

# ASSIGNMENT 2

1. a) It is important for the trees to be correlated up to a certain extent while maintaining diversity because of the following reasons-  
**Bias-Variance:** Highly correlated trees tend to have low variance but high bias, while highly diverse trees have low bias but high variance. By striking the right balance, you can achieve a lower prediction error by combining both low-bias and low-variance models.  
**Overfitting:** Introducing diversity through bootstrapping and feature randomization helps reduce overfitting, while having some level of correlation among the trees helps in making accurate predictions based on the learning from the training data.
- b) The curse of dimensionality can become a problem in Naive Bayes when we have a large number of features(dimensions) in comparison to the amount of data available.  
Some strategies that can be employed to mitigate this problem in practice are as follows
  1. Feature Selection
  2. Dimensionality reduction using PCA, LDA T-Sne
  3. Using L1 regularization
- c) Yes it will affect the inference results. It can face the following problems
  1. If the new value is not present in the training dataset, its count will be zero. As a result, the conditional probability for that value will be zero, which can lead to issues in the classification process as that class will **never be classified**.
  2. The absence of the new value in the training data means that the model does not have any information about the relationship between that value and the class labels, hence resulting in **inefficient training**.

To mitigate the problem of encountering new values in Naive Bayes, we can use the following approaches:

1. Laplace Smoothing
2. Discounted Smoothing
3. Bayesian Estimation

Now, let's illustrate this issue with an example:

Let us consider a simple dataset with one feature(ice-cream) and one target label(taste) with values 0 and 1 now consider ice-cream cassata which is not present in the training data the conditional probability for both classes (0 and 1) will be 0, hence it won't be possible to classify the ice-cream cassata.

d) Yes, Information Gain can be biased if some attributes have more cardinality (distinct values) than others. This is because Information Gain measures the reduction in uncertainty or entropy, and attributes with higher cardinality tend to have a higher chance of reducing uncertainty by having more values. This can lead to biases towards attributes with more distinct values during the tree-building process.

To mitigate this problem we can use criterias that are less sensitive to cardinality of attributes such as

1. Gini-impurity
2. Chi-squared Statistic

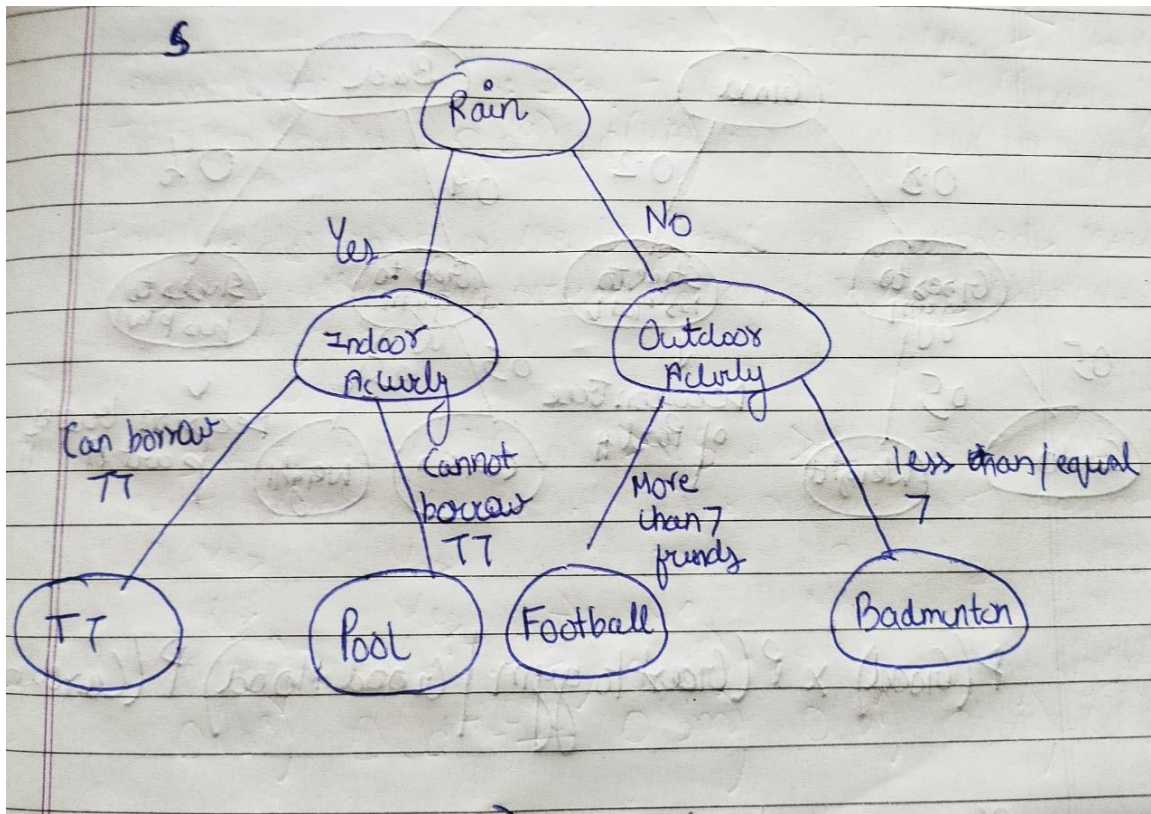
To explain the issues with information gain let us consider an example of a dataset with two attributes “Age group” and “Pin Code” to predict whether someone will like a movie or not.

