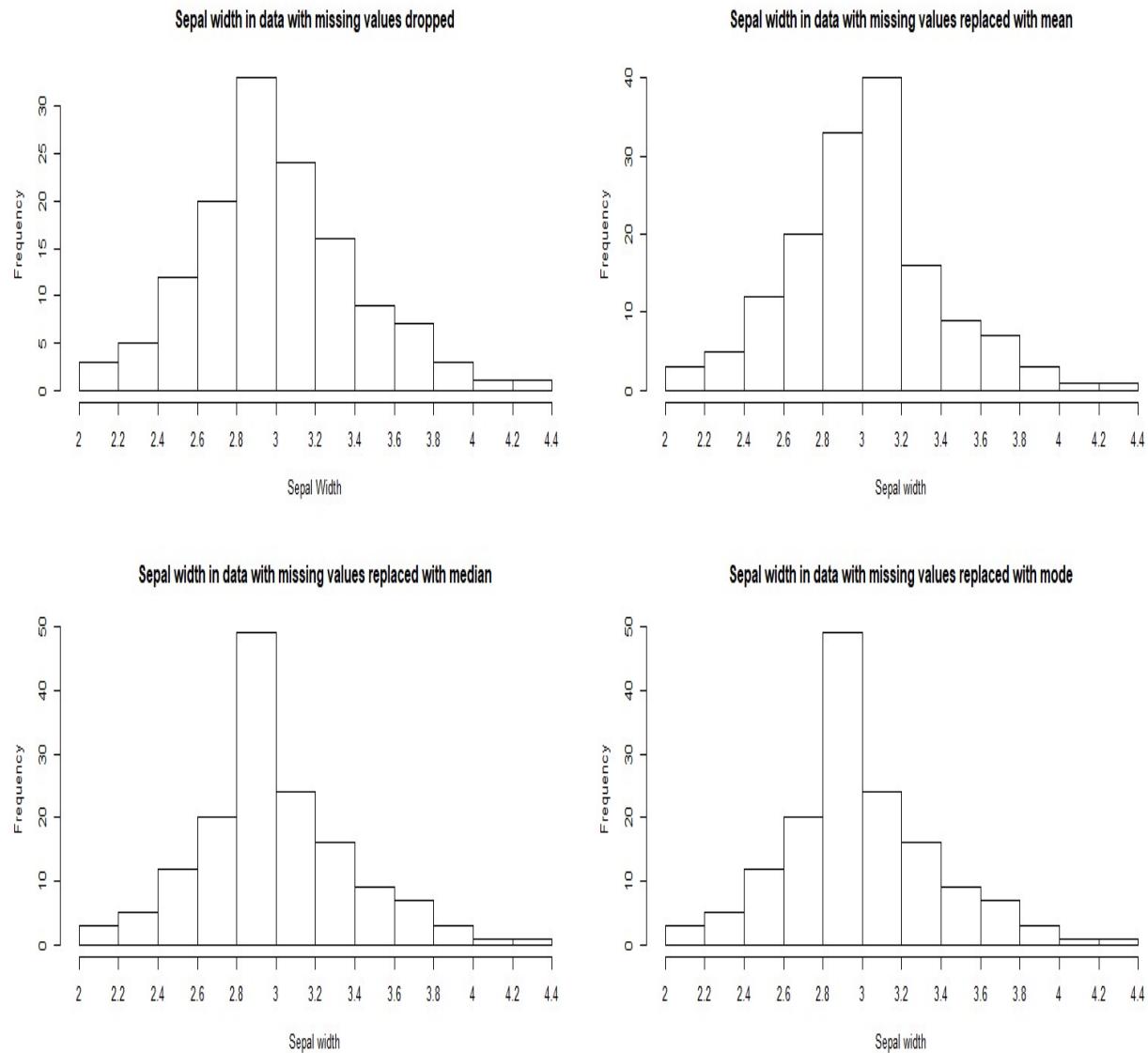


DATA MINING

Deepak kumar(1793606)

Part B Question 4: Histograms



Comment on Histograms:

It can be clearly observed from the set of histograms that :

- a) When the Sepal width rows with missing values is dropped and histogram is plotted then, its observed that values between range 2.8 to 3 have maximum frequency in the data i.e about 34, values between range 3 to 3.2 have second highest frequency i.e 23, and the values with the third highest frequency are

those between range 2.6 to 2.8. Values with lowest frequency are those with their value between 4 to 4.4 as their frequency is hardly 2.

