

Shourya Jindal
2020336
ML Assignment 2

Section A

(Q1)

(RF)

(a) In random forests, correlation refers to the degree to which individual decision trees, in the RF are similar to each other.

High Correlation = Trees make similar error and decision on same data pts

Low " = Make diff. decisions on same data pts.

Diversity refers to differences among the trees in the RF.

↳ diff. feature selections, thresholds, rules.

High diversity \Rightarrow Reduces overfitting, reduces variance.

So, tradeoff b/w these 2 is that

- if we increase correlation too much, it is similar to having single decision tree. leads to overfitting \Rightarrow ~~High Bias~~ Low Bias
High Variance
does not perform well on unseen data
- ~~But~~ And, if increase diversity too much, it leads to lack of consensus in the predictions. Reduces overall accuracy. It reduces overfitting, thus reducing variance.

(b) In 'naive Bayes', the "curse of dimensionality" becomes an issue, when, 'no. of features' i.e. the dimensionality of data becomes very large. It leads to very high computational ~~complex~~ complexity, overfitting. In higher dimensions data becomes sparse making it difficult to calculate probabilities in 'Naive Bayes' also leading to increased computational complexity. Also, leads to overfitting, thus poor performance on unseen data.

Strategies to Mitigate this:

- Feature Selection: Use only relevant features and discard redundant ".
- Dimensionality Red: Use methods like PCA, TSNE to reduce dimensionality
- Regularization: Use L_1, L_2 methods to prevent overfitting
- ensemble Methods: Use ensemble methods like boosting, random forest \Rightarrow More effective in higher dimensions.

(c) yes, if some value of attributes which was not present in the training set is encountered, it will affect the inference results.

Problems:

- when unseen value of attri. is encountered, it will assign it prob. 0.
- This can lead to wrong predictions
- Thus leading to high variance i.e. low accuracy on unseen data.
- Accuracy reduces

Soln: • Laplace Smoothing: Add a small value (generally 1, to count of each attribute value pair. \Rightarrow No zero prob.

- Use prior prob: If we have access to prior prob. for attri values, these can be given to the model to help estimate unseen values prob.

20

eg: Email ~~classification~~ classification $\begin{cases} \text{Spam} \\ \text{not spam} \end{cases}$

Training:

Email	to Lottery	free	Spam / Not
1	yes	yes	yes
2	yes	No	yes
3	No	yes	No
4	No	No	No

Test sample: emails: "Lottery free hello". = x^{new}

In this, $P(\text{"hello"} | \text{"spam"}) = 0$

$\Rightarrow P(\text{Spam} | x^{\text{new}}) = 0$ although it contains both "lottery" and "free".

but since "hello" is

(\Rightarrow) unseen, it is not declared as spam.

but if Laplace Smoothing will be used it will prob > 0 .

So, prob. are:
possible outcomes and prob. Using formulas
from last part (c)

* Cardio:

$$P[\text{Cardio}] = (0.42)(0.8)(0.5) + (0.18)(0.4)(0.5) \\ = 0.168 + 0.036 = 0.204$$

* weights:

$$P[\text{weights}] = \text{same as cardio} \\ = 0.204$$

* Plan:

$$P[\text{Plan}] = (0.42)(0.2) + (0.18)(0.6) \\ = 0.084 + 0.108 \\ = 0.192$$

So, To calculate exactly which Sport he will play in
info. is not available.

Most likely he will do:

So, He will visit gym and would either do

Cardio or weights with equal probability.

