. DT

214. LR

215. SVC

216. KNN

217. QDA

218. GNB

219. SGD

220. PCA

221. NN

222. RF

223. DT

224. LR

225. SVC

226. KNN

227. QDA

228. GNB

229. SGD

230. PCA

231. NN

232. RF

233. DT

234. LR

235. SVC

236. KNN

237. QDA

238. GNB

239. SGD

240. PCA

241. NN

242. RF

243. DT

244. LR

245. SVC

246. KNN

247. QDA

248. GNB

249. SGD

250. PCA

251. NN

252. RF

253. DT

254. LR

255. SVC

256. KNN

257. QDA

258. GNB

259. SGD

260. PCA

261. NN

262. RF

263. DT

264. LR

265. SVC

266. KNN

267. QDA

268. GNB

269. SGD

270. PCA

271. NN

272. RF

273. DT

274. LR

275. SVC

276. KNN

277. QDA

278. GNB

279. SGD

280. PCA

281. NN

282. RF

283. DT

284. LR

285. SVC

286. KNN

287. QDA

288. GNB



INPUT

Function
Data Point

Date _____
Page _____

Output
Label

Unknown, real-world process

$$\begin{array}{l} \text{True Labels} \\ y_1 = f(x_1) \\ y_2 = f(x_2) \end{array}$$

$$\begin{array}{l} \text{Model: } f \\ \text{Predicted Labels} \\ y_1 = f(x_1) \\ y_2 = f(x_2) \end{array}$$

Our Goal:- Find f such that

$$y_1 \approx y_1$$

$$y_2 \approx y_2$$

ADVANTAGES:- * With the help of supervised learning, the model can predict the O/P on basis of prior experience.

* We can have an exact idea about the classes of object

* It helps us to solve various real-world problems such

as fraud detection, spam filtering etc.

DISADVANTAGES:- * Are not suitable for handling the complex task

* Cannot predict the correct O/P if test data is diff from the training dataset

* Training required a lots of computations times

* We need enough knowledge about classes of object

Regression

Classification

1. The system attempts to predict a value from SP based on past data

2. Real No./Continuous No - Regression problem

3. Ex:- Temp. for Tomorrow

Classification

1. In classification, predictions are made by classifying them into diff categories.

2. Discrete / Categorical Variable -

3. Ex:- Type of cancer

Concern Y/N.

Supervised Learning

Unsupervised Learning

- These are trained using labelled data
- In this, O/P data is provided to the model along with the O/P

- It predicts the O/P
- In this, only O/P data is provided to the model.

- Finds the hidden pattern in data
- It does not require any supervision to train the model.

- It produces an accurate result
- It is less complex and less time consuming fast
- The goal is to identify relationship b/w the O/P and O/P Variables and then use it to map new unlabeled data
- It is used for price predictions, weather forecasting
- It gives less accurate result
- It is more complex and relatively slower
- The goal of this is to find pattern/meaningful pattern in data based on the relationship b/w data point themselves.
- Used for customer personal recommendation engines.

- ### # BOOLEAN FUNCTIONS:-
- A function $f: A^n \rightarrow A$ is called Boolean function if it can be specified by a boolean expression of n variables. $f(A) = \{x_1, x_2, \dots, x_n\}$ & f operations t, \neg, \wedge, \vee
- $$(x_1 \wedge x_2 \wedge x_3) + x_2 \cdot x_1 \quad \text{OR} \quad u, v (x_2 \wedge x_3) \vee (x_2 \wedge x_1)$$
- Boolean expression \rightarrow Boolean expression in n variables of 3 variables
- Boolean Algebra**- is a branch of algebra in which values T/F usually of variable are truth values T/F usually



→ Problems in ARN:-

* levels of frequency of appearance determination

* finding strong association among frequent items.

→ Functions of ARN:-

* finding set of items that has significant impact on business.

* collating information from numerous transactions

* Generating rules from counts in Transactions

→ Strength of ARN → Weakness

1. Long interpretation

2. Exponential Growth in Complications

3. Flexible data formats

4. Simplicity

5. Rule selection

6. Rare items → not applicable for

*** APRIORI ALGORITHM :-

Idea is to generate candidate itemsets of a given size.

and then scan dataset to check if their counts are within range.

