

Assignment 2

Section - A

- a) In a random forest, as the diversity increases, the correlation decreases and vice versa. This is the tradeoff between diversity and correlation in a random forest. If the correlation is too high and the diversity is too low, the variance will be high, thus causing overfitting. Similarly, if the trees are completely uncorrelated or they have very low correlation, the diversity will be high, however it will not be possible to form aggregates easily as all decision trees will act independently and hence we will not get a reliable output. That is why, It is important for the trees to be correlated to some extent in order to both minimize the variance and also reach a consensus among trees, to form reliable aggregate output, thus having both generalization and reliability.

- b) The curse of dimensionality can become an issue in Naive Bayes when we have data with a very large number of features and correspondingly not enough data points. This will result in overfitting, as the model is too complex and is not generalizable with the amount of data available.

To solve this problem, we can use dimensionality reduction. Techniques like Principal component analysis can be used to reduce the dimensionality of the data and choose those features that are most strongly correlated with the output and thus have the most effect on it.

- c) If the Naive bayes classifier encounters a sample which contains attribute values that is not present in the training dataset, it can lead to the zero probability problem. In other words, the classifier will assign it 0 probability of occurring, as it has not encountered this value earlier while training.

To resolve this problem, we can use Laplace smoothing. In laplace smoothing, using a parameter p (not equal to 0), we assign a non zero probability to all data points.

For example: let's say we have a feature X which predicts a target variable Y , X consists of values $\{0, 1, 2\}$, Y consists of values $\{0, 1\}$. We don't have any value such that, $\{Y = 0 \mid X = 0\}$. When we get such a sample without laplace smoothing, $P(Y = 1 \mid X = 0) = 0$.

However, if we use laplace smoothing with parameter k , then $P(Y = 1 \mid X = 0) = (\text{Count}(Y = 1, X = 0) + k) / (\text{Count}(Y = 1) + k * n)$, where n = number of classes. Thus, in this way we have resolved the zero probability problem.

- d) Yes, splitting a decision tree node using Information gain is biased if some attributes have more cardinality than others. This can lead to overfitting as an attribute with high cardinality will split the data further using several branches, and this will result in lack of generalizability.

To resolve this issue, we also need to take cardinality into account when splitting a node. One such criterion for that is gain ratio. Gain ratio mitigates the

effect of cardinality by also taking into consideration the information gained from a split along with its cardinality, thus giving a more balanced result.

For example, let's consider a case where we have two features, F1 with values {0, 1, 2} and F2 with values {3, 4, 5}, and an output variable O with values {0, 1}. In this case, the formulae for Information Gain and Gain Ratio can be expressed as:

$$\text{Gain}(F1) = \text{Entropy}(O) - \sum_{v \in F1} P(F1 = v) \cdot \text{Entropy}(O|F1 = v) \quad \text{GainRatio}(F1) = \text{Gain}(F1) / \text{IntrinsicInformation}(F1)$$

In these formulas, 'Entropy(O)' represents the entropy of the output variable 'O,' and 'P(F1 = v)' represents the probability that the attribute 'F1' takes on the value 'v.' The Gain Ratio effectively takes both the information content of 'F1' and its cardinality into account, ensuring a more balanced attribute selection process.

2)

