

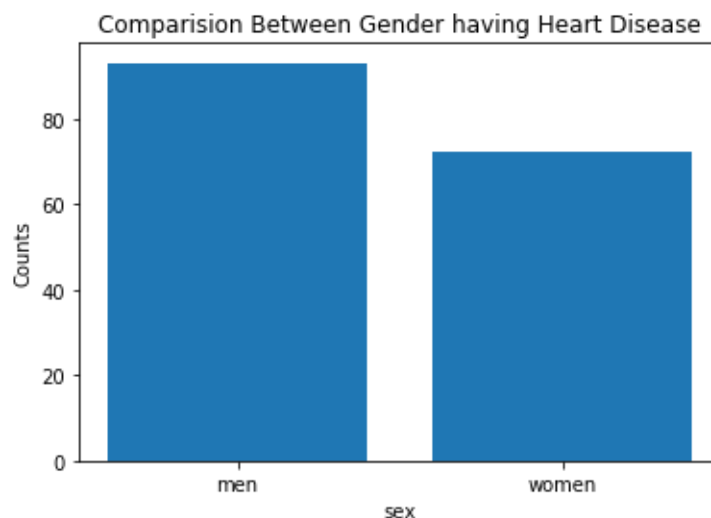
## 1. Data Exploration

- Sentiment Analysis Dataset Consist of total 18,389 records
- Dataset is divided into training and test set in the ratio of 80-20%
- There are two attributes in which Phrases represent review of movies and Sentiment represents the target corresponding to reviews.
- Targets are 0,1,2,3 and 4 where 0,1 are considered as negative reviews, 2 is considered as neutral and 3 and 4 are considered as positive reviews.

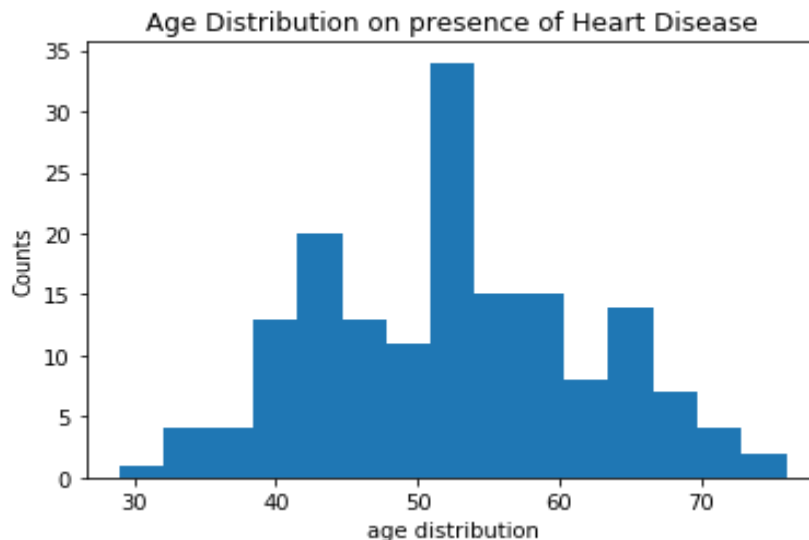
❖ **df['target'].value\_counts()**

1 165(present of heart disease)

0 138(absence of heart disease)



**Fig:** Counts of male and female having Heart Disease which clearly shows that risk of heart disease in male is higher than the female.



**Fig:** above figure is the age distribution according to the presence of Heart Disease. It clearly shows below the age 40 there is a minimal chance of Heart Disease whereas between the ages 50 to 55 there is a higher chance of having Heart Disease.

## 2. Feature Extraction and preprocessing

### ❖ `df.isnull().sum()`

There were no null values in all the records of the variables.

- ❖ Training and test data are splitted to `X_train`, `y_train`, `X_test` and `y_test` where `X_train` and `X_test` hold the features and `y_train` and `y_test` hold the target variable.
- ❖ Since before training we have encoded the categorical datas using `OneHotEncoder`.
- ❖ `OneHotEncoder` encodes categorical features as a one-hot numerical array.
- ❖ Later we do the same procedure to preprocess the categorical features of test data.
- ❖ After encoding, the total number of both test and training dataset columns contain 23 features.

### 3. Grid Search

- Here, we used the Grid search for tuning the hyperparameter.
- Grid Search uses the Cross Validation technique to find the best hyperparameter.

**Different hyperparameters defined for each classifier is shown below:**

S.No	Classifier	Hyperparameters
1	<a href="#">SVC</a>	<b>Kernel:</b> linear, rbf <b>C:</b> 1,10,5
2	<a href="#">DecisionTreeClassifier</a>	<b>max_depth:</b> 2,3,5,8,6,4 <b>min_samples_leaf :</b> 12, 8, 2,14
3	<a href="#">RandomForestClassifier</a>	<b>N_estimators:</b> 10,50,100 <b>max_features:</b> 1,0.5,0.8, 0.1 <b>Min_samples_leaf:</b> 12,1,5,8
4	<a href="#">ExtraTreesClassifier</a>	<b>N_estimators:</b> 10,50,100 <b>max_features:</b> 1,0.5,0.8, 0.1, <b>Min_samples_leaf:</b> 12,1,5,8
5	<a href="#">AdaBoostClassifier</a>	<b>N_estimators:</b> 5,10,20,30 <b>Learning_rate:</b> 0.5,0.1,0.001,0.0001
6	<a href="#">GradientBoostingClassifier</a>	<b>N_estimators:</b> 10,20,40,45 <b>learning_rate:</b> 0.05,0.1, 0.01, 0.2, 0.5 <b>Min_samples_split:</b> 2,3,5

S.No	Classifier	Best Hyperparameters
1	SVC	<b>Kernel:</b> linear <b>C:</b> 1
2	DecisionTreeClassifier	<b>max_depth:</b> 3 <b>min_samples_leaf :</b> 14
3	RandomForestClassifier	<b>N_estimators:</b> 50 <b>max_features:</b> 1 <b>Min_samples_leaf:</b> 5
4	ExtraTreesClassifier	<b>N_estimators:</b> 10 <b>max_features:</b> 0.1, <b>Min_samples_leaf:</b> 5
5	AdaBoostClassifier	<b>N_estimators:</b> 30 <b>Learning_rate:</b> 0.1
6	GradientBoostingClassifier	<b>N_estimators:</b> 10 <b>learning_rate:</b> 0.1 <b>Min_samples_split:</b> 2

#### 4. Model Evaluation and Comparison

Below table shows the accuracy of each classifier for the dataset of Heart Disease.

Classifier	F1-macro
SVC	
DecisionTreeClassifier	
RandomForestClassifier	
ExtraTreesClassifier	
AdaBoostClassifier	
GradientBoostingClassifier	

The accuracy rate of GradientBoostingClassifier is higher compared to other five classifiers. Hence, GradientBoostingClassifier with best hyperparameter **N\_estimators**: 10, **learning\_rate**: 0.1 and **Min\_samples\_split**: 2 is chosen as the best classification algorithm for the Heart Disease Dataset.