• the process is iterative

Steps :- 1. All singleton items are candidates in the first pass. Items with less than specified support value is eliminated.

2. Two member candidate itemset.

3. Then make 3 member candidate itemsets upto 'n' no. of members.

4. frequent itemsets conditions, set of frequent itemset.

5. we generate association Rules which have confidence values greater than or equal to specified min_r confidence

Confidence

The Apriori algorithm is a popular method used for mining frequent itemsets.

Eg:- Tid Items Support

Tid	Items	Support
1	2,3	1
2	1,3,5	1
3	1,2,4	2
4	2,3	3
5		1

1	2,3	2
2	1,3,5	1
3	1,2,4	2
4	2,3	3

1	2,3	1
2	1,3,5	1
3	1,2,4	2
4	2,3	3

1. Eg:-

2,3	2
1,3,5	1
1,2,4	2
2,3	3

Real-life Ex:- Market Basket Analysis

A supermarket wants to understand the purchasing behaviour of its customers to optimize product placement & promotion.

Process → Data collection :- the supermarket collects transaction data,

where each transaction is a list of items purchased together

→ frequent Itemset :- set of items that frequently

appear together in transactions.

Application:- * Product Placement

* Promotion Strategy

* Inventory Management

2. Eg:- Transaction ID . Items Purchased

T₁ Milk, Bread, Butter

T₂ Bread, Butter, Jam

T₃ Milk, Bread, Butter, Jam

T₄ Butter, Jam

T₅ Milk, Bread

3. Initialize the Candidate 1 - itemset c₁

Items support

Milk → 3

Bread → 4

Butter → 4

Jam → 2

3. Make pair to generate C₃

$$C_3 = \text{Itemset Sup. Count}$$

$$M_{1,K,E} \rightarrow 1 \times$$

$$M_{1,K,V} \rightarrow 2 \times$$

$$O_{1,K,E} \rightarrow 3$$

$$O_{1,K,V} \rightarrow 2 \times$$

Now again compare with min support

$$L_3 = \text{Itemset Supp. Count}$$

$$O_{1,K,E} \rightarrow 3$$

* Now create association rules with Support & Confidence

$$\text{for } O_{1,K,E}$$

$$\text{Association Rule} \quad \text{Support} \quad \text{Confidence}$$

$$O_{1,K} \rightarrow E \quad 3 \quad 3/3 = 1$$

$$O_{1,E} \rightarrow K \quad 3 \quad 3/3 = 1$$

$$K \wedge E \rightarrow O \quad 3 \quad 3/4 = 0.75$$

$$E \rightarrow O \wedge K \quad 3 \quad 3/4 = 0.75$$

$$K \rightarrow O \wedge E \quad 3 \quad 3/5 = 0.6$$

$$O \rightarrow K \wedge E \quad 3 \quad 3/5 = 0.75$$

Compare this with min confidence = 80%.

$$O \wedge K = E \quad 3 \quad 100$$

$$O \wedge E = K \quad 3 \quad 100$$

Never final association rules are

$$O \wedge K = E \rightarrow \text{market basket analysis.}$$

Ex 4:- T_{10} item bought

$$T_1 \quad \{ \text{Bread, Butter, Milk}\}$$

$$T_2 \quad \{ \text{Bread, Butter}\}$$

$$T_3 \quad \{ \text{Bread, Diapers}\}$$

$$T_4 \quad \{ \text{Milk, Diapers, Bread}\}$$

$$T_5 \quad \{ \text{Bread, Diapers}\}$$

For given dataset, apply Apriori algo. to discover stronger association rules among item tags.

min support = 40%.

min confidence = 70%.

$$\text{Sol:-} \quad \text{Support} = \frac{240}{400} \times 5 = \underline{\underline{2}}$$

