

INTRODUCTION

This notebook has been made as an attempt to analyze the data about heart related conditions of various patients. After analyzing the data, an attempt has been made to predict the possibility of a person having heart disease based on the given attributes. The starting point of this journey is the analysis of the given data using different visualization techniques. Gradually, we move onto the process of modelling in which we try out different algorithms to predict the occurrence of heart disease based on the given attributes. A number of models are tried and one of them is shortlisted for further fine-tuning based on its accuracy. The selected model is used to find out the most significant variables that contribute to a correct prediction. The final model uses only the most important variables to predict whether a person has a heart disease or not.

DATA

Data for this project is collected from the link mentioned below. The entire dataset is split into training and validation parts. The training part of the dataset is used to train the model and the validation part of the dataset is used to test the accuracy of the model. Validation data is necessary to gauge the performance of the model when it deals with completely unknown data. This gives a rough idea about how the model will perform in real life.

The link for the complete dataset is:
<https://www.kaggle.com/chenrgs/heart-disease-cleveland-uci>

DATA DICTIONARY

Data dictionary is like the legend. It describes what each variable in the dataset stands for. The data dictionary of the dataset is explained as follows:

- There are 14 features. The first 13 are the variables. The 14th feature, "condition", is our **prediction target**.
- 1. **age**: age in years
 - 1. **sex**: sex (1 = male; 0 = female)
 - 1. **cp**: chest pain type
 - Value 0: typical angina
 - Value 1: atypical angina
 - Value 2: non-anginal pain
 - Value 3: asymptomatic
 - 1. **trestbps**: resting blood pressure (in mm Hg on admission to the hospital)
anything above 130-140 is typically cause for concern.
 - 1. **chol**: serum cholesterol in mg/dl
serum = LDL + HDL * .2 * triglycerides,
above 200 is cause for concern
 - 1. **fbbs**: (fasting blood sugar > 120 mg/dl) (1 = true; 0 = false)
>126 mg/dl. signals diabetes.
 - 1. **restecg**: resting electrocardiographic results
 - Value 0: normal
 - Value 1: having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of > 0.05 mV)
 - can range from mild symptoms to severe problems.
 - = signals non-normal heart beat.
 - Value 2: showing probable or definite left ventricular hypertrophy by Estes' criteria.
 - Enlarged heart's main pumping chamber.
 - 1. **thalach**: maximum heart rate achieved
 - 1. **exang**: exercise induced angina (1 = yes; 0 = no)
 - 1. **oldpeak**: ST depression induced by exercise relative to rest
Looks at stress of heart during exercise.
Unhealthy heart will stress more.
 - 1. **slope**: the slope of the peak exercise ST segment
 - 0: Upsloping: better heart rate with exercise (uncommon).
 - 1: Flat/sloping: minimal change (typical healthy heart).
 - 2: Downsloping: signs of unhealthy heart.
 - 1. **ca**: number of major vessels (0-3) colored by fluoroscopy
Colored vessel means the doctor can see the blood passing through.
The more blood movement the better (no clots).
 - 1. **thal**: Thallium stress result.
 - 0 = normal;
 - 1 = fixed defect;
 - 2 = reversible defect
 - 1. **condition**:
 - 0 = no disease;
 - 1 = disease

END GOAL

The 2 main purposes of this project are:

- To create a model which can predict the presence of heart disease with a reasonable accuracy (TARGETED ACCURACY:85%).
- To find out which of the given variables have a critical contribution towards a heart ailment.

IMPORTING REQUIRED LIBRARIES AND METHODS

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
sns.set()

%matplotlib inline

from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.metrics import gradient_boosting_classifier,
BaggingClassifier, AdaBoostClassifier, ExtraTreesClassifier

from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.metrics import precision_score, recall_score, f1_score
from sklearn.metrics import plot_roc_curve
```

LOADING THE DATA

```
In [2]: df = pd.read_csv('heart_cleveland_upload.csv')
# Reading the CSV file as a pandas DataFrame
df.shape
# Checking the dimensions of the DataFrame (Rows:297,Columns:14)
```

Out[2]: (297, 14)

EXPLORATORY DATA ANALYSIS

POINTS TO BE ADDRESSED:

- What kind of data do we have and how do we treat them differently?
- Is there any missing data? If yes then how do we deal with it?
- How does the presence of outliers affect our prediction?
- How do some individual features affect the prediction target?

