

## TITLE - Heart Disease Prediction

### Group Members -

Mr. Vaibhav Pawar 2127052

Mr. Avdhoot Kumbhar 2127037

Mr. Yashraj Devrat 2127011

Mr. Shubham Keskar 2127029

Mr. Tushar Mamadge 2127040

Explanation – The Objective of this project is to create a model that can predict the patient's heart disease Status. Another Objective is to explore the data we have been given and find key insights into Heart Disease that could be helpful for the Medical community going forward. Dataset for this project is taken from Kaggle website <https://www.kaggle.com/datasets/priyanka841/heart-disease-prediction-uci>. Here we use Logistic Regression model. In this Dataset there are 303 Rows and 14 Columns, means 302 Persons Data out of which 241 are use for Training the Data and 61 use for Testing the Data.

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
import matplotlib.pyplot as plt
```

### Data Collection and Processing

```
#loading the csv data to a Pandas DataFrame
heart_data = pd.read_csv('/content/heart disease.csv')

# print first 5 rows of the dataset
heart_data.head()
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak
0	63	1	3	145	233	1	0	150	0	2.3
1	37	1	2	130	250	0	1	187	0	3.5
2	41	0	1	130	204	0	0	172	0	1.4
3	56	1	1	120	236	0	1	178	0	0.8
4	57	0	0	120	354	0	1	163	1	0.6

	ca	thal	target
0	0	1	1
1	0	2	1
2	0	2	1
3	0	2	1
4	0	2	1

*# print last 5 rows of the dataset*  
heart\_data.tail()

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang
oldpeak \									
298	57	0	0	140	241	0	1	123	1
0.2									
299	45	1	3	110	264	0	1	132	0
1.2									
300	68	1	0	144	193	1	1	141	0
3.4									
301	57	1	0	130	131	0	1	115	1
1.2									
302	57	0	1	130	236	0	0	174	0
0.0									

	slope	ca	thal	target
298	1	0	3	0
299	1	0	3	0
300	1	2	3	0
301	1	1	3	0
302	1	1	2	0

*# number of rows and columns in the dataset*  
heart\_data.shape

(303, 14)

*# getting some info about the data*  
heart\_data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   age         303 non-null   int64
1   sex         303 non-null   int64
2   cp          303 non-null   int64
3   trestbps    303 non-null   int64
4   chol        303 non-null   int64
5   fbs         303 non-null   int64
6   restecg     303 non-null   int64
7   thalach     303 non-null   int64
```

```

8   exang      303 non-null   int64
9   oldpeak    303 non-null   float64
10  slope      303 non-null   int64
11  ca         303 non-null   int64
12  thal       303 non-null   int64
13  target     303 non-null   int64
dtypes: float64(1), int64(13)
memory usage: 33.3 KB

```

```

# checking for missing values
heart_data.isnull().sum()

```

```

age          0
sex          0
cp           0
trestbps     0
chol         0
fbs          0
restecg      0
thalach      0
exang        0
oldpeak      0
slope        0
ca           0
thal         0
target       0
dtype: int64

```