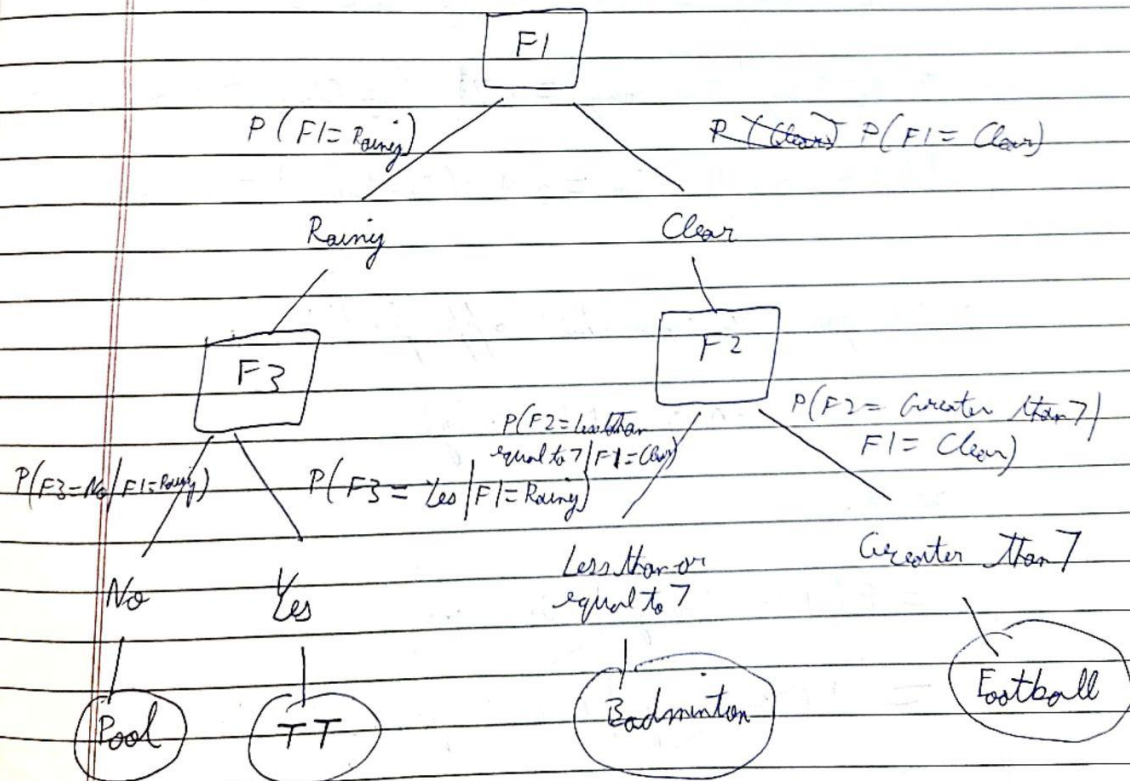
A2) (a) Initially, there are 3 features for classification

F1 \Rightarrow Weather $\begin{cases} \text{Raining} \\ \text{Clear} \end{cases}$

F2 \Rightarrow Outdoor activity friends $\begin{cases} \text{Greater than 7} \\ \text{Less than or equal to 7} \end{cases}$

F3 \Rightarrow Can borrow racket? $\begin{cases} \text{Yes} \\ \text{No} \end{cases}$

Based on these 3 features and the possibilities mentioned in the problem, ~~the~~ here is the decision tree.



Outcome = { Pool, TT, Badminton, Football }

$$P(\text{Pool}) = P(F_2 = \text{No} \mid F_1 = \text{Raining}) \cdot P(F_1 = \text{Raining})$$

$$P(\text{TT}) = P(F_2 = \text{Yes} \mid F_1 = \text{Raining}) \cdot P(F_1 = \text{Raining})$$

$$P(\text{Badminton}) = \frac{P(F_3 = \text{less than or equal to } 7 \mid F_1 = \text{Clear})}{P(F_1 = \text{Clear})}$$

$$P(\text{Football}) = \frac{P(F_3 = \text{Greater than } 7 \mid F_1 = \text{Clear})}{P(F_1 = \text{Clear})}$$

(b) let P_c be the app's prediction,

$$P(P_c = \text{Raining}) = 0.3, \quad P(P_c = \text{Clear}) = 0.7$$

$$P(P_c = \text{Raining} \mid F_1 = \text{Raining}) = 0.8, \quad P(P_c = \text{Clear} \mid F_1 = \text{Raining}) = 0.2$$

$$P(P_c = \text{Clear} \mid F_1 = \text{Clear}) = 0.9, \quad P(P_c = \text{Raining} \mid F_1 = \text{Clear}) = 0.1$$