Age group: It has relatively low cardinality(4).

Pin Code: It has high cardinality. There are thousands of unique pin codes in a large geographic area.

Using Information Gain as the attribute selection criterion might lead the decision tree algorithm to prefer the “Pin Code” attribute due to its higher cardinality. However, the geographic location represented by the pin code may not be very informative in predicting whether someone will like a movie or not in contrast “Age group” could be a more relevant attribute for this prediction task, as people within the same age group may have similar movie preferences..

2.



**Outcome:** Football, Conditional

Probability:  $P(\text{Football} \mid \text{No Rain}, >7 \text{ Friends}) = P(\text{No Rain}) * P(>7 \text{ Friends} \mid \text{No Rain})$

**Outcome:** Badminton

Probability:  $(P(\text{Badminton} \mid \text{No Rain}, \leq 7 \text{ Friends}) = P(\text{No Rain}) * P(\leq 7 \text{ Friends} \mid \text{No Rain})$

**Outcome:** Table Tennis

Probability:  $P(\text{Table Tennis} \mid \text{Rain}, \text{Can Borrow Rackets}) = P(\text{Rain}) * P(\text{Can Borrow Rackets} \mid \text{Rain})$

**Outcome:** Pool

Probability:  $P(\text{Pool} \mid \text{Rain}, \text{Can't Borrow Rackets}) = P(\text{Rain}) * P(\text{Can't Borrow Rackets} \mid \text{Rain})$

b)  $P(\text{Predicted Rain}) = 0.3$

$P(\text{Predicted Clear}) = 0.7$

$P(\text{Predicted Rain} \mid \text{Rain}) = 0.80$

$P(\text{Predicted Clear} \mid \text{Clear}) = 0.90$

$P(\text{Predicted Clear} \mid \text{Rain}) = 0.20$

$P(\text{Predicted Rain} \mid \text{Clear}) = 0.10$

$$P(\text{Predicted Rain}) = P(\text{Rain}) * P(\text{Predicted Rain} | \text{Rain}) + P(\text{Clear}) * P(\text{Predicted Rain} | \text{Clear})$$