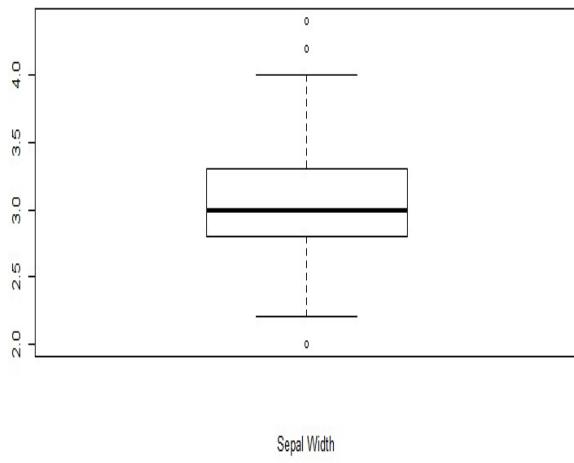
b) When the Sepal width rows with missing values is replaced with mean and a histogram is plotted then, it's observed that values between range 3 to 3.2 have maximum frequency in the data i.e about 40, values between range 2.8 to 3.0 have second highest frequency i.e 33, and the values with the third highest frequency(20) are those between range 2.6 to 2.8. Values with lowest frequency are those with their value between 4 to 4.4 as their frequency is hardly 2.

c) When the Sepal width rows with missing values is replaced with median and a histogram is plotted then, it's observed that values between range 2.8 to 3.0 have maximum frequency in the data i.e about 48, values between range 3 to 3.2 have second highest frequency i.e around 23, and the values with the third highest frequency(20) are those between range 2.6 to 2.8. Values with lowest frequency are those with their value between 4 to 4.4 as their frequency is just 1.

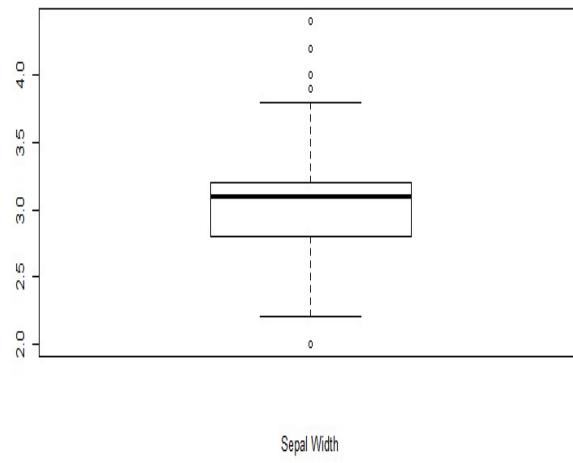
d) When the Sepal width rows with missing values is replaced with mode and a histogram is plotted then, it's observed that the plot is almost completely similar to that of histogram plotted with missing values is replaced with median.

Part B Question 4: Boxplot

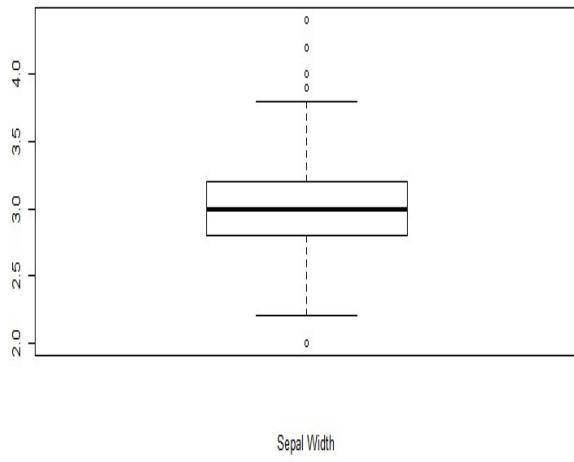
Sepal width in data with missing values dropped