~~$P(\text{It's going to rain but app predicts raining})$~~

~~$$\Rightarrow P(P_c = \text{Raining} \mid F_1 = \text{Raining})$$~~

~~$\Rightarrow 0$~~

~~$$= P(F_1 = \text{Raining} \mid P_c = \text{Raining})$$~~

~~$$= \frac{P(P_c = \text{Raining} \mid F_1 = \text{Raining}) \cdot P(F_1 = \text{Raining})}{P(P_c = \text{Raining})}$$~~

~~$$= \frac{0.8 \times P(F_1 = \text{Raining})}{0.3}$$~~

$$P(FI = \text{Raining}) =$$

P (It is going to rain but app predicts raining)

$$= P(BI = \text{Raining} | FI = \text{Raining})$$

$$= P(FI = \text{Raining} | BI = \text{Raining})$$

$$= \frac{P(BI = \text{Raining} | FI = \text{Raining}) \cdot P(FI = \text{Raining})}{P(BI = \text{Raining})}$$

$$= \frac{0.8 \times P(FI = \text{Raining})}{0.3}$$

Assumption: If $P(FI = \text{Raining}) = P(BI = \text{Raining}) = 0.3$,
Then

$$P(FI = \text{raining} | BI = \text{raining}) = 0.3$$

$$(Q2) \quad P(\text{Raining}) =$$

$$P(B_1 = \text{Raining}) = P(F_1 = \text{Raining}) \cdot P(B_1 = \text{Raining} \mid F_1 = \text{Raining}) \\ + P(F_1 = \text{Clear}) \cdot P(B_1 = \text{Raining} \mid F_1 = \text{Clear})$$

$$0.3 = P(F_1 = \text{Raining}) \cdot (0.8) + (1 - P(F_1 = \text{Raining})) \cdot (0.2)$$

$$P(F_1 = \text{Raining}) = \frac{0.2}{0.7} = \frac{2}{7}$$

$$P(F_1 = \text{Raining} \mid B_1 = \text{Raining})$$

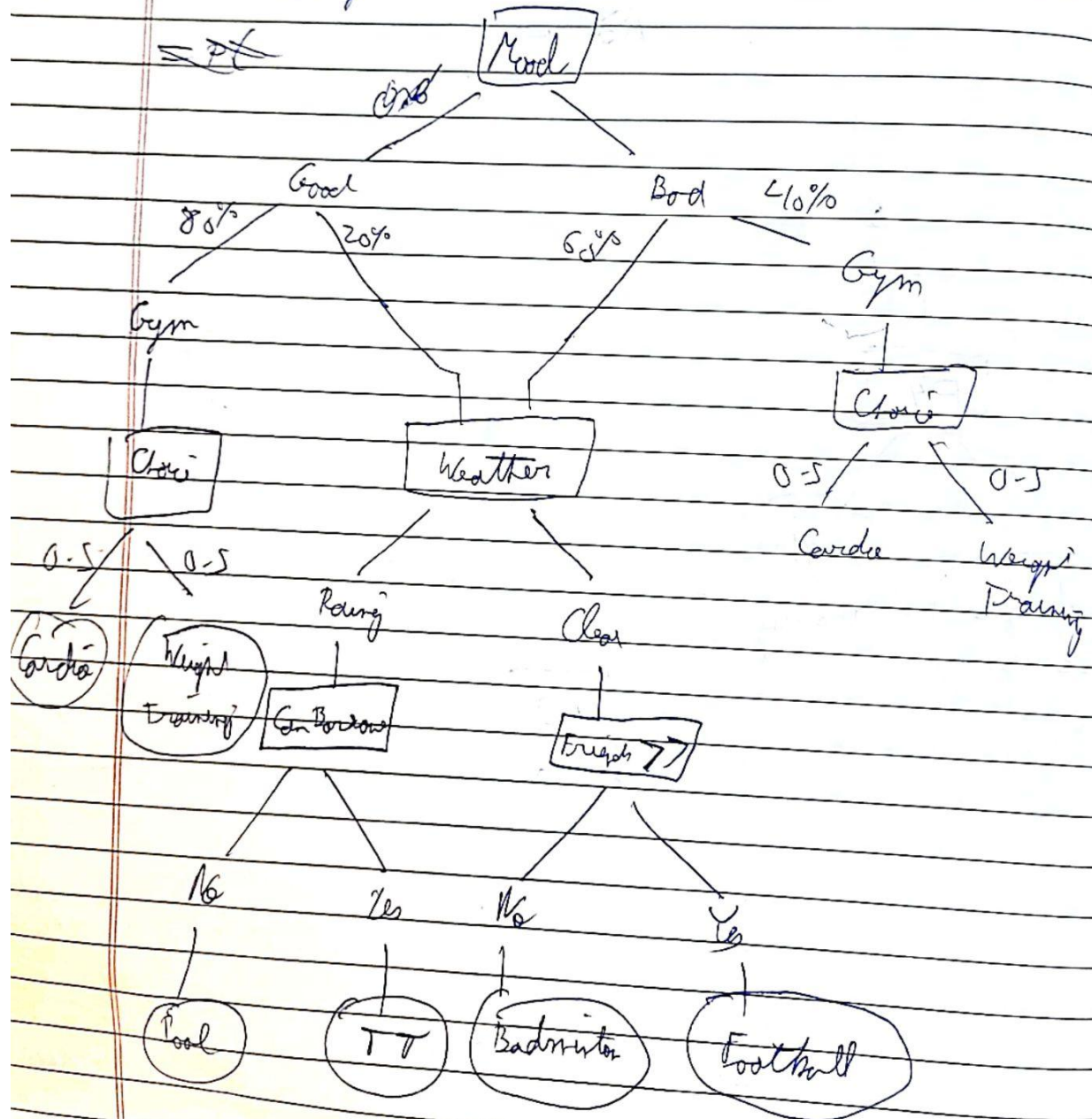
$$= \frac{P(B_1 = \text{Raining} \mid F_1 = \text{Raining}) \cdot P(F_1 = \text{Raining})}{P(B_1 = \text{Raining})}$$