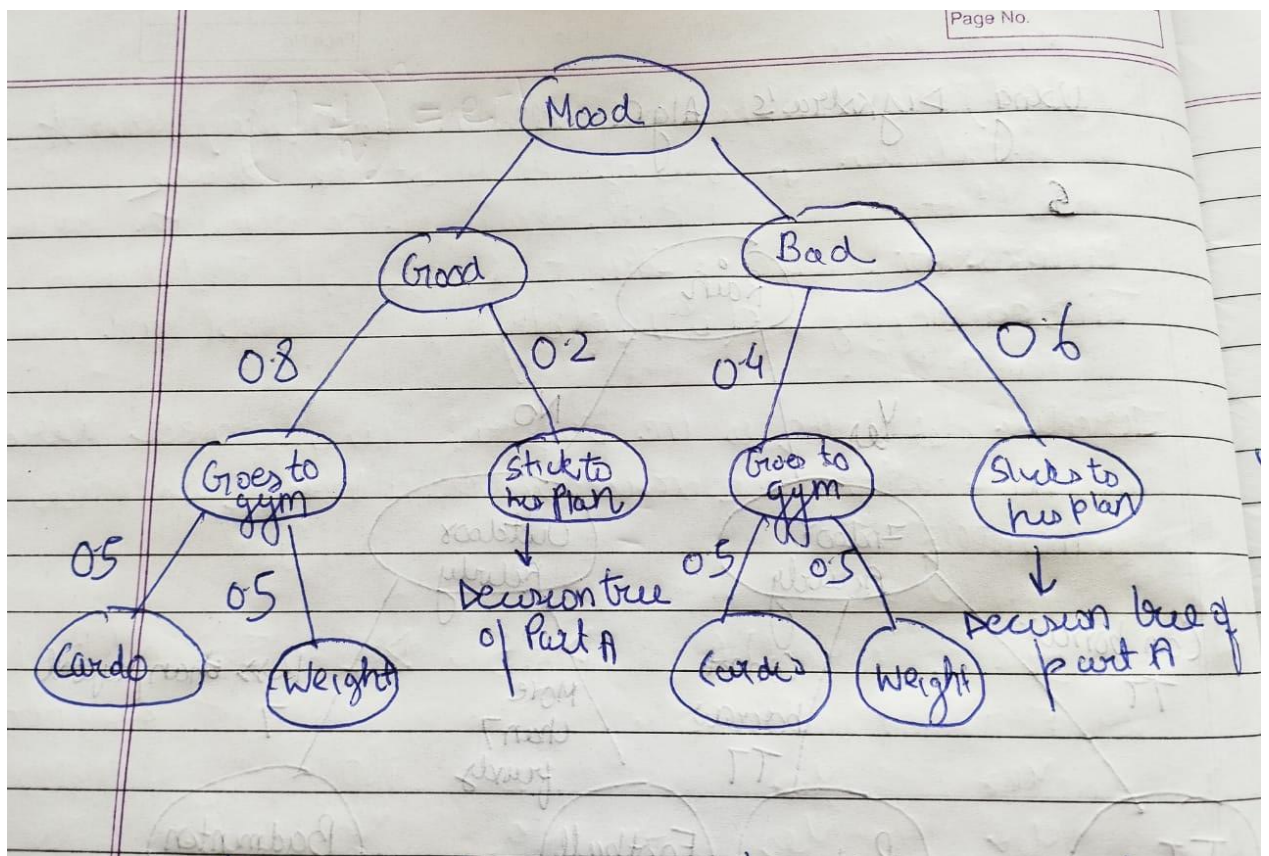
$$0.3 = P(\text{Rain}) * 0.8 + (1 - P(\text{Rain})) * 0.1$$

$$0.2 = 0.7 * P(\text{Rain})$$

$$2/7 = P(\text{Rain})$$

$$\begin{aligned} P(\text{Rain} | \text{Predicted Rain}) &= P(\text{Predicted Rain} | \text{Rain}) * P(\text{Rain}) / P(\text{Predicted Rain}) \\ &= 0.8 * P(\text{Rain}) / 0.3 \\ &= 16/21 \\ &= 0.76 \end{aligned}$$

c)