```
In [3]: df.head()
# Looking at the first 5 rows of the entire dataset
```

Out[3]:

	age	sex	cp	trestbps	chol	fbbs	restecg	thalach	exang	oldpeak	slope	ca	thal	condition
0	69	1	0	160	234	1	2	131	0	0.1	1	1	0	0
1	69	0	0	140	239	0	0	151	0	1.8	0	2	0	0
2	66	0	0	150	226	0	0	114	0	2.6	2	0	0	0
3	65	1	0	138	282	1	2	174	0	1.4	1	1	0	1
4	64	1	0	110	211	0	2	144	1	1.8	1	0	0	0

```
In [4]: df.condition.value_counts()
# Checking the number of cases of heart disease in the entire dataset
```

Out[4]:

	0	1
count	160	137
Name: condition, dtype: int64		

```
In [5]: df['condition'].value_counts().plot(kind='bar', color=['teal','crimson']);
plt.xlabel('Condition')
plt.ylabel('No. of patients')
```

Out[5]: Text(0, 0.5, 'No. of patients')

INFERENCE: Out of the 297 entries in the dataset, 160 patients don't have any heart disease and 137 patients suffer from heart ailments

```
In [6]: df.info() # Checking some general information
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 297 entries, 0 to 296
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   age         297 non-null    int64
 1   sex         297 non-null    int64
 2   cp          297 non-null    int64
 3   trestbps    297 non-null    int64
 4   chol        297 non-null    int64
 5   fbs         297 non-null    int64
 6   restecg     297 non-null    int64
 7   thalach     297 non-null    int64
 8   exang       297 non-null    int64
 9   oldpeak     297 non-null    float64
10  slope       297 non-null    int64
11  ca          297 non-null    int64
12  thal        297 non-null    int64
13  condition   297 non-null    int64
dtypes: float64(1), int64(13)
memory usage: 32.6 KB
```

INFERENCE: Since the non-null count of each variable is 297, there are no missing values in the dataset

```
In [7]: df.describe() # A brief statistical summary of the dataset
```

Out[7]:

	age	sex	cp	trestbps	chol	fbbs	restecg	thalach	exang	oldpeak	slope	ca	thal	condition
count	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000	297.000000
mean	54.542088	0.676768	2.158249	131.693603	247.350168	0.144781	0.996633	149.599327	0.3261	0.949736	0.468500	0.964850	0.964850	0.559158
std	9.049736	0.468500	0.964850	17.762806	51.997583	0.352474	0.994914	22.941562	0.4689	1.000000	0.000000	0.000000	0.000000	0.500000
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	48.000000	0.000000	2.000000	120.000000	211.000000	0.000000	0.000000	153.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	56.000000	1.000000	2.000000	130.000000	243.000000	0.000000	1.000000	153.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	61.000000	1.000000	3.000000	140.000000	276.000000	0.000000	2.000000	166.000000	1.000000	0.000000	0.000000	0.000000	0.000000	1.000000
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000	202.000000	1.000000	0.000000	0.000000	0.000000	0.000000	1.000000

INFERENCES:

- Average age of all patients in the dataset is 54.54 years.
- Average of the sex column is 0.68 which indicates that male patients slightly outnumber the female patients in the dataset.
- The age range of the patients in the dataset is from 29 to 77.