$$= \boxed{0.204}$$

Now, we have to find

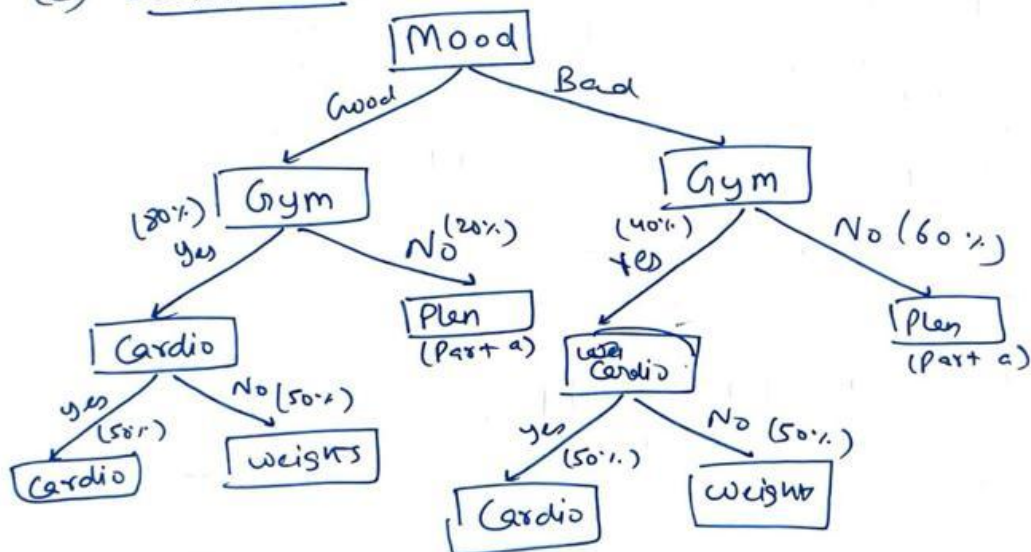
$$P[\text{Rain} | \text{Predicts Rainy}] = ??$$

$$= \frac{P[\text{Predicts Rainy} | \text{Rain}] \times P[\text{Rain}]}{P[\text{Predicts Rainy}]}$$

(Using given values and eqn ②)

$$= \frac{0.8 \times \frac{2}{3}}{0.3} = \frac{1.6}{2.1} = \boxed{0.76} \rightarrow \text{Ans}$$

(C) Decision Tree:



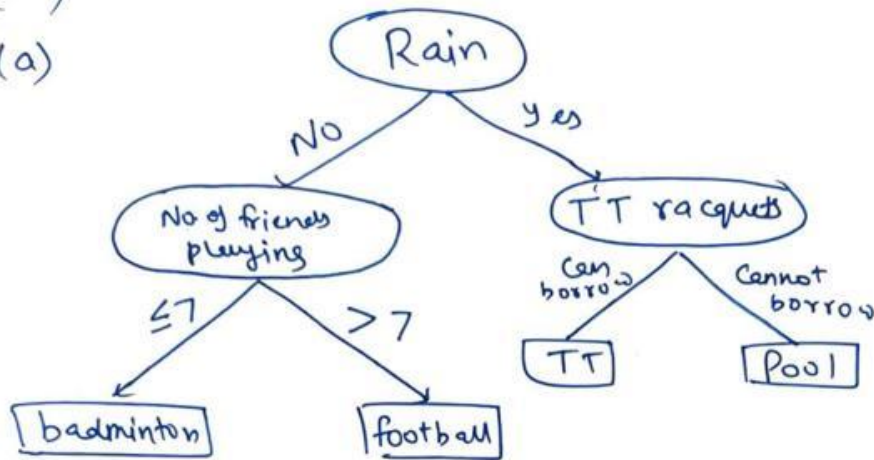
Possible Outcomes and prob.:

$$* P[\text{Cardio}] = P(\text{Good Mood}) \cdot \underbrace{P(\text{Gym} | \text{Good Mood})}_{= 0.8} \cdot \underbrace{P(\text{Cardio} | \text{Gym})}_{= 0.5} \\ + P(\text{Bad Mood}) \cdot \underbrace{P(\text{Gym} | \text{Bad Mood})}_{= 0.4} \cdot \underbrace{P(\text{Cardio} | \text{Gym})}_{= 0.5}$$

$$* P[\text{weights}] = P(\text{Good Mood}) \cdot \underbrace{P(\text{Gym} | \text{Good Mood})}_{= 0.8} \cdot \underbrace{P(\text{weights} | \text{Gym})}_{= 0.5} \\ + P(\text{Bad Mood}) \cdot \underbrace{P(\text{Gym} | \text{Bad Mood})}_{= 0.4} \cdot \underbrace{P(\text{weights} | \text{Gym})}_{= 0.5}$$

