

**Teacher Initial:**  
**NA**

**Daffodil International University**

**Fall 2021**

**Department of Computer Science and Engineering**

**Midterm Open Book Examination Answer Script**

**Full Marks: 25 Allowed, Time: 2:30 hrs.**

**Date: Thursday 18, November 2021**

**General Information (must be filled by the student)**

COURSE CODE: **CSE 450**

SECTION: **PC-B**

PROGRAM: **DAY EVEN**

STUDENT ID: **181-15-1905**

TIME STARTED:

TIME ENDED:

**[Student must either TYPE or HAND WRITE the answers in this template; In case needed just write your detail on the paper using hand]**

**\*\* Plagiarism will be checked while you submit your response. You are advised to be honest during the open book exam.**

Date .....

### Answers to the question no. 1

a) Corona create a huge change in our daily life, economy, our personal life. In this new normal we need to reconsider a lot of activities we did previous.

By data analysis we can take a lot of decision. Analysing our future data we can predict future step and upcoming effect of predicted activities.

For detecting post corona challenges, data analysis method are highly impactful. By collecting recent data our data analyst will take care of that data and can give us impactful

**Pantonix®**  
Pantoprazole

Date .....

information from that data. Now here some data analysis method that help to gain predict our information.

cluster analysis that helps to identify a group. and from that group we can gain information.

regression analysis is get a continuous value that is dependent on some other independent value.

Neural network are learn like human. it segment data into different part and learn from that give us useful information. from using those methodology we can get useful information that help's to fight

**Pantonix®**  
Pantoprazole

• Of all we need to collect data. Without collecting huge number of data we will not get a perfect prediction. Again only data can not give us perfect information. For get a good and productive information we need a productive data. For getting a productive data we need to preprocess data. By preprocess data we remove noise from data, remove duplicate data, reduce attribute.

-There are some data preprocessing technique,

- i) Aggregation
- ii) Sampling
- iii) Dimensionality Reduction
- iv) Feature subset selection
- v) Feature creation
- vi) Discretization and Binarization
- vii) Attribute transformation.

Without preprocessing we will get messy numbers of data, with unnecessary information. That will unhealthy for gain information from data. A set of data may contain extra attribute (ex. id, tag, phone etc) those are not necessary for information gain.

Date .....

So data preprocessing is very important.

Data attributes are mainly numbered and symbols to form

for health section our data attribute may be, Sex, Age, Phone, Area of living, Occupation, Salary, Total monthly cost, etc. Those are the data attribute and those attribute have different types of values.

For preprocessing those value we may need to clean some values. Need to re populate missing values. May remove some aggregation, attributes and also aggregate some attributes to one attribute.

**Pantonix®**  
Pantoprazole

Answers to the question no 2

Attribute class → is Transportation Mode,  
Hence,

Bus = 3

Train = 3

Car = 3

$$\text{Entropy}(q_B, q_T, q_C) = \left( -\frac{4}{10} \log_2 \frac{4}{10} - \frac{3}{10} \log_2 \frac{3}{10} - \frac{3}{10} \log_2 \frac{3}{10} \right)$$

$$= -0.53 + 0.52 + 0.52$$

$$= 0.57$$

Date .....

Entropy for child attribute  
Gender,

$$\text{Male} = \frac{1}{5}$$

$$\text{Female} = \frac{4}{5}$$

for male,

$$\text{Entropy } (B, C, T)_2 = -\frac{3}{5} \log_2(\frac{3}{5}) - \frac{1}{5} \log_2(\frac{1}{5}) - \frac{1}{5} \log_2(\frac{1}{5})$$

$$= 0.49 + 0.96 + 0.96$$

$$= 1.36$$

$$\text{for female } (B, D, T, C)_2 = -\frac{1}{5} \log_2(\frac{1}{5}) - \frac{2}{5} \log_2(\frac{2}{5}) - \frac{2}{5} \log_2(\frac{2}{5})$$

$$= 0.96 + 0.53 + 0.53$$

$$= 1.52$$

$$\begin{aligned}\text{Gain (for Gender)} &= \left( 1.57 - \left( \frac{5}{10} \times 1.36 + \frac{5}{10} \times 1.52 \right) \right) \\ &\geq (1.57 - (0.68 + 0.76)) \\ &\geq 0.13\end{aligned}$$

**Pantonix®**  
Pantoprazole

Date .....