**Outcome :** Cardio

Probability:  $P(\text{Good Mood}) * P(\text{Go to Gym} | \text{Good Mood}) * P(\text{Cardio} | \text{Good Mood, Go to Gym}) + P(\text{Bad Mood}) * P(\text{Go to Gym} | \text{Bad Mood}) * P(\text{Cardio} | \text{Bad Mood, Go to Gym})$

**Outcome :** Weight Training

Probability:  $P(\text{Good Mood}) * P(\text{Go to Gym} | \text{Good Mood}) * P(\text{Weight Training} | \text{Good Mood, Go to Gym}) + P(\text{Bad Mood}) * P(\text{Go to Gym} | \text{Bad Mood}) * P(\text{Weight Training} | \text{Bad Mood, Go to Gym})$

**Outcome** : Stick to Plan

**Probabilities of Outcomes** of previous decision tree will be multiplied by  $(P(\text{Bad Mood}) * P(\text{Stick to Plan} | \text{Bad Mood}) + P(\text{Good Mood}) * P(\text{Stick to Plan} | \text{Good Mood}))$  as this represents  $P(\text{Stick to Plan})$

d) Rahul had 7 hours of sleep

$P(\text{Good Mood}) = 0.6$ ,  $P(\text{Bad Mood}) = 0.4$

$P(F = 7 | \text{Good Mood}) = 0.7$ ,  $P(F = 7 | \text{Bad Mood}) = 0.45$

$P(F=7) = P(\text{Good Mood}) * P(F = 7 | \text{Good Mood}) + P(\text{Bad Mood}) * P(F = 7 | \text{Bad Mood})$

$= 0.6 * 0.7 + 0.4 * 0.45 = 0.72$

However his mood is also dependent on number of hours he slept

$P(\text{Good Mood} | F = 7) = P(F = 7 | \text{Good Mood}) * P(\text{Good Mood}) / P(F = 7) = 0.58$

$P(\text{Bad Mood} | F = 7) = P(F = 7 | \text{Bad Mood}) * P(\text{Bad Mood}) / P(F = 7) = 0.25$

In the expressions calculated in previous parts  $P(\text{Good Mood})$  and  $P(\text{Bad Mood})$  will be replaced by  $P(\text{Good Mood} | F = 7)$  and  $P(\text{Bad Mood} | F = 7)$  respectively