(Q2)

(a)



Possible Outcomes & their Probabilities:

* Badminton: $P(\text{Badminton}) = P(\text{No Rain}) * P(\text{No. of friends playing} \leq 7 | \text{No Rain})$

* Football: $P(\text{Football}) = P(\text{No Rain}) * P(\text{No. of friends playing} > 7 | \text{No Rain})$

* TT: $P(\text{TT}) = P(\text{Rain}) * P(\text{Can borrow racquets} | \text{Rain})$

* Pool: $P(\text{Pool}) = P(\text{Rain}) * P(\text{Cannot borrow racquets} | \text{Rain})$

(b) $P(\text{Predicts Rainy}) = 0.3$

$$P(\text{Predicts Clean}) = 0.7$$

$$P(\text{Predicts Rainy} | \text{Rain}) = 0.8$$

$$P(\text{Predicts Clean} | \text{Clean}) = 0.9$$

$$\hookrightarrow P(\text{Predicts Rain} | \text{Clean}) = 0.1 \quad \text{--- (1)}$$

Now, first we will calculate $P(\text{Rain})$

for this we will use total probability theorem:

$$P[\text{Predicts Rainy}] = P[\text{Rain}] * P[\text{Predicts Rainy} | \text{Rain}] + P[\text{Clean}] * P[\text{Predicts Rainy} | \text{Clean}]$$

(Using given a and eqn (1))

$$\Rightarrow 0.3 = 0.8 P[\text{Rain}] + 0.1 P[\text{Clean}]$$

$$\Rightarrow 0.3 = 0.8 P[\text{Rain}] + 0.1 (1 - P[\text{Rain}])$$

$$\Rightarrow \boxed{P(\text{Rain}) = \frac{2}{7}} \quad \text{--- (2)}$$

Now, we have to find

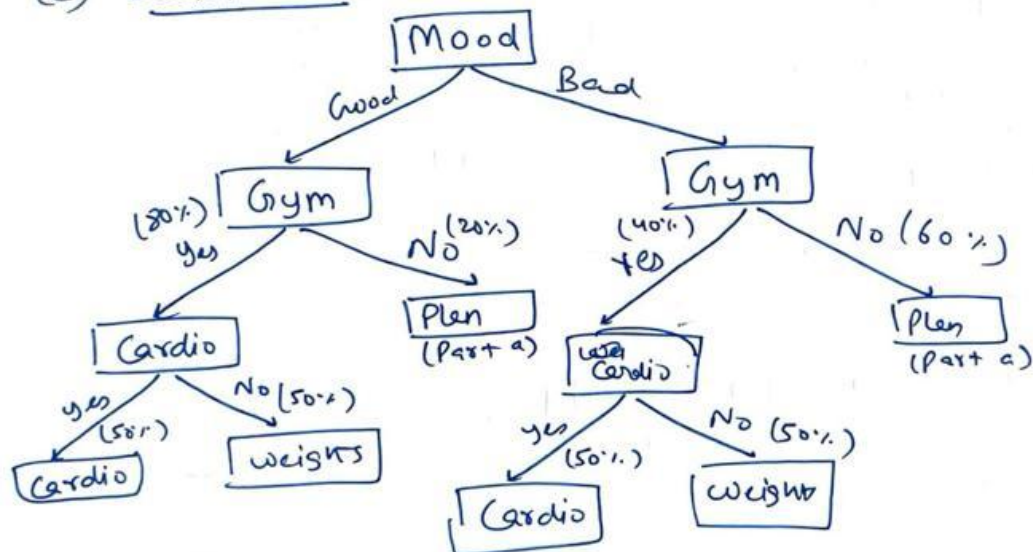
$$P[\text{Rain} | \text{Predicts Rainy}] = ??$$

$$= \frac{P[\text{Predicts Rainy} | \text{Rain}] \times P[\text{Rain}]}{P[\text{Predicts Rainy}]}$$