Entropy using rule of weighted average,

Here,

$$0 = 3$$

$$1 = 5$$

$$2 = 2$$

Entropy for 0,

$$\text{Entropy } (2B, 1T, 0C) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) - 0 \log_2 (0)$$

$$= 0.38 + 0.52$$

$$= 0.9$$

$$\rightarrow \text{Entropy } (2B, 2T, 1C) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{1}{5} \log_2 \left(\frac{1}{5}\right)$$

$$= 0.53 + 0.53 + 0.96$$

$$= 2.02$$

$$2 \rightarrow \text{Entropy } (0, 0, 2C) = -0 - 0 - \frac{2}{2} \log_2 \left(\frac{1}{2}\right)$$

$$= 4.0$$

**Pantonix®**  
Pantoprazole

$$\text{Gain (for CO)} = 1.57 - \left( \frac{3}{10} \times 0.7 + \frac{5}{10} \times 1.52 + \frac{2}{15} \times 0 \right)$$

$$= 1.57 - 1.03$$

$$= 0.54$$

Entropy for Travel cost,

$$\text{Cheap} = 5$$

$$\text{Standard} = 2$$

$$\text{Expensive} = 3$$

$$\text{Cheap} \Rightarrow \text{Entropy}(q_{3,15,0}) = \frac{4}{5} \log_2 \left(\frac{4}{5}\right) - \frac{1}{5} \log_2 \left(\frac{1}{5}\right) - \frac{6}{5} \log_2 \left(\frac{6}{5}\right)$$

$$= 0.25 + 0.96$$

$$= 0.71$$

$$\text{Standard} \Rightarrow \text{Entropy}(2,10,6) = -\frac{2}{2} \log_2 \left(\frac{1}{2}\right)$$

$$= 0$$

Date .....

Expensive  $\Rightarrow$  Entropy  $(30, 0, 0) = 0$

$$\begin{aligned} \text{Gain} &= 1.57 - (5/10 \times 0.7 + 0.10) \\ &= 1.57 - 0.35 \\ &= 1.22 \end{aligned}$$

Entropy for Income Level,

$$\begin{aligned} \text{Low} &\approx 2 \\ \text{medium} &\approx 6 \\ \text{High} &\approx 2 \end{aligned}$$