```

# statistical measures about the data
heart_data.describe()

```

	age	sex	cp	trestbps	chol
count	303.000000	303.000000	303.000000	303.000000	303.000000
mean	54.366337	0.683168	0.966997	131.623762	246.264026
std	9.082101	0.466011	1.032052	17.538143	51.830751
min	29.000000	0.000000	0.000000	94.000000	126.000000
25%	47.500000	0.000000	0.000000	120.000000	211.000000
50%	55.000000	1.000000	1.000000	130.000000	240.000000
75%	61.000000	1.000000	2.000000	140.000000	274.500000
max	77.000000	1.000000	3.000000	200.000000	564.000000

	restecg	thalach	exang	oldpeak	slope
count	303.000000	303.000000	303.000000	303.000000	303.000000
mean	0.000000	149.633696	0.000000	0.000000	0.000000
std	0.000000	27.868852	0.000000	0.000000	0.000000
min	0.000000	79.000000	0.000000	0.000000	0.000000
25%	0.000000	127.000000	0.000000	0.000000	0.000000
50%	0.000000	147.000000	0.000000	0.000000	0.000000
75%	0.000000	163.000000	0.000000	0.000000	0.000000
max	0.000000	201.000000	0.000000	0.000000	0.000000

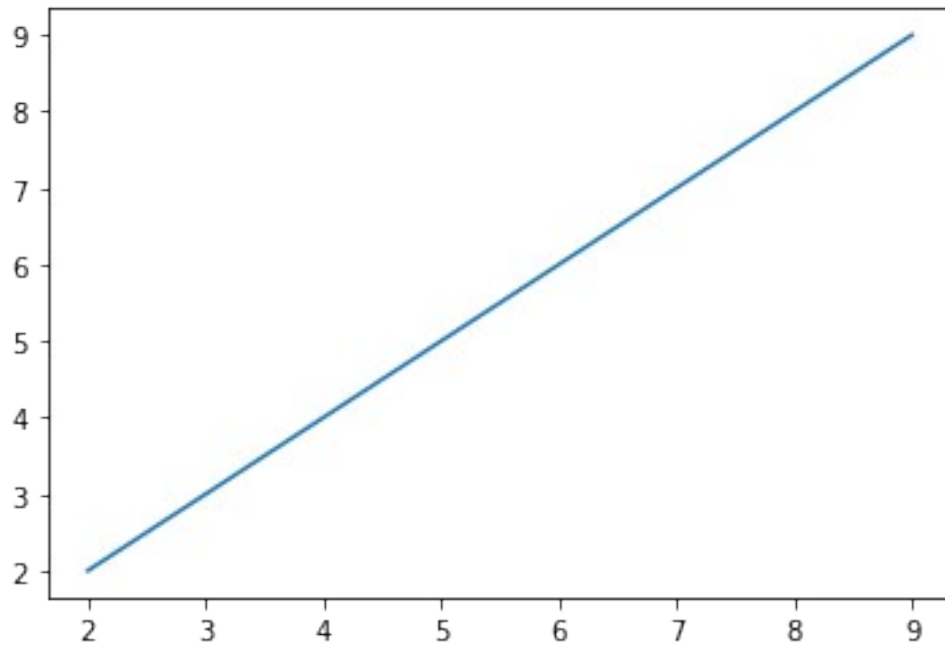
ca \					
count	303.000000	303.000000	303.000000	303.000000	303.000000
mean	0.528053	149.646865	0.326733	1.039604	1.399340
std	0.525860	22.905161	0.469794	1.161075	0.616226
min	0.000000	71.000000	0.000000	0.000000	0.000000
25%	0.000000	133.500000	0.000000	0.000000	1.000000
50%	1.000000	153.000000	0.000000	0.800000	1.000000
75%	1.000000	166.000000	1.000000	1.600000	2.000000
max	2.000000	202.000000	1.000000	6.200000	2.000000

	thal	target
count	303.000000	303.000000
mean	2.313531	0.544554
std	0.612277	0.498835
min	0.000000	0.000000
25%	2.000000	0.000000
50%	2.000000	1.000000
75%	3.000000	1.000000
max	3.000000	1.000000

### Data Visualization

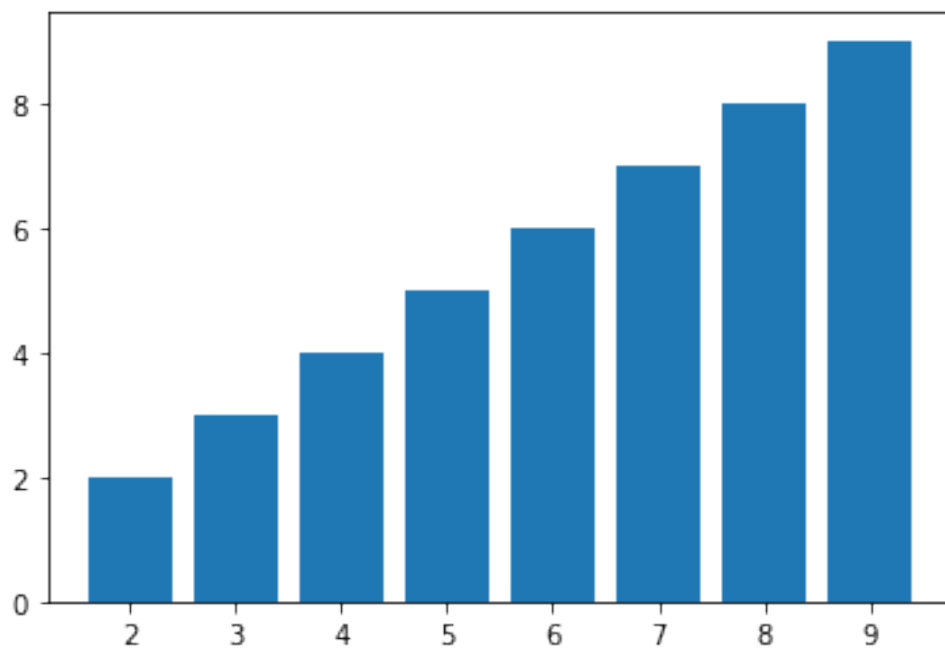
```
x = [2,3,4,5,6,7,8,9]
y = [2,3,4,5,6,7,8,9]
```