(Using given values and eqn (2))

$$= \frac{0.8 \times \frac{2}{3}}{0.3} = \frac{1.6}{2.1} = \boxed{0.76} \rightarrow \text{Ans}$$

(C) Decision Tree:



Possible Outcomes and prob.:

$$* P[\text{Cardio}] = P(\text{Good Mood}) \cdot \underbrace{P(\text{Gym} | \text{Good Mood})}_{= 0.8} \cdot \underbrace{P(\text{Cardio} | \text{Gym})}_{= 0.5} \\ + P(\text{Bad Mood}) \cdot \underbrace{P(\text{Gym} | \text{Bad Mood})}_{= 0.4} \cdot \underbrace{P(\text{Cardio} | \text{Gym})}_{= 0.5}$$

$$* P[\text{Weights}] = P(\text{Good Mood}) \cdot \underbrace{P(\text{Gym} | \text{Good Mood})}_{= 0.8} \cdot \underbrace{P(\text{Weights} | \text{Gym})}_{= 0.5} \\ + P(\text{Bad Mood}) \cdot \underbrace{P(\text{Gym} | \text{Bad Mood})}_{= 0.4} \cdot \underbrace{P(\text{Weights} | \text{Gym})}_{= 0.5}$$

(d) Yes, if some attri have more cardinality, there can be bias in splitting a decision tree node using Info. Gain.

This is because cardinality means that no. of distinct values an attribute can take.

If an attribute has high cardinality \Rightarrow

\Rightarrow they can be split into more subset (more branchy in a tree)

\Rightarrow low impurity

\Rightarrow More information Gain

\Rightarrow Thus, biasedness.

* Other criterion that can be used:

- Gini Index: $Gini(s) = 1 - \sum p_i^2$
it is less sensitive to cardinality

- Gain Ratio: It takes into cardinality and penalizes attributes with high ".

Example:

Suppose, we are building a decision tree to predict, ~~which type of ice-cream is sold~~ ice-Cream Sales.
~~Attributes Considered~~

Consider 3 attributes:

① Age Group: ~~low~~ $\{ <18, 18-35, 36-60, >60 \}$

\Rightarrow low cardinality

② Weather conditions: $\{ \text{Rainy, hot, cold} \}$

\Rightarrow low cardinality

③ Flavour = $\{ \text{Chocolate, butterscotch, berry, ...} \}$

\Rightarrow high cardinality

If we use IG as criterion, it can favour 'flavour' attribute (due to its high cardinality).

But, this may not be ~~low~~ optimal approach, as sales of ice-cream are generally affected most by the factors like 'weather conditions'.

* Rest of the outcomes are same
just use previous formula ~~and multiply them with~~ ~~also~~
following:

$$\underbrace{P(\text{Event})}_{\text{new prob}} = P(\text{Plan}) \cdot \underbrace{P(\text{Event})}_{\text{older prob}}$$

$$\text{where } P(\text{Plan}) = P(\text{Good Mood}) \cdot \overbrace{P(\text{No Gym} | \text{Good Mood})}^{= 0.2} \\ + P(\text{Bad Mood}) \cdot \underbrace{P(\text{No Gym} | \text{Bad Mood})}_{= 0.6}$$

$$(d) \quad P(\text{Good Mood}) = 0.6$$

$$P(\text{Bad Mood}) = 0.4$$

$$P(F=7 | \text{Gm}) = 0.7$$

$$P(F=7 | \text{Bm}) = 0.45$$

Assuming he slept 7 hrs
one night before

$$\therefore P(F=7) = 1$$

~~We are assuming that he slept $P(\text{Gm}) = 0.6$ & $P(\text{Bm}) = 0.4$
are given of one day earlier.~~