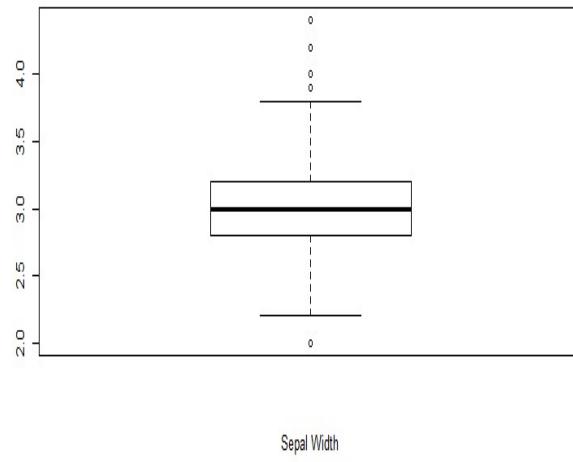
Sepal width in data with missing values replaced with mean



Sepal width in data with missing values replaced with median



Sepal width in data with missing values replaced with mode



Comment on boxplot:

It can be clearly observed from the set of boxplots that:

- When the Sepal width rows with missing values is dropped and boxplot is plotted then, its observed that 1st quartile here is around 2.8, 2nd quartile or median of the data is 3.0 and the 3rd quartile is around 3.3. The minimum value is around 2.2 and the maximum value is around 4.0.
- When the Sepal width rows with missing values is replaced with mean and a boxplot is plotted then, its observed that 1st quartile here is around 2.8 i.e. same as previous boxplot, 2nd quartile or median of the data is 3.1 and the 3rd quartile is around 3.2. The minimum value is around 2.2 and the maximum value is around 3.75.

c) When the Sepal width rows with missing values is replaced with median and a boxplot is plotted then, its observed that 1st quartile here is around 2.8, 2nd quartile or median of the data is 3.0 and the 3rd quartile is around 3.2. The minimum value is around 2.2 and the maximum value is around 3.75.

d) When the Sepal width rows with missing values is replaced with mode and a boxplot is plotted then, its observed that the boxplot is almost completely like that of boxplot plotted with missing values is replaced with median.

Part C question 2.

Class	Sex	Age	Survived	Freq
1 st	Male	Child	No	0
2 nd	Male	Child	No	0-
3 rd	Male	Child	No	35-
1 st	Female	Child	No	0
2 nd	Female	Child	No	0-
3 rd	Female	Child	No	17-
1 st	Male	Adult	No	118
2 nd	Male	Adult	No	154-
3 rd	Male	Adult	No	387-
1 st	Female	Adult	No	4
2 nd	Female	Adult	No	13-
3 rd	Female	Adult	No	85-
1 st	Male	Child	Yes	5
2 nd	Male	Child	Yes	11+
3 rd	Male	Child	Yes	13-
1 st	Female	Child	Yes	1
2 nd	Female	Child	Yes	13+
3 rd	Female	Child	Yes	14-
1 st	Male	Adult	Yes	57
2 nd	Male	Adult	Yes	14+
3 rd	Male	Adult	Yes	75-
1 st	Female	Adult	Yes	140
2 nd	Female	Adult	Yes	80+
3 rd	Female	Adult	Yes	76-

For the above given dataset we have survived as the Class attribute.

Now, as there are 12 yes for survived and 12 no as well

so there are 12 positive(P)yes and 12 negative(N)no's also.

Entropy for Class:

$$= -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

Information gain for each attribute:

$$I(P_i, N_i) = -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \left(\frac{N}{P+N} \right) \log_2 \left(\frac{N}{P+N} \right)$$

Entropy (Attribute):

$$E = \frac{\sum P_i + N_i}{P+N} (I(P_i, N_i))$$

Gain:

$$E(\text{Class}) - E(\text{Attribute})$$

Now, we will find the entropy for the class ie survived

$$\begin{aligned} E(\text{survived}) &= \frac{-12}{12+12} \log_2\left(\frac{12}{12+12}\right) - \frac{12}{12+12} \log_2\left(\frac{12}{12+12}\right) \\ &= -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \\ &= + \end{aligned}$$

E(Survived) given that from 1316 instances of data provided in titanic survival data is 817 didn't survive and 500 survived.

$$\begin{aligned} E(\text{survived}) &= \frac{203}{325} \log_2\left(\frac{203}{325}\right) - \frac{122}{325} \log_2\left(\frac{122}{325}\right) \\ &= (-0.625 \times -0.69) = 0.9574 \end{aligned}$$

Now, calculating the information gain for Class.