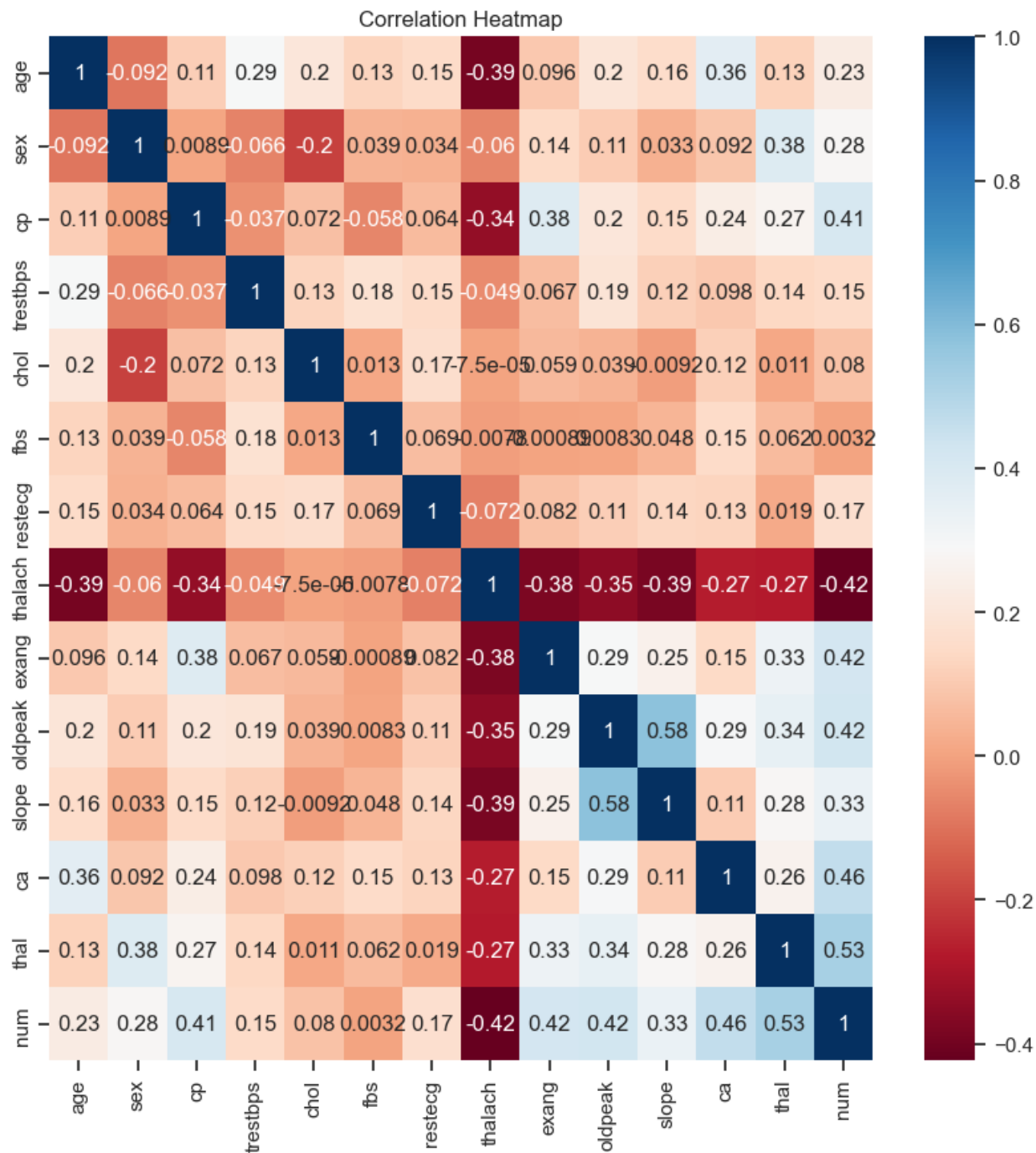
$P(\text{Cardio}) = 0.58 * 0.8 * 0.5 + 0.25 * 0.4 * 0.5 = 0.282$