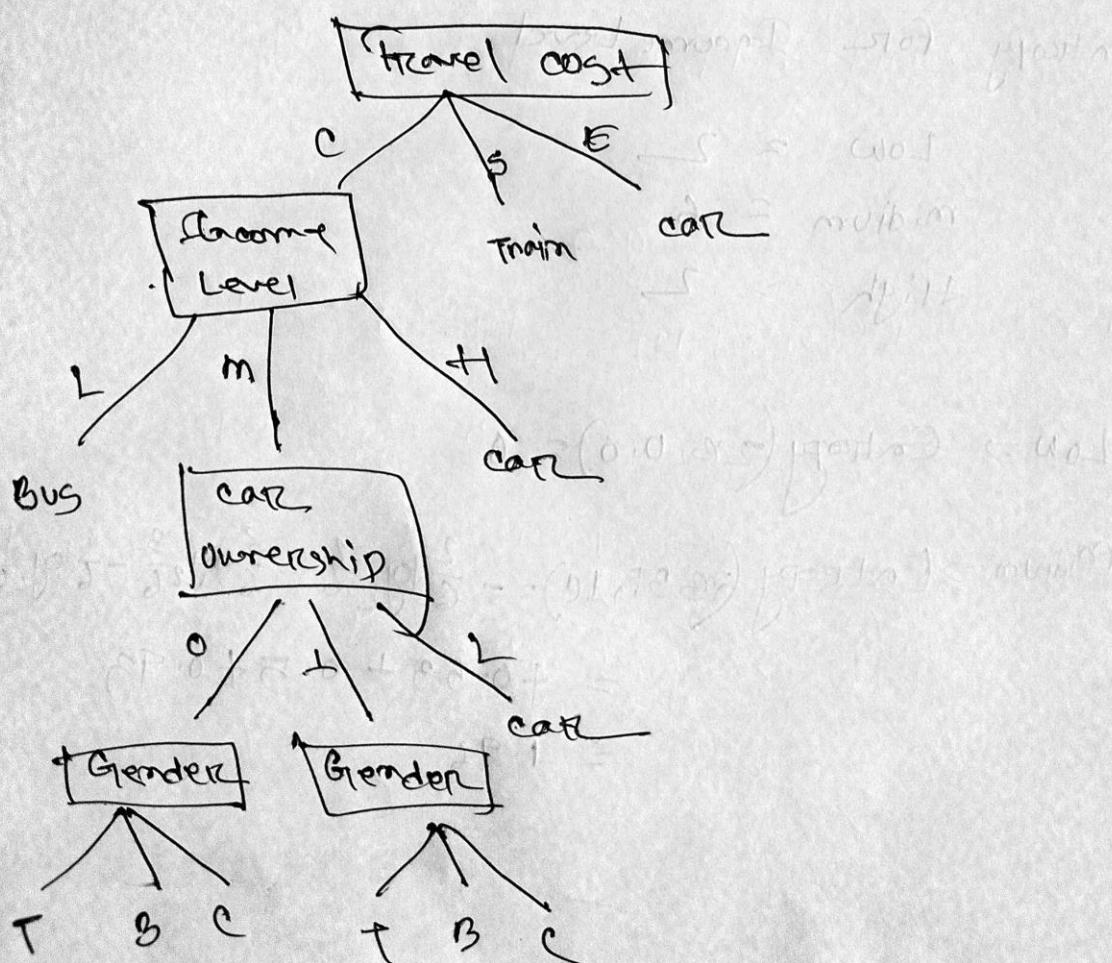
Low  $\Rightarrow$  Entropy  $(23, 0, 0) = 0$

$$\begin{aligned} \text{Medium} \Rightarrow \text{Entropy} &(23, 31, 10) = -\frac{2}{6} \log_2 \frac{2}{6} - \frac{3}{6} \log_2 \frac{3}{6} - \frac{1}{6} \log_2 \frac{1}{6} \\ &= +0.53 + 0.5 + 0.93 \\ &= 1.96 \end{aligned}$$