	P _i	N _i	I(P _i , N _i)
1 st	203	122	0.96
2 nd	118	167	0.97
3 rd	178	528	0.81

$$\begin{aligned} \text{For 1st Class} &= \frac{203}{325} \log_2\left(\frac{203}{325}\right) - \frac{122}{325} \log_2\left(\frac{122}{325}\right) \\ &= (-0.625 \times -0.69) (-0.375 \times -1.42) \\ &= 0.43 + 0.53 \\ &= 0.96 \quad (0.9574 \text{ accurate}) \end{aligned}$$

$$\begin{aligned} \text{2nd Class: } & -\frac{118}{285} \log_2\left(\frac{118}{285}\right) - \frac{167}{285} \log_2\left(\frac{167}{285}\right) \\ & = -0.41 \times -1.28 - 0.58 \times -0.28 \\ & = 0.9785 \end{aligned}$$

Information gain for 3rd class:

$$= \frac{178}{706} \log_2 \left(\frac{178}{706} \right) - \frac{528}{706} \log_2 \left(\frac{528}{706} \right)$$

$$= 0.8146$$

Now, if we calculate the entropy for Class:

$$E(\text{Class}) = \frac{203+122}{1316} \times 0.96 + \frac{118+167}{1316} \times 0.97$$

$$+ \frac{178+528}{1316} \times 0.81$$

$$= 0.8847$$

Information gain for Sex			
	P _i	N _i	I(P _i , N _i)
male	175	694	0.72
female	324	123	0.85

Information gain for male:

$$= \frac{175}{324} - \frac{175}{(175+694)} \log_2 \left(\frac{175}{(175+694)} \right) - \frac{324}{(324+123)} \log_2 \left(\frac{324}{324+123} \right)$$

$$= 0.7246$$

~~R(Female)~~ = Information gain for female:

$$= -\frac{324}{324+123} \times \log_2 \left(\frac{324}{324+123} \right) - \frac{123}{324+123} \log_2 \left(\frac{123}{447} \right)$$

$$= -0.72 \log_2 (0.72) - 0.27 \log_2 (0.28)$$

$$= 0.8482$$

Entropy for Sex is given by:

$$\frac{175+694}{1316} \times 0.72 + \frac{324+123}{1316} \times 0.85$$

$$= 0.7668$$

Considering the Age attribute

	P_i	N_i	$I(P_i, N_i)$
Child	57	52	0.9981
Adult	442	765	0.9477

Information gain for Child:

$$E(\text{Child}) = -\frac{57}{109} \log_2\left(\frac{57}{109}\right) - \frac{52}{109} \log_2\left(\frac{52}{109}\right)$$

$$= 0.9981$$

$$E(\text{Adult}) = -\frac{765}{1207} \log_2\left(\frac{765}{1207}\right) - \frac{442}{1207} \log_2\left(\frac{442}{1207}\right)$$

$$= 0.9477$$

Entropy of Age is given as:

$$\left(\frac{57+52}{1316}\right) \times (0.9981) + \frac{442+765}{1316} \times (0.95)$$

$$= 0.082 + 0.87$$

$$= 0.9519$$

Final Information gain for Class:

$$\rightarrow 0.9574 - 0.8847$$

$$= 0.0727$$

Final Information gain for Sex:

$$= 0.9574 - 0.7668$$

$$= 0.1906$$

Final Information gain for Age:

$$= 0.9574 - 0.9519$$

$$= 0.0055$$

As the information gain is maximum for Sex
so, it becomes the root node



Now we need to categorize the tree as per the male and female members:

Entropy (Female) as calculated before is 0.8487

Now, we calculate the entropy of class for female members:

	P_i	N_i	$I(P_i, N_i)$
1st	141	4	0.1821
2nd	93	13	0.5364
3rd	90	106	0.9951

Information gain

Entropy of 1st class:

$$-\frac{141}{145} \log_2 \left(\frac{141}{145} \right) - \frac{4}{145} \log_2 \left(\frac{4}{145} \right)$$

$$= 0.1821$$

Information gain for 2nd class:

$$-\frac{93}{106} \log_2 \left(\frac{93}{106} \right) - \frac{13}{106} \log_2 \left(\frac{13}{106} \right)$$

$$= 0.5364$$

Information gain for 3rd class:

$$-\frac{90}{196} \log_2 \left(\frac{90}{196} \right) - \frac{106}{196} \log_2 \left(\frac{106}{196} \right)$$

$$= -0.9951$$

for female sex

Entropy of the attribute Class is given as:

$$\text{EC(Class1)} = \left(\frac{145}{447} \right) \times (0.18) + \left(\frac{106}{447} \right) \times (0.53)$$

$$+ \left(\frac{196}{447} \right) \times (0.49)$$

$$= 0.6227$$

Now, we consider the attribute Age:

Information gain for adult is calculated as

$$\text{E(Adult)} = \frac{296}{402} \log_2\left(\frac{296}{402}\right) - \frac{106}{402} \log_2\left(\frac{106}{402}\right)$$

$$= 0.8322$$

Information gain for child is calculated as

$$\begin{aligned} \text{E(Child)} &= \frac{17}{45} \log_2\left(\frac{17}{45}\right) - \frac{28}{45} \log_2\left(\frac{28}{45}\right) \\ &= 0.9564 \end{aligned}$$

Entropy for the attribute age given sex is female is given as:

$$\begin{aligned} \text{E(Age)} &= \frac{402}{447} \times 0.83 + \frac{45}{447} \times 0.9564 \\ &= 0.8447 \end{aligned}$$

Final gain of ^{attribute} class given members are female:

$$0.8487 - 0.6227$$

$$= 0.2259$$

Final gain of ^{attribute} age given members are female:

$$0.8487 - 0.8444$$

$$= 0.0039$$

Now we consider calculation of gain when the sex is male.

Information gain calculated for male in part 1 was E(Male) 0.7247.

Information gain for Class 1st when sex is male

$$\begin{aligned} \text{E(1st)} &= \frac{62}{180} \log_2\left(\frac{62}{180}\right) - \frac{118}{180} \log_2\left(\frac{118}{180}\right) \\ &= 0.9280 \end{aligned}$$

Information gain for class 2nd when sex is female

$$= -\frac{25}{179} \log_2\left(\frac{25}{179}\right) - \frac{154}{179} \log_2\left(\frac{154}{179}\right)$$

$$= 0.5833$$

Information gain for class 3rd when sex is male:

$$= -\frac{88}{510} \log_2\left(\frac{88}{510}\right) - \frac{422}{510} \log_2\left(\frac{422}{510}\right)$$

$$= 0.6635$$

Entropy of the class when sex is male is given as:

$$= \frac{180}{869} \times 0.92 + \frac{179}{869} \times 0.58 + \frac{510}{869} \times 0.66$$

$$= 0.7019$$

Now we consider attribute age:

Information gain for Adult is given as (when sex is male)

~~$$E_{\text{Adult}} = -\frac{659}{805} \log_2$$~~

$$= -\frac{146}{805} \log_2\left(\frac{146}{805}\right) - \frac{659}{805} \log_2\left(\frac{659}{805}\right)$$

$$= 0.6830$$

Information gain for Children with male sex:

$$-\frac{35}{64} \log_2\left(\frac{35}{64}\right) - \frac{29}{64} \log_2\left(\frac{29}{64}\right)$$