$P(\text{Weight Training}) = 0.58 * 0.8 * 0.5 + 0.25 * 0.4 * 0.5 = 0.282$

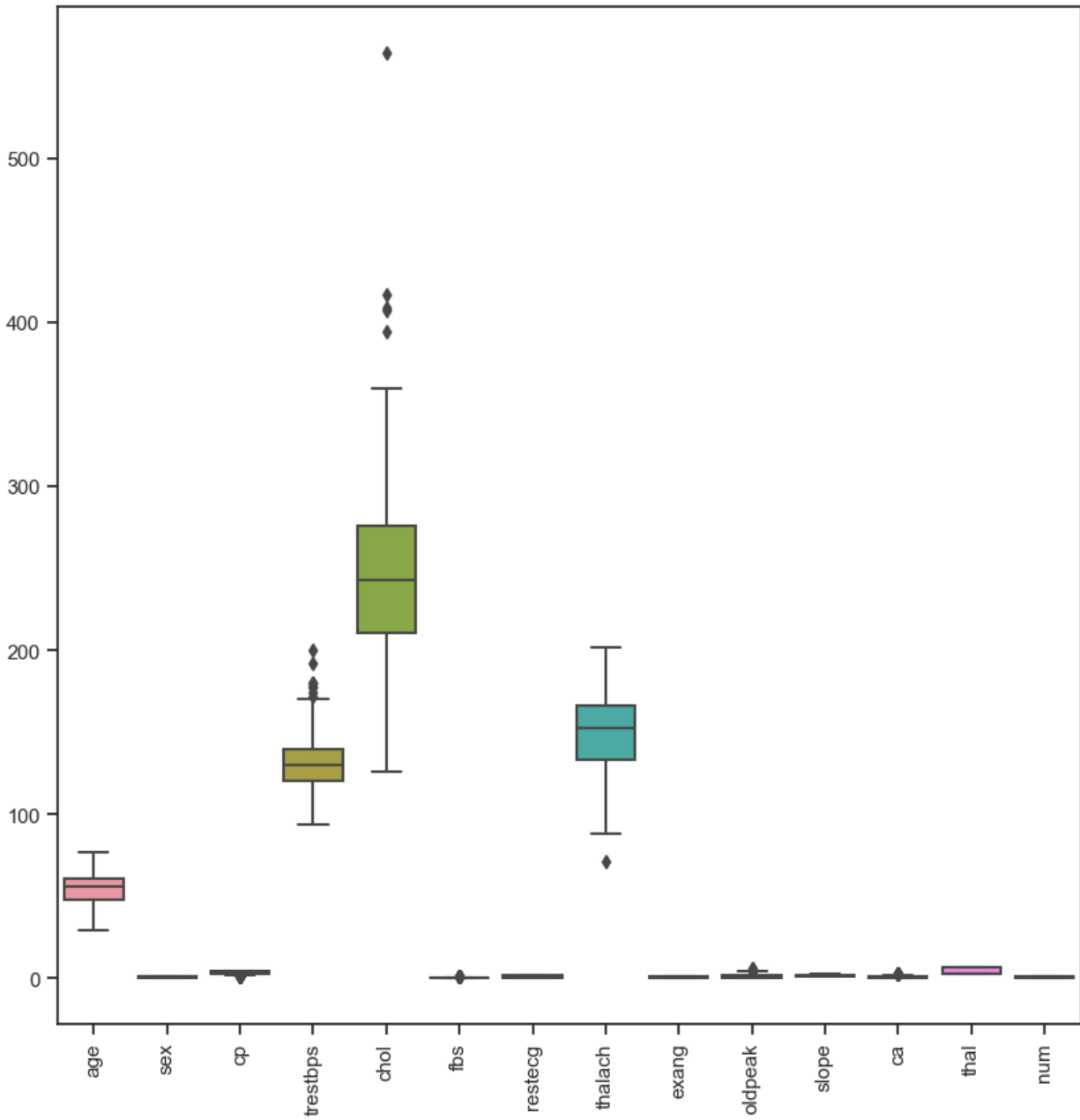
$P(\text{Sticks to Plan}) = 0.58 * 0.6 + 0.25 * 0.2 = 0.266$