(i) Generate frequent itemset

$$C_1 = \text{Itemset Sup. Count}$$

$$\text{Bread} \rightarrow 3$$

$$\text{Nuk} \rightarrow 2$$

$$\text{Beer} \rightarrow 2$$

$$\text{Butter} \rightarrow 3$$

$$\text{Diapers} \rightarrow 3$$

C₂ Itemset Sup. Count

$$\text{Bread, Butter} \rightarrow 3$$

$$\text{Bread, Nuk} \rightarrow 2$$

$$\text{Bread, Beer} \rightarrow 2$$

$$\text{Bread, Diapers} \rightarrow 1 \times$$

$$\text{Bread, Nuk} \rightarrow 2$$

$$\text{Butter, Beer} \rightarrow 0 \times$$

$$\text{Butter, Diapers} \rightarrow 1 \times$$

$$\text{Butter, Nuk} \rightarrow 2 \times$$

$$\text{Butter, Beer} \rightarrow 0 \times$$

$$\text{Diapers, Nuk} \rightarrow 1 \times$$

$$\text{Diapers, Beer} \rightarrow 0 \times$$

$$\text{Nuk, Beer} \rightarrow 0 \times$$

(iii) Generate pair to generate C₃

$$C_3 = \text{Itemset Sup. Count}$$

$$\text{Bread} \rightarrow 3$$

$$\text{Butter} \rightarrow 2$$

$$\text{Diapers} \rightarrow 2$$

$$\text{Nuk} \rightarrow 2$$

$$\text{Beer} \rightarrow 2$$

Need of closed and maximal itemset :-

Bread, Milk, Ben \rightarrow 0 X
 Bread, Diaper, Beer \rightarrow 0 X
 Butter, Milk, Diaper \rightarrow 1 X

Butter, Milk, Ben \rightarrow 0 X
 Butter, Diaper, Beer \rightarrow 0 X

Milk, Diaper, Ben \rightarrow 0 X
 So L3 =

Bread, Butter, Milk \rightarrow 2 X

* Now create association rule with support & confidence.
 few Butter, Bread, Milk.

Association Rule

Support

Confidence

confidence.

Bread, Butter \rightarrow Milk 2
 Bread \wedge Milk \rightarrow Butter 2

$2/3 = 0.66$

Go X

Bread \wedge Milk \rightarrow Bread 2

$2/2 = 1$

Go X

Bread \wedge Milk \rightarrow Bread 2
 Milk \rightarrow Bread \wedge Butter 2

$2/2 = 1$

Go X

Butter \rightarrow Bread Milk 2
 Butter \rightarrow Bread Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Bread 2
 Butter \wedge Milk \rightarrow Butter 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Milk 2
 Butter \wedge Milk \rightarrow Butter 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

Butter \wedge Milk \rightarrow Butter 2
 Butter \wedge Milk \rightarrow Milk 2

$2/3 = 0.66$

Go X

\rightarrow These closed frequent itemset and maximal frequent itemset are useful when huge amount of data is used in association rule mining.

\rightarrow If the length of frequent itemset is ' k ', then by downward closure property all of its 2^k subsets are also frequent.

\rightarrow When the computation is very expensive & time is no interest to find additional subsets. This can be avoided by frequent itemset with max-length.

\rightarrow One disadvantage with maximal frequent itemsets is that when all its subsets are frequent, we do not know their support for mining next, support information is very important.

\therefore closed freq. itemset is preferred.

Frequent itemsets

Closed freq. itemset

Min. support count = 3
 find frequent, closed, Maximal Itemset

Tid List of items

A, B, C, D

Maximal freq. itemset

closed freq. itemset

All items that is A, B, C, D are frequent because they sup. count is greater than or eq. to minimum support count.

C1 : \rightarrow Item Count
 A \rightarrow 3
 B \rightarrow 4
 C \rightarrow 5
 D \rightarrow 4

A, B \rightarrow 3
 A, C \rightarrow 3
 A, D \rightarrow 2
 B, C \rightarrow 4
 B, D \rightarrow 3
 C, D \rightarrow 4

A, B, C \rightarrow 3
 A, B, D \rightarrow 3
 A, C, D \rightarrow 2
 B, C, D \rightarrow 4

A, B, C, D \rightarrow 3
 B, C, D \rightarrow 4

A (Count) is not greater than its immediate superset.

A (Count) is not closed.

T₅ bread, juice, milk.

Freq. Itemset

Items

frequency support

Bread → 4
juice → 3
milk → 3
egg → 1

$\frac{4}{15} = 80\%$

$\frac{3}{15} = 60\%$

$\frac{3}{15} = 60\%$

$\frac{1}{15} = 20\%$

X eliminate

Juice → 3
milk → 3
Yogurt → 1

$\frac{3}{15} = 60\%$

$\frac{1}{15} = 20\%$

X eliminate

Support(Bread) = $\frac{4}{15} = 26.7\% = 0.8 = 80\%$

$m = 5$

* Make 2-item candidate set & write them Frequency.