$$= 0.9936$$

Entropy of attribute age when the sex is male is

$$= \frac{805}{869} \times 0.68 + \frac{64}{869} \times 0.99$$

$$= 0.7059$$

Total Information gain of Class when sex is male is

$$0.7947 - 0.7019$$

$$= 0.0928$$

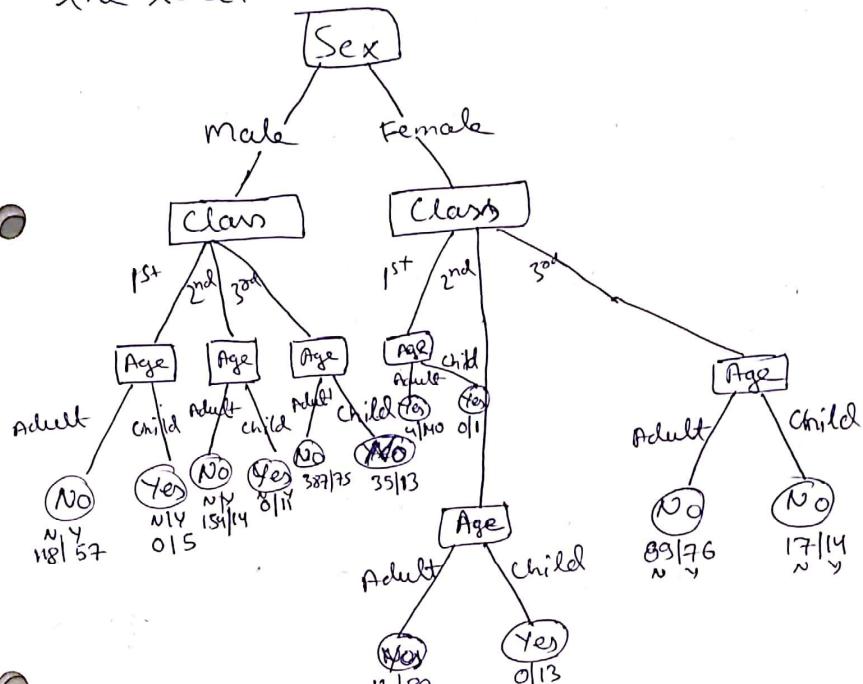
Total Information gain of Age when the sex is male is :

$$= 0.7247 - 0.7059$$

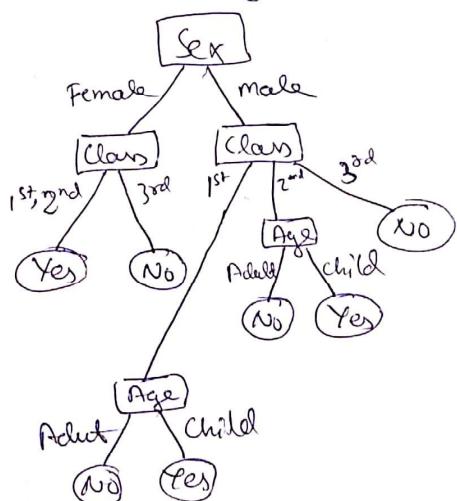
$$= 0.0187$$

Information gain of Age

So on the basis of gains found, when the Sex is male and when sex is female, we split the tree:



Above tree can be generalized as :



Rules obtained from the tree:

- a) IF Sex = Female \wedge (Class = 1st \vee Class = 2nd)
THEN Survived = Yes
- b) IF Sex = Male \wedge (Class = 1st \vee Class = 2nd)
 \wedge Age = Adult THEN Survived = No
- c) IF Sex = Male \wedge (Class = 1st \vee Class = 2nd)
 \wedge Age = Child THEN Survived = Yes
- d) IF (Sex = Male \vee Sex = Female) \wedge
Class = 3rd THEN Survived = No

PART D Question 1:

4(a)
Q

Question : 4(a)

Five draw tables for the deer
provided for Hair Eye color data

When Sex = Male

Hair	Brown	Blue	Hazel	Green
Black	32	11	10	3
Brown	53	50	25	15
Red	10	10	7	7
Blond	3	30	5	8

When Sex = Female

Hair	Brown	Blue	Hazel	Green
Black	36	9	5	2
Brown	66	34	29	14
Red	16	7	7	7
Blond	4	64	5	8

Drawing an individual table for Hair with Male and Female columns we have (We add males and females with different coloured hair.)