**Pantonix®**  
Pantoprazole

High  $\Rightarrow$  Entropy ( $20,0,0$ ) = 0

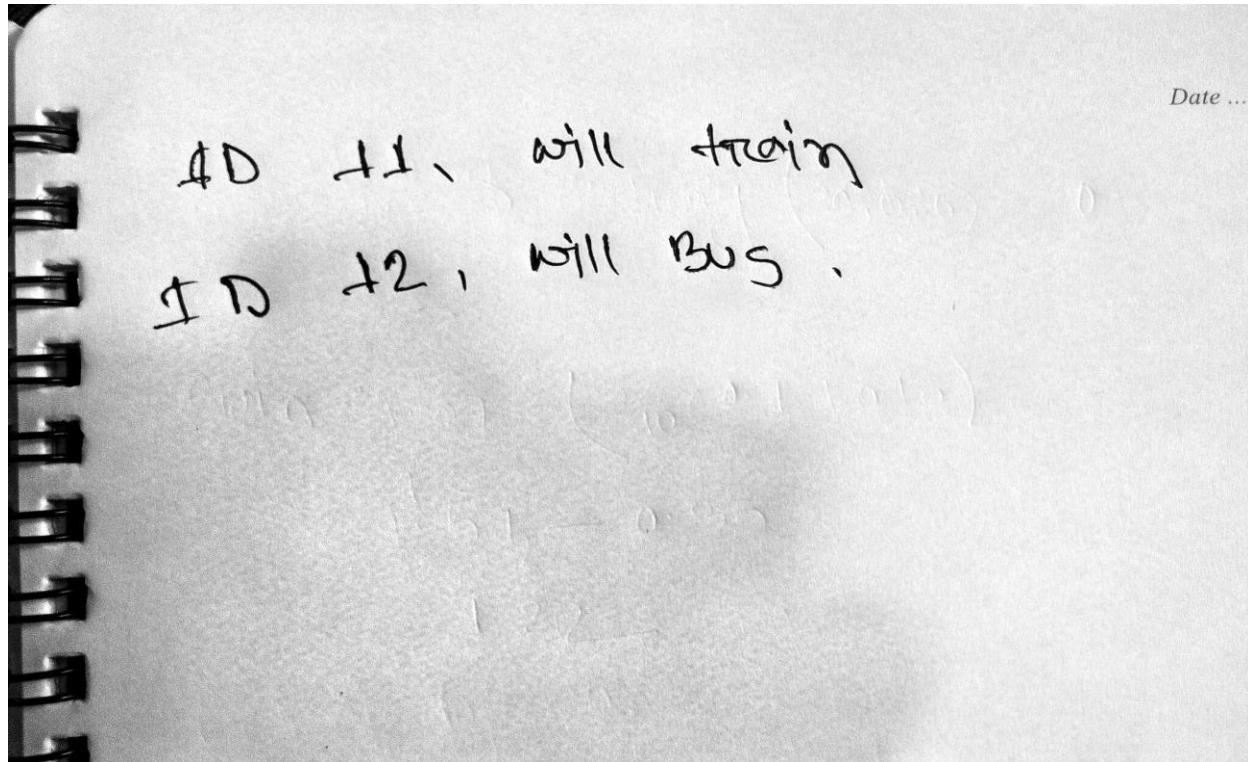
$$\text{Gini (for Income)} = 1 - \left( \frac{2}{10} \times 0 + \frac{6}{10} \times 0.46 + 0 \right) \\ = 1 - 0.276 \\ = 0.724$$



Date ....

AD 11, will train

ID 12, will bus.



Date .....

b) Naive Bayes algorithm work using the bayes theorem. This algorithm is very good for big data sets.

first of all we need to convert the data set into a frequency table. Then we need to create likelihood table by calculating the probability of attributes. For every attribute we need to calculate the probability based on the class attribute class.

Using the Bayes equation we need to calculate the posterior probability for each class. Hence the class with the highest

probability → is the outcome of prediction.

Here bayes equation for the classes,  
for id 11,

~~sol~~

$$\text{likelihood of Bus} = P(\text{Gender}=\text{Male} \mid \text{Bus}) * P(\text{CO}=1 \mid \text{Bus}) * \\ P(\text{TC}=\text{standard} \mid \text{Bus}) * P(\text{IC}=\text{medium}) * \\ P(\text{Bus})$$

$$\text{likelihood of Train} = P(\text{Gender}=\text{Male} \mid \text{Train}) * P(\text{CO}=1 \mid \text{Train}) * \\ P(\text{TC}=\text{standard} \mid \text{Train}) * P(\text{IC}=\text{medium} \mid \text{Train}) * \\ P(\text{Train})$$

$$\text{likelihood of Cart} = P(\text{Gender}=\text{Male} \mid \text{Cart}) * P(\text{CO}=1 \mid \text{Cart}) * \\ P(\text{TC}=\text{standard} \mid \text{Cart}) * P(\text{IC}=\text{medium} \mid \text{Cart}) * \\ P(\text{Cart})$$

Like id 11 we the equation similar to 10.