```
In [8]: df.hist(figsize=(18,20)) # Checking the overall distribution of all variables
```

Out[8]:

INFERENCES:

- The dataset contains both numerical and categorical data.

NUMERIC DATA:

- **age**: Follows normal distribution approximately. Most of the patients are 50 to 70 years old.
- **trestbps**: Follows normal distribution approximately but there are a few outliers in the range of 170 to 200.
- **chol**: Follows normal distribution approximately. Very few outliers are present, some in the range 350-420 and very few above 500
- **thalach**: This too can be approximated as a normal distribution. Outliers are present in the range 75-100 and some around 200

CATEGORICAL DATA:

- **sex**: The number of male patients is more than twice the number of female patients in the dataset. This non-uniformity can introduce some deviations in our analysis if sex and condition are found to be strongly related.
- **fbs**: Most of the patients in our dataset are non-diabetic. For increasing our accuracy, we might have to include more diabetic patients in our analysis.
- **exang**: Most of the patients do not suffer from exercise induced angina. This is yet another skewness in our data which might affect the accuracy of our model.
- **slope**: Very few patients have a downsloping condition.

This brief overview of our data suggests that there is still scope for improving the dataset. Some suggestions for improvement are as follows:

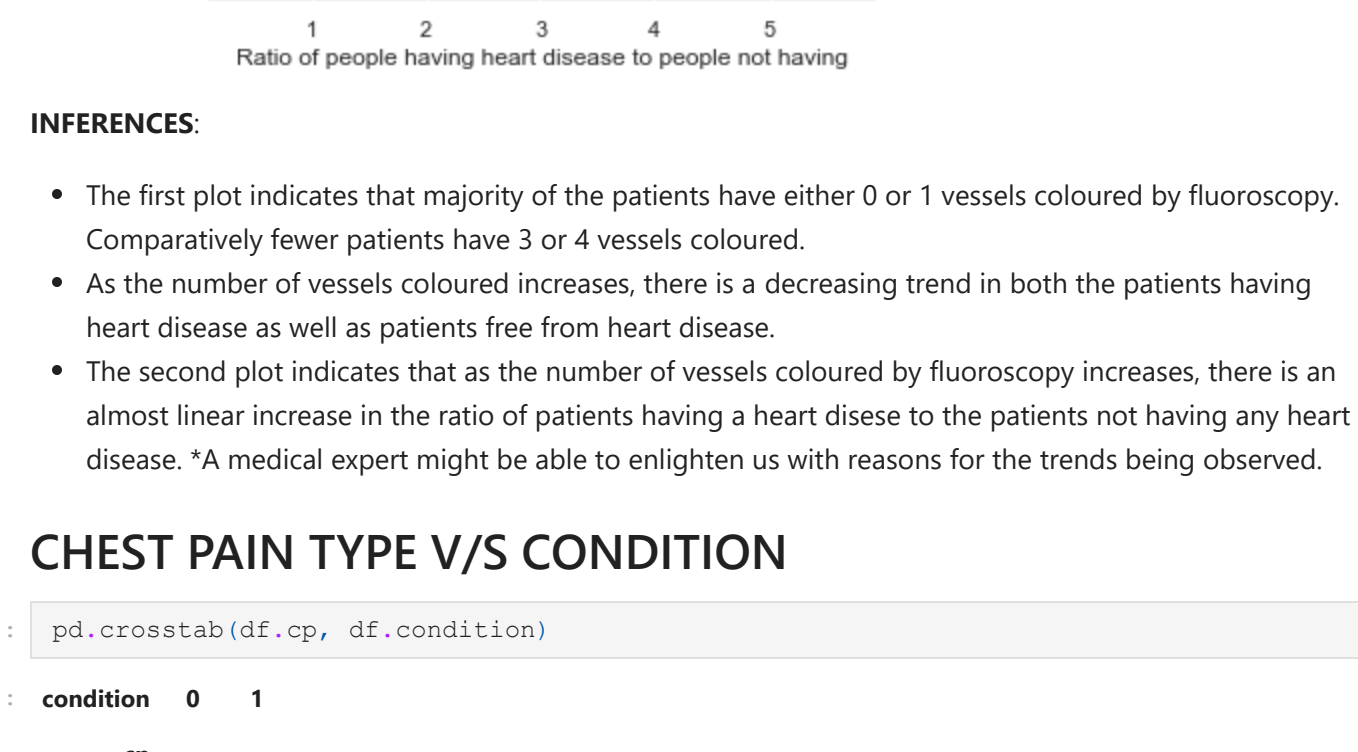
- Removal of outliers from numeric data wherever possible.
- Inclusion of a greater number of female participants in the dataset.
- Inclusion of more diabetic patients in the dataset.
- Inclusion of more patients having exercise induced angina.

CORRELATION MATRIX AND ITS VISUALISATION

The correlation matrix helps us to understand the interrelation between our features as well as between different features and our prediction target. This will give us some preliminary insight into the most and the least significant features. There are two types of correlations that can occur:

- **POSITIVE CORRELATION**: Two variables are said to be positively correlated when increase in one leads to an increase in the other and a decrease in one leads to a decrease in the other.
- **NEGATIVE CORRELATION**: Two variables are said to be negatively correlated when increase in one leads to a decrease in the other and a decrease in one leads to an increase in the other.