$$= \frac{0.8 \times \frac{2}{7}}{\frac{2}{7}} = 0.8$$

$$= \frac{16}{21} = 0.762$$

Outcome: (Cardiological Exercise, Weight Training, Pool, T.V., Badminton, Football)

P (Cardiological Exercise)



$$P(\text{Cardiological Exercise})$$

$$= P(\text{Cardiological Exercise} \mid \text{Gym}) \cdot P(\text{Gym} \mid \text{Mood} = \text{Good})$$

$$+ P(\text{Cardiological Exercise} \mid \text{Gym})$$

$$= 0.5 \cdot P(\text{Gym} \mid \text{Mood} = \text{Bad}) \cdot P(\text{Mood} = \text{Bad})$$

$$= (0.5 \times 0.8 \times 0.4) + ($$

$$P(\text{Weight Training})$$

$$= P(\text{Weight Training} \mid \text{Gym}) \cdot P(\text{Gym} \mid \text{Good Mood})$$

$$+ P(\text{Weight Training} \mid \text{Gym}) \cdot P(\text{Gym} \mid \text{Bad Mood})$$

$$P(\text{Bad Mood})$$

$$P(\text{Pool}) = P(\text{Can Borrow} \mid \text{Raining}) \cdot P(\text{Raining} \mid \text{Good Mood}) \cdot P(\text{Good Mood})$$

$$P(\text{Pool}) = P(\text{Can Borrow} \mid \text{Raining}) \cdot P(\text{Raining} \mid \text{Good Mood}) \cdot P(\text{Good Mood})$$

$$P(\text{TT}) = P(\text{Can Borrow} \mid \text{Raining}) \cdot P(\text{Raining} \mid \text{Good Mood}) \cdot P(\text{Good Mood})$$

$$P(\text{Fail}) = \frac{P(\text{Cannot Borrow} | \text{Rainy}) \cdot P(\text{Rainy} | \text{Good Mood})}{P(\text{Good Mood})}$$