So/ On each last night

$$P(F=7) = P(\text{Gm}) \cdot P(F=7 | \text{Gm}) + P(\text{Bm}) \cdot P(F=7 | \text{Bm}) \\ = (0.6)(0.7) + (0.4)(0.45)$$

$$P(F=7) = 0.6$$

$$\text{Now, for next day: } P(\text{Gm} | F=7) = \frac{P(F=7 | \text{Gm}) \cdot P(\text{Gm})}{P(F=7)}$$

$$= \frac{(0.7)(0.6)}{0.6} = 0.7$$

$$P(\text{Bm} | F=7) = \frac{P(F=7 | \text{Bm}) \cdot P(\text{Bm})}{P(F=7)}$$

$$= \frac{(0.45)(0.4)}{0.6} = 0.3$$

Section B

PART A: Preprocessing and EDA

Read the data and store it in the pandas' data frame.

In total, 13 features and 1 target column.

Total rows: 303.

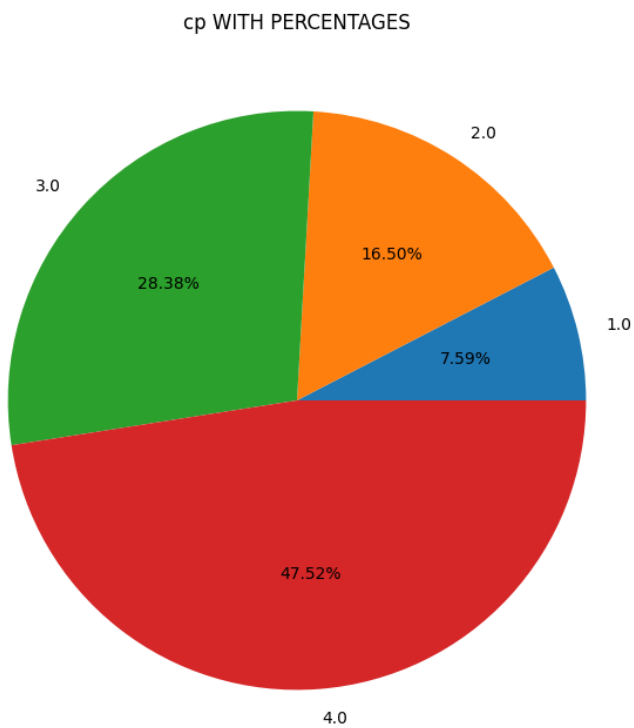
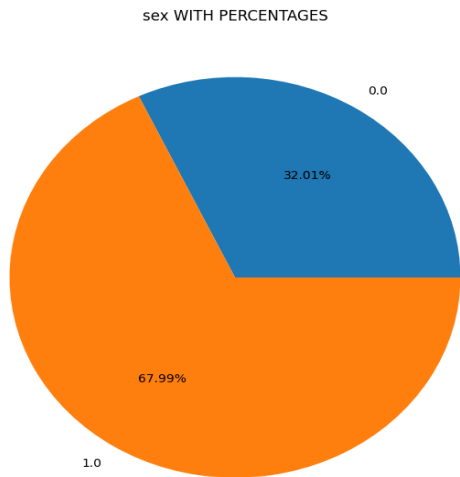
Preprocessing:

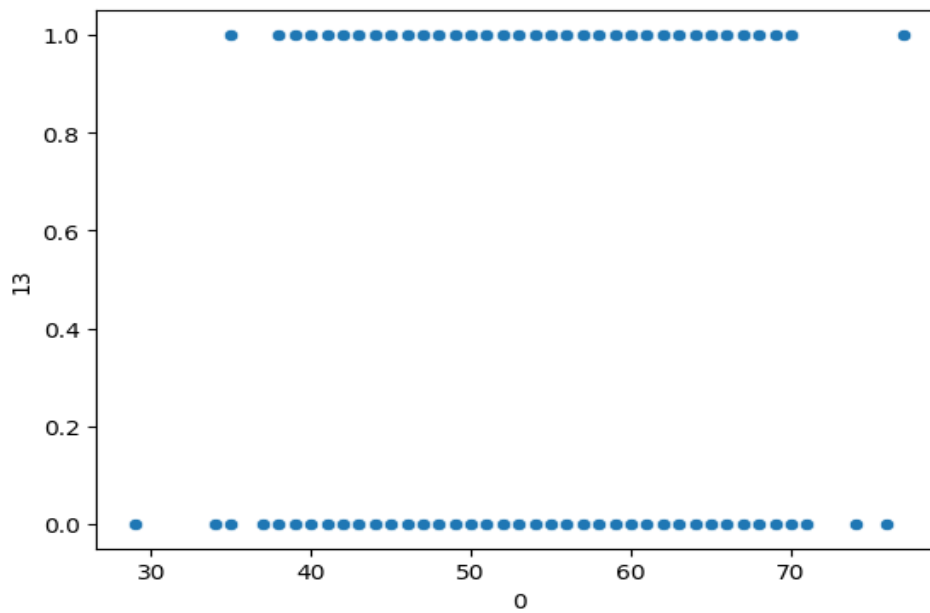
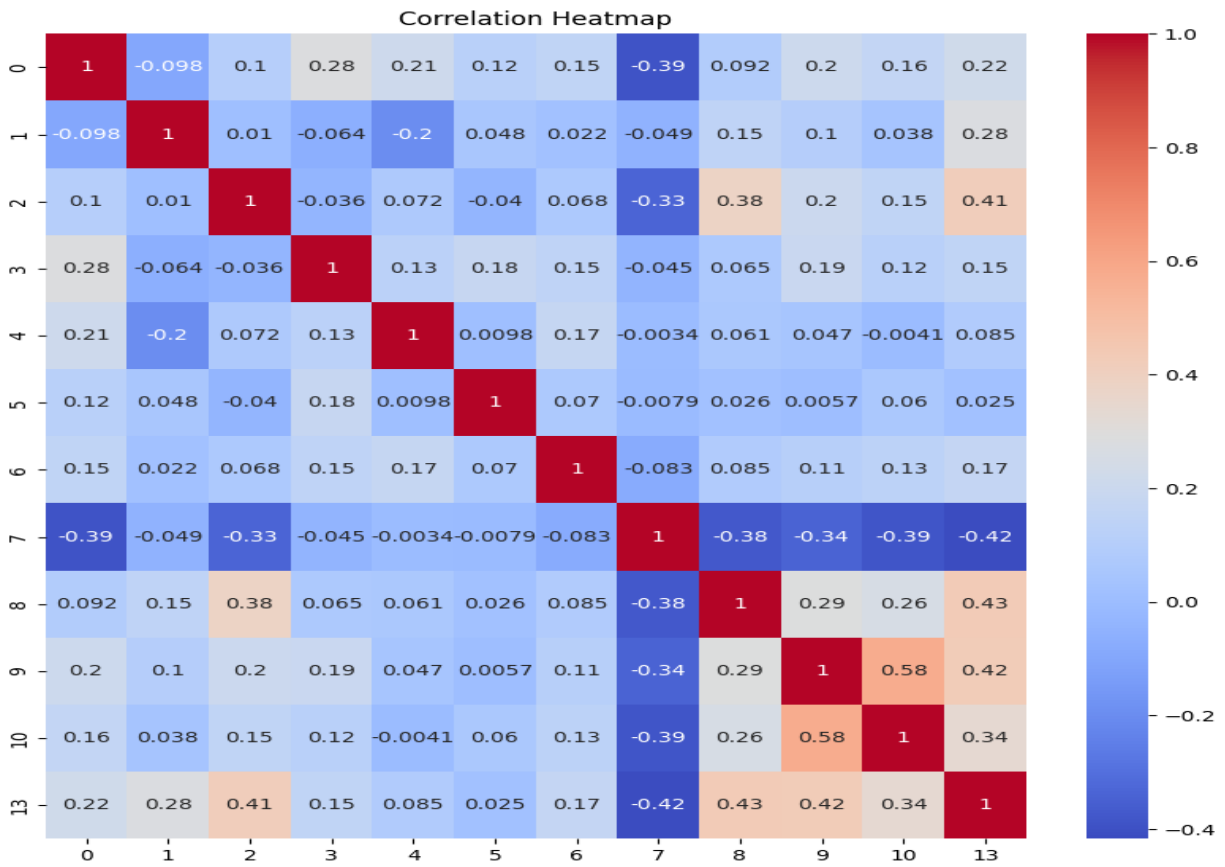
- Converted the classification to binary classification.
Target values with value > 0 were converted to 1
- In the data, missing values were marked as '?'. Replaced them with the mode value of their respective columns, as most values were categorical data.