```
In [9]: # Constructing the correlation matrix
cor_matrix = df.corr()
fig, ax = plt.subplots(figsize=(16,12))
sns.heatmap(cor_matrix, annot=True, linewidth=0.5, fmt=".2f", cmap='YlGn')
```



The darker shades of green represent strong positive correlation whereas the lighter shades tend towards white indicate strong negative correlation. The colourbar on the right shows that 1 is the maximum correlation value attained whereas the minimum value lies somewhere around -0.4. Correlation of 1 is obtained when a variable is compared to itself.

NOTABLE INFERENCES:

- Thallium induced stress(thal) has a positive correlation (0.52) with condition.
- Other features showing a strong positive correlation with condition are:
 1. Number of major vessels coloured by fluoroscopy(ca):0.46
 2. Exercise induced angina(exang):0.42
 3. ST depression induced by exercise(oldpeak):0.42
 4. Type of chest pain(cp):0.41
- Maximum heart rate achieved(thalach) has a strong negative correlation (-0.42) with condition.

MAJOR VESSELS COLOURED BY FLUOROSCOPY V/S CONDITION

```
In [10]: pd.crosstab(df.ca, df.condition) # making a cross-table
```

Out[10]:

condition	0	1
ca		
0	129	45
1	21	44
2	7	31
3	3	17

```
In [11]: # Visualizing the crosstables
pd.crosstab(df.ca, df.condition).plot(kind='line', figsize=(10,6), color=['teal','crimson'])
plt.title('Major Vessels coloured by Fluoroscopy vs Heart Disease')
plt.xlabel('Amount')
plt.ylabel('Amount')
```

Out[11]:

```
In [12]: # Visualizing cross-table
df1 = pd.DataFrame(pd.crosstab(df.ca, df.condition))
disease = df1.loc[:,1]
no_disease = df1.loc[:,0]
df1['Ratio: 1/0'] = disease/no_disease
# Ratio of people having heart disease to people not having df1
```

Out[12]:

condition	0	1	Ratio: 1/0
ca			
0	129	45	0.348837
1	21	44	2.095238
2	7	31	4.428571
3	3	17	5.666667

```
In [13]: plt.plot(df1['Ratio: 1/0'],range(0,4))
# Plotting this ratio against major vessels coloured by fluoroscopy
plt.title('Major vessels coloured by fluoroscopy v/s Ratio if people having disease to people not having')
plt.xlabel('Ratio of people having heart disease to people not having')
plt.ylabel('Major vessels coloured by fluoroscopy')
```

Out[13]:

INFERENCES:

- The first plot indicates that majority of the patients have either 0 or 1 vessels coloured by fluoroscopy. Comparatively fewer patients have 3 or 4 vessels coloured.
- As the number of vessels coloured increases, there is a decreasing trend in both the patients having heart disease as well as patients free from heart disease.
- The second plot indicates that as the number of vessels coloured by fluoroscopy increases, there is an almost linear increase in the ratio of patients having a heart disease to the patients not having any heart disease. A medical expert might be able to enlighten us with reasons for the trends being observed.

CHEST PAIN TYPE V/S CONDITION

```
In [14]: pd.crosstab(df.cp, df.condition)
```

Out[14]:

condition	0	1
cp		
0	16	7
1	40	9
2	65	18
3	39	103

```
In [15]: # Crosstab Visualization
pd.crosstab(df.cp, df.condition).plot(kind='bar',
figsize=(10,6),
color=['teal','crimson'])
plt.title('Chest Pain Type vs Heart Disease')
plt.xlabel('Chest Pain Type')
plt.ylabel('Amount')
plt.legend(['No Disease', 'Disease'])
```

Out[15]:

Chest Pain Types(FROM DATA DICTIONARY)

- Value 0: typical angina
- Value 1: atypical angina
- Value 2: non-anginal pain
- Value 3: asymptomatic

INFERENCES:

- As seen in the above plot, having asymptomatic chest pain greatly increases a person's chance of having a heart disease

We know that the older a person becomes, the older his/her heart becomes. Hence it might be possible to get some insights into the data if we plot the age and chest pain type together and check the variation in condition.

RELATION BETWEEN AGE AND CHEST PAIN TYPE FOR DETERMINING CONDITION

```
In [16]: plt.figure(figsize=(10,8))
plt.scatter(df.age[df.condition==0],df.cp[df.condition==0], color='teal')
plt.scatter(df.age[df.condition==1],df.cp[df.condition==1], color='crimson')
```

Out[16]:

INFERENCES:

- From the above plot, 2 groups of people can be identified who tend to be highly susceptible to heart diseases:
 1. People in the age group (40-70) having asymptomatic chest pain.
 2. People in the age group (55-70) having non-anginal chest pain.

RELATION BETWEEN AGE AND MAXIMUM HEART RATE ACHIEVED