Item Freq. Support

bread, juice → 2

$\frac{2}{15} = 40\%$

X

bread, milk → 2

$\frac{2}{15} = 40\%$

X

bread, juice, milk → 1

$\frac{1}{15} = 20\%$

X

juice, milk → 2

$\frac{2}{15} = 40\%$

X

juice, milk → 3

$\frac{3}{15} = 60\%$

X

juice, milk, milk → 1

$\frac{1}{15} = 20\%$

X

Bread, juice → 3

$\frac{3}{15} = 60\%$

X

Bread, juice, milk → 1

$\frac{1}{15} = 20\%$

X

Bread, juice, milk, milk → 1

$\frac{1}{15} = 20\%$

X

Bread, juice, milk, milk, milk → 1

$\frac{1}{15} = 20\%$

X

Juice → 3

$\frac{3}{15} = 60\%$

X

milk → 3

$\frac{3}{15} = 60\%$

X

juice, milk → 2

$\frac{2}{15} = 40\%$

X

juice, milk, milk → 1

$\frac{1}{15} = 20\%$

X

Frequent Itemsets'

- An itemset x is frequent only if $\text{support}(x) \geq \text{min-sup}$.
- $\text{Support}(x) = \text{No. of transac. contains } x$.

Itemset Support F/I

A 3/5 F I

ABD 0/5 I

2/5 F

ACD 1/5 I

1/5 F

ACE 2/5 I

1/5 F

BCD 0/5 I

0/5 F

BCE 3/5 F

3/5 F

BCD 0/5 I

0/5 F

CD 1/5 I

1/5 F

CDE 0/5 I

0/5 F

ABC 0/5 I

0/5 F

ABCDE 2/5 F

2/5 F

BCD 3/5 F

3/5 F

BD 0/5 I

0/5 F

BE 4/5 F

4/5 F

CD 1/5 I

1/5 F

ABDE 0/5 I

0/5 F

CE 3/5 F

3/5 F

DE 0/5 I

0/5 F

ABC 2/5 F

2/5 F

ABCD 0/5 I

0/5 F



- Total itemsets = 31
- Frequent Itemset = 15
- In frequent Itemsets = 16.

- $\text{support}_{\text{min}} = 2/5$
- $\text{Total minsup} = 2/5$

- $\text{support}(J) = 5/4$
- $\text{Juice} \rightarrow \text{Juice}$ OR $\text{Juice} \rightarrow \text{chew}$
- $\frac{3}{8} \cdot \frac{5}{4} = 75\%$

- all rules are good.

- $Q^5 - 1$



Prediction— It is a process of identifying the missing or available numerical data for new obj.

- ↳ Algorithm which use training data-set to derive a model, that model is predictor. When the new data is given, this model should find a numerical op.
- ↳ Unlike in classification, this model does not have the class label, thus model predicts a continuous valued function or output value.

Issues regarding classification & Prediction.

- * Preparing the data for classification & prediction
- * Performance evaluation measures.

↳ Preparing the data for classification & prediction!

Initially, data should be pre-processed to improve the accuracy, efficiency, scalability of the classification or prediction process.

→ Pre-processing techniques:

1. Data cleaning: Removal of noise data.

2. Relevance analysis: Removal of redundant attributes using correlation

• Removing irrelevant attributes using Attribute subset selection.

3. Data transformation and reduction:-

• Normalization:- Scaling of values from a given attribute so that they fall within a small specified range.

Ex:- min., max., Z-score.

• Generalization:- Generalizing the data to higher conceptual levels. Particularly it can be useful for continuous valued attributes.

Ex:- attribute "income" can be generalized as: low, medium, high.

4. Performance Evaluation measures:-

1. Accuracy of a classifier refers to the ability of a given classifier to correctly predict the class labels of unseen data.

• Accuracy of predictor refers to how well a given predictor.

can guess the value of the predicted attributes for unseen data.

• Accuracy can be estimated using the test set that is independent of training set.

2. Speed:- refers to the computational cost involved in generating

2. using the given classifier or predictor.

3. Robustness:- the ability of classifier or predictor to make correct predictions given noisy data or data with missing values.