Columns: {0: 'age', 1: 'sex', 2: 'cp', 3: 'trestbps', 4: 'chol', 5: 'fbs', 6: 'restecg', 7: 'thalach', 8: 'exang', 9: 'oldpeak', 10: 'slope', 11: 'ca', 12: 'thal', 13: 'num'}

Where the target variable is column no. 13 or num.

EDA: I visualised the data by creating several graphs and plots. Scatter plot, heatmap, pairplots, piecharts. Some of these are:





X-axis: age. Y-axis: target.

Some insights are:

- From the correlation heatmap, we can see that most of the features are not correlated, except features no. 9 and 10. Also, target variable no. 13 somewhat correlates with features 2, 8, and 9.
- Made scatter plots and pair plots to find outliers. But there was not much of any outlier.
- On a surprising, from the scatter plot above, we can see that age does not much affect the outcome of heart disease.
- Data has an unequal distribution of sex. One gender data is almost double the other, but it is not specified.

Part B:

Split the data into training and testing in the ratio of 80:20.

Done using sklearn library:

```
x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size = 0.2)
```

Part C:

Trained the decision tree using entropy as the splitting criterion.

Got the following accuracy:

```
#Entropy
Entropy = DecisionTreeClassifier(criterion='entropy')
Entropy.fit(x_train, y_train)
ypredEntropy = Entropy.predict(x_test)
accuracyEntropy = accuracy_score(y_test, ypredEntropy)
print(accuracyEntropy)

0.7049180327868853
```

When trained using 'gini impurity' as the splitting criterion, it got the following accuracy:

```
#Gini
Gini = DecisionTreeClassifier(criterion='gini')
Gini.fit(x_train, y_train)
ypredGini = Gini.predict(x_test)
accuracyGini = accuracy_score(y_test, ypredGini)
print(accuracyGini)

0.6557377049180327
```

So, based on the accuracy score, **‘Entropy’** was the better criterion for node splitting.

Part D:

Performed hyperparameter search for the parameters `min_samples_split` and `max_features` using Grid search, with the help sklearn.

Values passed in the grid were:

```
'min_samples_split': [2, 3, 5, 7, 9, 10],
'max_features': ['sqrt', 'log2', None]
```

The criterion used for splitting was: **‘Entropy’** (from part c).

After running this and using the best parameters that came out were:

```
Hyperparameters Values:
min_samples_split:9
max_features:log2
Accuracy:0.7213114754098361
```

With an accuracy of 0.72.

Part E: Random Forests

Performed hyperparameter search for the parameters `n_estimators`, `max_depth` and `min_samples_split` using Grid search, with the help of sklearn.

Values passed in the grid were:

```
'n_estimators': [50, 100, 150],  
'max_depth': [None, 10, 20, 30],  
'min_samples_split': [2, 3, 5, 7, 9, 10]
```

After running a grid search, the following were the results:

Hyperparameters Values:

n_estimators:50

max_depth:20

min_samples_split:9

Classification Report:

	precision	recall	f1-score	support
0	0.72	0.91	0.81	32
1	0.86	0.62	0.72	29
accuracy			0.77	61
macro avg	0.79	0.76	0.76	61
weighted avg	0.79	0.77	0.76	61

Accuracy:0.7213114754098361

This shows the combination of best hyperparameters, classification reports, and accuracy.