$$- \frac{P(\text{Cannot Borrow} | \text{Rainy}) \cdot P(\text{Rainy} | \text{Bad Mood})}{P(\text{Bad Mood})}$$

$$P(\text{Badminton}) = P(\text{Fail})$$

$$P(\text{TT}) = \frac{P(\text{Can Borrow} | \text{Rainy}) \cdot P(\text{Rainy} | \text{Good Mood}) \cdot P(\text{Good Mood})}{P(\text{Good Mood})}$$

$$+ \frac{P(\text{Can Borrow} | \text{Rainy}) \cdot P(\text{Rainy} | \text{Bad Mood}) \cdot P(\text{Bad Mood})}{P(\text{Bad Mood})}$$

$$P(\text{Badminton}) = \frac{P(\text{Friends} \leq 7 | \text{Clear}) \cdot P(\text{Clear} | \text{Good Mood}) \cdot P(\text{Good Mood})}{P(\text{Good Mood})}$$

$$+ \frac{P(\text{Friends} \leq 7 | \text{Clear}) \cdot P(\text{Clear} | \text{Bad Mood}) \cdot P(\text{Bad Mood})}{P(\text{Bad Mood})}$$

$$P(\text{Football}) = \frac{P(\text{Friends} > 7 | \text{Clear}) \cdot P(\text{Clear} | \text{Good Mood}) \cdot P(\text{Good Mood})}{P(\text{Good Mood})}$$

$$+ \frac{P(\text{Friends} > 7 | \text{Clear}) \cdot P(\text{Clear} | \text{Bad Mood}) \cdot P(\text{Bad Mood})}{P(\text{Bad Mood})}$$

$$0.42 \div 0.17 = 2.47$$

$$P(\text{Good Mood} | F=7) = \frac{0.42}{0.6} = 0.7$$

$$P(\text{Bad Mood} | F=7) = \frac{0.18}{0.6} = 0.3$$

$$P(\text{Gym}) = P(\text{Gym} | F=7) = P(\text{Good Mood} | F=7) \cdot P(\text{Gym} | \text{Good Mood}, F=7)$$

$$+ P(\text{Bad Mood} | F=7) \cdot P(\text{Gym} | \text{Bad Mood}, F=7)$$

$$= 0.7 \times 0.8 + 0.3 \times 0.4$$

$$= 0.68$$

$$P(\text{Not Gym} | F=7) = P(\text{Good Mood} | F=7) \cdot P(\text{Not Gym} | \text{Good Mood}, F=7)$$

$$+ P(\text{Bad Mood} | F=7) \cdot P(\text{Not Gym} | \text{Bad Mood}, F=7)$$

$$= 0.7 \times 0.2 + 0.3 \times 0.6$$

$$= 0.32$$

Therefore gym is the most likely option,
as $P(\text{Gym} | F=7) > P(\text{Not Gym} | F=7)$

Also, as weight training and cardiovascular
sources are equally probable, they
are both equally likely outcomes