Hair	Male	Female
Black	56	52
Brown	143	143
Red	34	37
Blond	46	81
Total :	279	313

Similarly we make another table for Eye with Male and Female columns by adding the males and females with eyes of different colours:

Eye	Male	Female
Brown	98	122
Blue	101	114
Hazel	47	46
Green	33	31
Total :	279	313

Now we will use the secondary tables that we have created to do the probability computation

Prior probability:

Prior Probability of males =

= Total number of males

(Total number of males + Total Number of Females)

$$= \frac{279}{279+313} = 0.4712$$

Prior probability of Females =

Total number of females

Total males + Total females

$$= \frac{313}{279+313} = 0.5287$$

Now we calculate the likelihood function for hair:

$$P(\text{Male} | \text{Black}) = \frac{58}{279}$$

$$P(\text{Black} | \text{Male})_{\text{Hair}} = \frac{56}{279} = 0.20$$

$$P(\text{Black} | \text{Female})_{\text{Hair}} = \frac{52}{313} = 0.17$$

$$P(\text{Brown} | \text{Male})_{\text{Hair}} = \frac{143}{279} = 0.51$$

$$P(\text{Brown} | \text{Female})_{\text{Hair}} = \frac{143}{313} = 0.46$$

$$P(\text{Red} | \text{Male})_{\text{Hair}} = \frac{34}{279} = 0.12$$

$$P(\text{Red} | \text{Female})_{\text{Hair}} = \frac{37}{313} = 0.12$$

$$P(\text{Blond} | \text{Male})_{\text{Hair}} = \frac{46}{279} = 0.16$$

$$P(\text{Blond} | \text{Female})_{\text{Hair}} = \frac{81}{313} = 0.26$$

Further, we calculate the likelihood function for Eye color:

$$P(\text{Brown} \mid \text{Male}) = \frac{98}{\text{Eye} \quad 279} = 0.35$$

$$P(\text{Brown} \mid \text{Female}) = \frac{122}{\text{Eye} \quad 313} = 0.39$$

$$P(\text{Blue} \mid \text{Male}) = \frac{101}{\text{Eye} \quad 279} = 0.36$$

$$P(\text{Blue} \mid \text{Female}) = \frac{114}{\text{Eye} \quad 313} = 0.36$$

$$P(\text{Hazel} \mid \text{Male}) = \frac{47}{\text{Eye} \quad 279} = 0.17$$

$$P(\text{Hazel} \mid \text{Female}) = \frac{46}{\text{Eye} \quad 313} = 0.15$$

$$P(\text{Green} \mid \text{Male}) = \frac{33}{\text{Eye} \quad 279} = 0.12$$

$$P(\text{Green} \mid \text{Female}) = \frac{31}{\text{Eye} \quad 313} = 0.10$$

The calculated values

PART D Question 3:

Part D - Question 3

If the number of attribute instances is low or zero, the probabilities estimated from frequency of the attribute will yield bad results, ie about 0 probability.

Other probabilities in the model will not even matter if one probability is 0, as the end product obtained after multiplication of probabilities will be 0.

Usually laplacean correction is applied to fix this kind of issue. Since in Naive bayes products of probabilities of the features is evaluated during the training of the model and we clearly don't want it to evaluate to zero. So, to get rid of this, we have to assign some non-zero probabilities to words (features) which do not occur in the particular sample. This is what Laplace smoothing does.

3) If one of the features like if there were no bald headed men in the dataset, then the probability for that feature would be 0.

$$\text{i.e., } P(\text{Red Hair} | \text{Male}) = \frac{0}{279} = 0$$

Laplacian correction is a solution for this kind of issue.

Now if $P(\text{Red Hair} | \text{male}) = 0$

then in the Naive Bayes function we will check the attribute for which the probability is 0, and then loop through the attribute to add 1 to every instance. The outcome of this discussion is discussed on next page

Let us assume,
If the probability initially was

$$P(\text{Red hair | Male}) = \frac{0}{279} = 0$$

$$P(\text{Black hair | Male}) = \frac{56}{279} = 0.20$$

$$P(\text{Brown hair | Male}) = \frac{143}{279} = 0.51$$

$$P(\text{Blond hair | Male}) = \frac{46}{279} = 0.16$$

The after the Laplacean

correction the probability should
look like:

$$P(\text{Red hair | Male}) = \frac{1}{283}$$

$$P(\text{Black hair | Male}) = \frac{57}{283}$$

$$P(\text{Brown hair | Male}) = \frac{144}{283}$$

$$P(\text{Blond hair | Male}) = \frac{47}{283}$$