```
fig, ax = plt.subplots()
ax.plot(x,y);
```



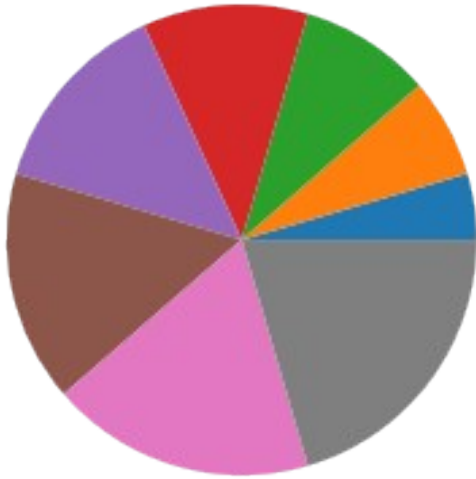
```
fig, ax = plt.subplots()
```

```
ax.bar(x,y);
```



```
fig, ax = plt.subplots()
```

```
ax.pie(y);
```



```
# checking the distribution of Target Variable
heart_data['target'].value_counts()
```

```
1    165
0    138
Name: target, dtype: int64
```

Splitting the Features and Target

```
X = heart_data.drop(columns='target', axis=1)
Y = heart_data['target']
```

```
print(X)
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang
0	63	1	3	145	233	1	0	150	0
1	37	1	2	130	250	0	1	187	0
2	41	0	1	130	204	0	0	172	0
3	56	1	1	120	236	0	1	178	0
4	57	0	0	120	354	0	1	163	1
...	...	...	...	...	...	...	...	...	...
298	57	0	0	140	241	0	1	123	1
299	45	1	3	110	264	0	1	132	0

```

300    68    1    0    144    193    1    1    141    0
3.4
301    57    1    0    130    131    0    1    115    1
1.2
302    57    0    1    130    236    0    0    174    0
0.0

```

```

      slope  ca  thal
0         0   0    1
1         0   0    2
2         2   0    2
3         2   0    2
4         2   0    2
..      ...  ..   ...
298       1   0    3
299       1   0    3
300       1   2    3
301       1   1    3
302       1   1    2

```

```
[303 rows x 13 columns]
```

```
print(Y)
```

```

0      1
1      1
2      1
3      1
4      1
..
298    0
299    0
300    0
301    0
302    0

```

```
Name: target, Length: 303, dtype: int64
```

Splitting the Data into Training data & Test Data

```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y,
test_size=0.2, stratify=Y, random_state=2)
```

```
print(X.shape, X_train.shape, X_test.shape)
```

```
(303, 13) (242, 13) (61, 13)
```

Model Training 1

Logistic Regression

```
model = LogisticRegression()
```

*# training the LogisticRegression model with Training data*

```
model.fit(X_train, Y_train)
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: 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](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

```
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
```

```
LogisticRegression()
```

Model Evaluation 1

Accuracy Score

*# accuracy on training data*

```
X_train_prediction = model.predict(X_train)
```

```
training_data_accuracy = accuracy_score(X_train_prediction, Y_train)
```

```
print('Accuracy on Training data : ', training_data_accuracy)
```

```
Accuracy on Training data : 0.8512396694214877
```

*# accuracy on test data*

```
X_test_prediction = model.predict(X_test)
```

```
test_data_accuracy = accuracy_score(X_test_prediction, Y_test)
```

```
print('Accuracy on Test data : ', test_data_accuracy)
```

```
Accuracy on Test data : 0.819672131147541
```

Model Training 2

RANDOM FOREST CLASSIFIER

```
from sklearn.ensemble import RandomForestClassifier
```

```
model = RandomForestClassifier()
```

```
model.fit(X_train, Y_train)
```

```
RandomForestClassifier()
```

Model Evaluation 2

Accuracy Score



```
# accuracy on training data
X_train_prediction = model.predict(X_train)
training_data_accuracy = accuracy_score(X_train_prediction, Y_train)

print('Accuracy on Training data : ', training_data_accuracy)
```

Accuracy on Training data : 1.0

```
# accuracy on test data
X_test_prediction = model.predict(X_test)
test_data_accuracy = accuracy_score(X_test_prediction, Y_test)

print('Accuracy on Test data : ', test_data_accuracy)
```

Accuracy on Test data : 0.819672131147541

Building a Predictive System

```
input_data = (62,0,0,140,268,0,0,160,0,3.6,0,2,2)
```

```
# change the input data to a numpy array
input_data_as_numpy_array= np.asarray(input_data)
```

```
# reshape the numpy array as we are predicting for only on instance
input_data_reshaped = input_data_as_numpy_array.reshape(1,-1)
```

```
prediction = model.predict(input_data_reshaped)
print(prediction)
```

```
if (prediction[0]== 0):
    print('The Person does not have a Heart Disease')
else:
    print('The Person has Heart Disease')
```

```
[0]
The Person does not have a Heart Disease
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/base.py:451:
UserWarning: X does not have valid feature names, but
RandomForestClassifier was fitted with feature names
"X does not have valid feature names, but"
```