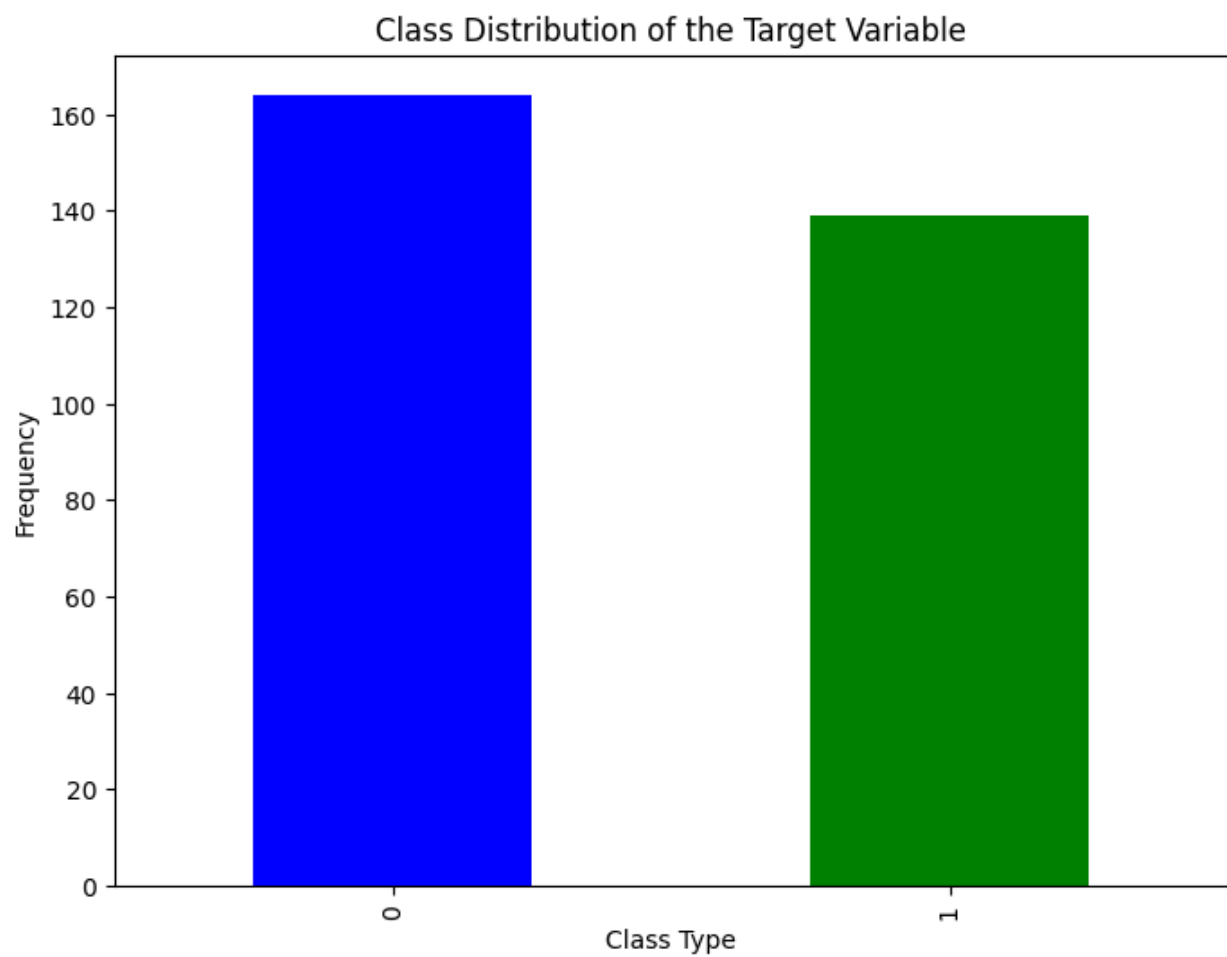
$$\text{Also, } P(\text{Cardio}) = 0.68 \times \frac{1}{2} = 0.34$$