4. Scalability refers to ability to construct the classifier or predictor efficiently given large amount of data.

5. Interpretability:- Refers to level of understanding and insight that is provided by classifier or the predictor.

Classification & Prediction issue:-

The major issue is preparing the data for classification & prediction. Preparing the data involves the following activities:

1. Data Cleaning:- It involves removing the noise, and treatment of missing values. The noise is removed by applying

smoothing techniques and the problem of missing values is solved by replacing a missing values with common occurring value.

2. Relevance analysis:- DB may also have the irrelevant attributes.

Correlation analysis is used to know whether any of the given attributes are related.

2. Given attributes are related by any of the following methods:

• Normalization:- It involves scaling all values for given attributes in order to make them fall

within a small specified range. It is used when in the learning step, the natural laws or the methods involving measurements are used.

• Generalization:- Generalizing data to higher concept for this purpose we can use the concept hierarchy.

decision-making purposes.

The decision tree creates classification or regression

models as a true structure.

• At this point, the decimal train is separated into two smaller subsets,

steadily developed. The final tree is a tall

with the decision nodes & leaf nodes. A decision

node has at least two branches. The leaf nodes show a

class reunion or graduation. We can't accomplish more split on seat nodes. The ultimate division is to

a tree that relates to the best predictor called

the next node. Decision trees can deal with both categorical and numerical data.

July 1900
Winnipeg
Sask.

impunity. In the decision time, it measures the

randomness or impurity in data sets.

information gain refers to the decrease in entropy after the datum is split. It is also called Entropy.

Reduction. Building a decision tree is all about.

Discovering attributes that return the highest d_{G} .

Shallowing → Piping

$$= E_1 - E_2$$

In Sport - decision making
Pens. → entropy E_2 where E_2 =
Pens. → entropy E_2 where E_2 =

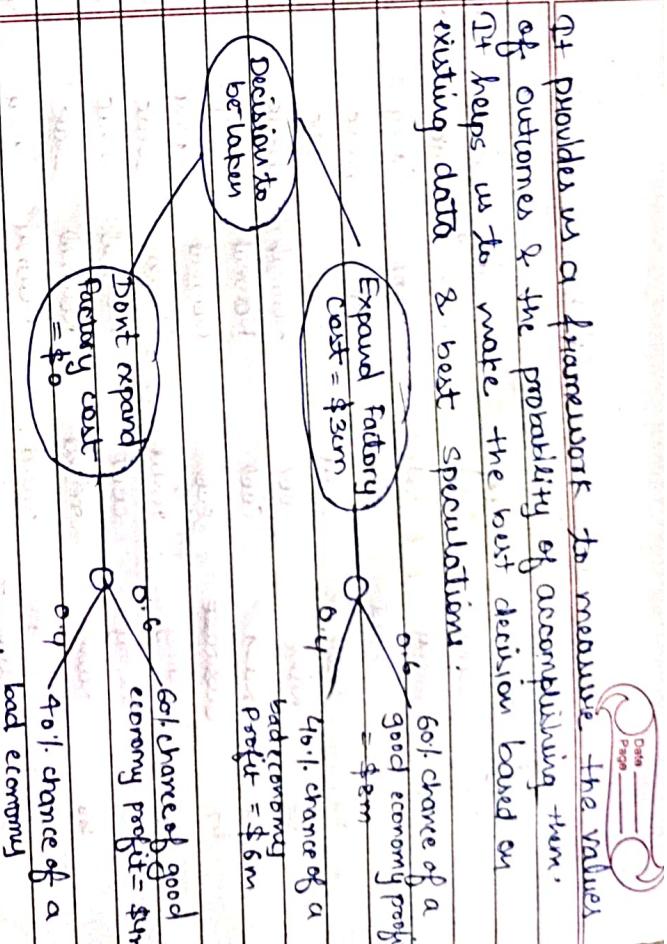
..., a decision tree is just like a flowchart.

Starting with d_{latent} , we can measure the entropy

To find a way to segment the set until the

~~shouldn't we have done this?~~

It enables us to analyse the possible consequences of a decision thoroughly.