**So he will either do cardio or weight training**

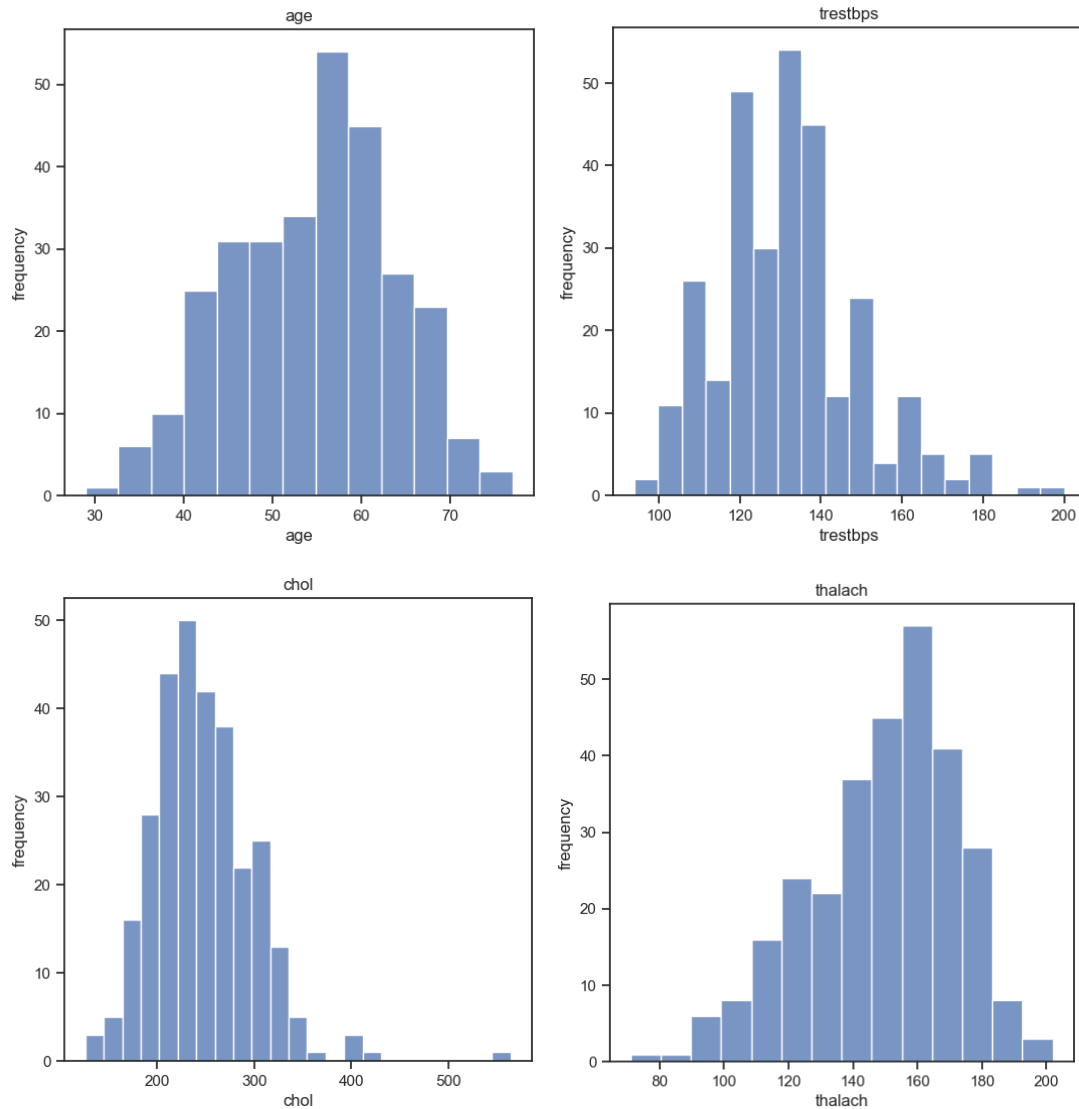
### 3. Exploratory data analysis



Box Plots of Features



## Histograms



c) Best criteria is entropy with training accuracy 1.0  
Testing accuracy for the same is : 0.7333333333333333

d) Best parameters are {'max\_features': 'log2', 'min\_samples\_split': 15}  
Training score is : 0.9029535864978903  
Testing score is : 0.75

e) Best hyperparameters for Random Forest: {'max\_depth': 20, 'min\_samples\_split': 30, 'n\_estimators': 250} with cross-validated accuracy 0.8524822695035462



Classification Report on Test Data:					
	precision	recall	f1-score	support	
0	0.76	0.93	0.84	30	
1	0.91	0.70	0.79	30	
accuracy			0.82	60	
macro avg	0.83	0.82	0.81	60	
weighted avg	0.83	0.82	0.81	60	