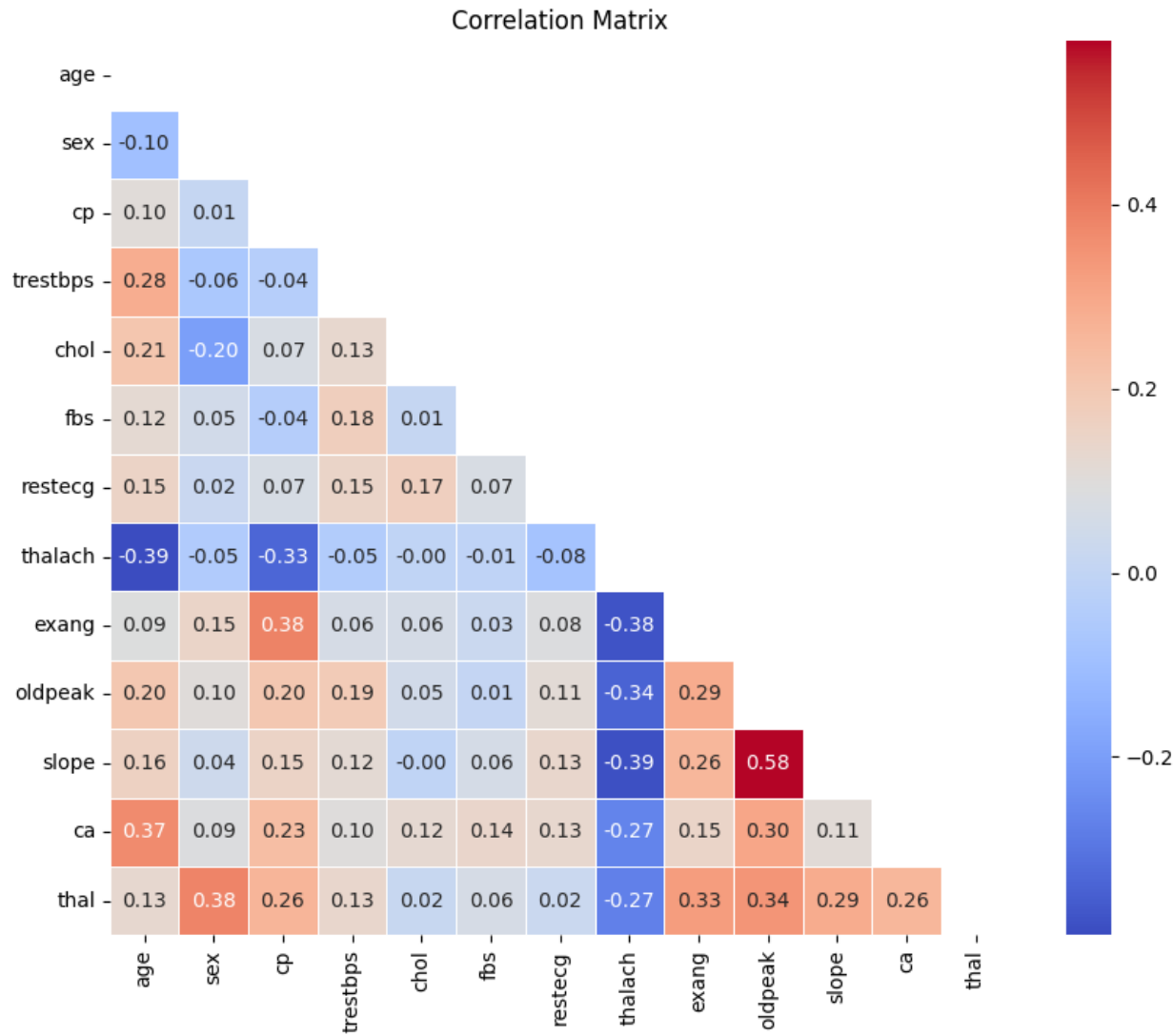
$$P(\text{Weight Training}) = 0.68 \times \frac{1}{2} = 0.34$$

$$P(\text{Total}) =$$

Section - C

- 1) In EDA, we found that, Features ca and thal have missing values. We used the median value to replace the missing values.
- 2) The class imbalance is very low as the percentage difference between classes 0 and 1 is only 16.5%. Therefore, we can carry on without resampling.
- 3) The data seems correlated. The feature thalach seems most strongly correlated with all the other features.
- 4) During Training and Testing, the following were our conclusions
Entropy (77%) gives a higher accuracy than Gini(72%) , therefore moving forward we will be using it as our criterion.
- 5) For decision tree, after hypertuning we got accuracy of 80%.
- 6) Random Forrest, gave an accuracy of upwards of 82%, slightly higher than decision tree.





Final Classification Report on test data:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.78 | 0.88 | 0.82 | 32 |
| 1 | 0.84 | 0.72 | 0.78 | 29 |
| accuracy | | | 0.80 | 61 |
| macro avg | 0.81 | 0.80 | 0.80 | 61 |
| weighted avg | 0.81 | 0.80 | 0.80 | 61 |