Naïve Bayes classification :- (Probability based)



N

Attributes : Target variable

	Outlook	Temp.	Humidity	Wind	Play Tennis
D ₁	Sunny	Hot	High	Weak	No
D ₂	Sunny	Hot	High	Strong	No
D ₃	Overcast	Hot	High	Weak	Yes
D ₄	Rain	Cool	Normal	Weak	Yes
D ₅	Rain	Cool	Normal	Strong	No
D ₆	Overcast	Cool	Normal	Strong	Yes
D ₇	Sunny	Mild	High	Weak	No
D ₈	Sunny	Cool	Normal	Weak	Yes
D ₉	Rain	Mild	Normal	Weak	Yes
D ₁₀	Sunny	Mild	Normal	Strong	Yes
D ₁₁	Overcast	Mild	Normal	Weak	Yes
D ₁₂	Rain	Mild	High	Strong	Yes
D ₁₃	Sunny	High	Normal	Weak	No
D ₁₄	Rain	High	Strong	Yes	

Humidity	Y	N
High	319	415
Normal	619	115
Wind	Y	N
Weak	619	215
Strong	319	315

$$V_{NB}(\text{Yes}) = V_{NB}(\text{No})$$

$$= 0.0206 + 0.0206 = 0.795$$

$$= 0.0259$$

$$V_{NB} = \arg \max_{j} P(Y|v_j) P_i(\text{Play } v_j)$$

$$\text{Naïve Bayes} = \arg \max_j P(v_j) P(\text{Sunny}/v_j) P(\text{Rain}/v_j)$$

$$P(\text{Play Tennis} = \text{yes}) = 9/14 = 0.64$$

$$P(\text{Play Tennis} = \text{No}) = 5/14 = 0.36$$

Now calculate conditional probability of ~~new~~ individual current attributes

$$V_{NB}(\text{Yes}) = P(\text{Yes}) P(\text{Sunny}/\text{Yes}) P(\text{Rain}/\text{Yes})$$

$$= 9/14 \times 2/9 \times 319/319 \times 3/9 = 0.0053$$

Attribute	Possible outcome
outlook	Y N
sunny	2/9 3/15
overcast	4/9 0
rain	3/9 2/15



$$V_{\text{large}} = P(\text{no}) P(\text{sunny no}) P(\text{cool/no}) \\ P(\text{high/no}) P(\text{strong/no})$$

$$\Rightarrow \frac{5}{14} \times \frac{3}{5} \times \frac{1}{15} \times \frac{4}{15} \times \frac{3}{15} = \frac{0.0206}{\text{large}}$$

$$V_{\text{large}} = \frac{V_{\text{large}}(\text{yes})}{V_{\text{large}}(\text{yes}) + V_{\text{large}}(\text{no})}$$

$$= \frac{0.0053}{0.0053 + 0.206} = .205.$$

$$V_{\text{large}}(\text{no}) = \frac{V_{\text{large}}(\text{no})}{V_{\text{large}}(\text{yes}) + V_{\text{large}}(\text{no})} = \frac{0.0206}{0.0053 + 0.206}$$

Large

- * Data Analytics
- * Data Engineering
- * Data Science
- * Data Visualisation (Software used for data visualisation)
- * KDD.

Plan it with No filter

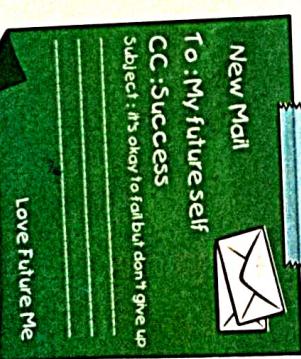
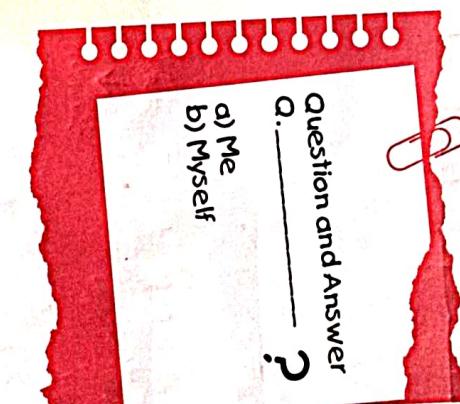
Goal

Checklist
(index)



Question and Answer
Q. ?

- a) Me
b) Myself



To : My future self
CC : Success
Subject : It's okay to fail but don't give up

Good Vibes