```
In [17]: plt.figure(figsize=(10,6))
plt.scatter(df.age[df.condition==0],df.thalach[df.condition==0], color='teal')
plt.scatter(df.age[df.condition==1],df.thalach[df.condition==1], color='crimson')
```

Out[17]:

This plot does not give us any concrete quantitative conclusion, however a general inference that can be made is:

- Older people having lower values of maximum heart rate are more likely to suffer from a heart ailment.

MODELLING

In this section, our goal is to make a predictive model for determining the value of condition based on the input features. Since this is a binary classification problem, we will try some classification algorithms provided by Sci-kit Learn python. The algorithms that we will try are:

- LOGISTIC REGRESSION
- K-NEIGHBORS CLASSIFIER
- RANDOM FOREST
- GRADIENT BOOSTING CLASSIFIER
- BAGGING CLASSIFIER
- ADA BOOST CLASSIFIER
- EXTRA TREES CLASSIFIER

```
In [18]: X = df.drop('condition', axis=1)
y = df.condition
```

```
In [19]: # Splitting data into training and test
np.random.seed(42)
test_X, train_X, test_y, train_y = train_test_split(X,y,test_size=0.2)
```

```
In [20]: # CREATING A DICTIONARY OF POSSIBLE MODELS TO BE USED
models = {'Logistic Regression':LogisticRegression(),
'KNN':KNeighborsClassifier(),
'Random Forest':RandomForestClassifier(),
'GBClassifier':GradientBoostingClassifier(),
'BaggingClassifier':BaggingClassifier(),
'AdaBoostClassifier':AdaBoostClassifier(),
'ETClassifier':ExtraTreesClassifier()}

# Defining a function containing all 3 models
```

```
def fit_and_score(models, test_X, train_X, test_y, train_y):
"""
For fitting and evaluating different machine learning models
models: A dictionary of sklearn ML models relevant to the problem
test_X: Testing data features
train_X: Training data features
test_y: Testing Target
train_y: Training Target
"""
# Setting a random seed
np.random.seed(42)

# Empty dictionary to store all model scores
model_scores={}
# Looping through all models
for name, model in models.items():
    # Fitting each model with training data
    model.fit(train_X,train_y)
    # Evaluating each model with test data
    model_scores[name]=model.score(test_X,test_y)
return model_scores
```

```
In [21]: models = fit_and_score(models, test_X, train_X, test_y, train_y)
models
```

Out[21]:

```
C:\Users\DELL\anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:762: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP. TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
('Logistic Regression': 0.835430379746836,
'KNN': 0.5738396624472584,
'Random Forest': 0.8016877631390801,
'GBClassifier': 0.754943671086076,
'BaggingClassifier': 0.784810126522784,
'AdaBoostClassifier': 0.6919831223628692,
'ETClassifier': 0.819345991561815)
```

MODEL COMPARISON

```
In [22]: model_compare = pd.DataFrame(models, index=['accuracy'])
model_compare.T.plot.bar(figsize=(10,6))
```

Out[22]:

Our aim is to create a model with at least 85% accuracy. Hence, from the plot, we will only proceed with the tuning of models that have a high accuracy for this dataset. From the chart, the three most highly accurate models for the given dataset are:

- Logistic Regression (83.54%)
- Extra Trees Classifier (81.43%)
- Random Forest Classifier (80.16%)

Hyperparameter tuning is necessary to improve the performance of our model. Each algorithm has a set of parameters which are assigned certain values. Changing the values of these parameters results in slight changes in the working of the algorithms and the results vary accordingly. **CROSS VALIDATION** is a useful technique to evaluate the performance of our model. Training the model on a given set of parameters and testing it on the same parameters is a mistake as we will always get the correct output for the training dataset, but we will never know how the model performed on data that is previously unknown to it.

In **CROSS VALIDATION**, the training dataset is divided into 'k' smaller sets and the following process is followed for each of the 'k' subsets:

- The model is trained using 'k-1' of the subsets.
- The resulting model is validated on the remaining 1 subset.

The performance metrics thus measured would be the average of the values computed in the loop. This process is a little expensive computationally, but it makes sure that not much of the data is wasted.

MODEL IMPROVEMENT BY HYPERPARAMETER TUNING

Hyperparameter tuning is done for the three models using GridSearchCV, which is a Cross Validation approach implemented in scikit learn.

