



Heart Disease

Done by : Rawan Al-Qahtani

Dr : Mejdal Al-Qahtani



Introduction

This project is to predict whether the patient has a **10-year risk of future coronary heart disease (CHD)**. The data was taken from Kaggle website as (.csv)file, and it is from a cardiovascular study on residents of the town of Framingham, Massachusetts.

After features checked and data splitting a comparison between (KNN – logistic regression - Random forest - Decision tree) Classifiers was made. First, the hypermeters where found then (Accuracy- Recall-Precision -F1- Confusion Matrix) for classifiers was shown. Then, the most compatible classifier was chosen to be the model for prediction which is **Logistic regression** since it had the highest recall score of 66.9% with the up-sampled data.



Extract The 10 Best Features

```
81]: 1 #extract the 10 best features
      2 features = SelectKBest(score_func=chi2, k=10)
      3 scores = features.fit(x,y).scores_
      4 columns = x.columns
      5
      6 df = pd.DataFrame()
      7 df['score'] = scores
      8 df['column'] = columns
      9
     10 best_feat = df.sort_values(by='score', ascending=False).head(10)
     11 best_feat
```

```
81]:
```

	score	column
10	728.028608	sysBP
14	370.026234	glucose
1	318.439180	age
9	233.735513	totChol
4	178.115009	cigsPerDay
11	151.897052	diaBP
7	92.181857	prevalentHyp
8	38.371005	diabetes
5	35.323423	BPMeds
0	16.707199	male

The highest
features with
effected people



Correlation Heatmap



Data splitting & Sampling

Data splitting

```
[82]: 1 split = StratifiedShuffleSplit(n_splits = 1, test_size=0.2, random_state=42)
      2 for train_index, test_index in split.split(heart_scaled, heart_scaled['TenYearCHD']):
      3     train = heart_scaled.loc[train_index]
      4     test = heart_scaled.loc[test_index]
      5
      6     x_test = test.drop("TenYearCHD", axis=1)
      7     y_test = test["TenYearCHD"]
```

Sampling

```
[83]: 1 #Separate majority and minority classes
      2 majority = train[train["TenYearCHD"] == 0]
      3 minority = train[train["TenYearCHD"] == 1]
```

```
[84]: 1 #upsample minority data
      2 minority_upsampled = resample(minority,
      3                             n_samples = 2500,
      4                             random_state=65)
      5
      6 #combine with majority dataframe
      7 upsampled = pd.concat([majority, minority_upsampled],axis=0)
```

```
[85]: 1 #downsample majority data
      2 majority_downsampled = resample(majority,
      3                                 n_samples=511,
      4                                 random_state=65)
      5
      6 downsampled = pd.concat([minority,majority_downsampled],axis=0)
```

Upsampled data:

0.0 2804

1.0 2500

Name: TenYearCHD, dtype: int64

Downsampled data:

0.0 511

1.0 502

Name: TenYearCHD, dtype: int64



Choosing The Best Hyperparameters



Classifier	The best score / Upsampled	The best score / Downsampled
KNN	0.9245852187028658	0.6811451135241856
Logistic regression	0.6830693815987934	0.6712734452122409
Random forest	0.8374811463046757	0.6712734452122409
Decision tree	0.8849924585218703	0.6238894373149062



Modelling

Accuracy

Logistic Regression	0.7230955259975816
KNN	0.7678355501813785
Decision Tree	0.7412333736396615
Random Forest	0.7799274486094316

Recall

Logistic Regression	0.6031746031746031
KNN	0.18253968253968253
Decision Tree	0.23809523809523808
Random Forest	0.4126984126984127

the logistic regression has the best parameters spicily the recall

Precision

Logistic Regression	0.2980392156862745
KNN	0.20535714285714285
Decision Tree	0.20270270270270271
Random Forest	0.325

F1

Logistic Regression	0.3989501312335958
KNN	0.19327731092436976
Decision Tree	0.21897810218978103
Random Forest	0.3636363636363636

CONFUSION MATRIX



CONFUSION MATRIX

```
In [99]: 1 print('Logistic Regression:\n ', confusion_matrix(y_test, lr.predict(x_test)))
          2 print('-----')
          3 print('KNN: \n', confusion_matrix(y_test, knn.predict(x_test)))
          4 print('-----')
          5 print('Decision Tree: \n', confusion_matrix(y_test, tree.predict(x_test)))
          6 print('-----')
          7 print('Random Forest: \n', confusion_matrix(y_test, clf.predict(x_test)))
          8 print('-----')
```

Logistic Regression:

[[522 179]

→ 50 76 ←

KNN:

[[612 89]

[103 23]]

Decision Tree:

[[583 118]

[96 30]]

Random Forest:

[[593 108]

[74 52]]

Also here at the matrix the
logistic tree has the least
FP and higher TN



Module training for data prediction



```
1 #initialize model
2 lr = LogisticRegression()

1 def predict_80_20(df, target):
2     split = StratifiedShuffleSplit(n_splits = 1, test_size=0.2, random_state=42)
3     for train_index, test_index in split.split(df, df[target]):
4         train = df.loc[train_index]
5         test = df.loc[test_index]
6
7         x = train.drop([target], axis=1)
8         y = train[target]
9
10        a = test.drop([target], axis=1)
11        b = test[target]
12
13        lr.fit(x,y)
14
15        predicted = lr.predict(a)
16        accuracy = accuracy_score(b, predicted)
17        c_m = confusion_matrix(b, predicted)
18
19        return print('Accuracy score: ', accuracy, '\n', 'Confusion matrix: \n', c_m)
```

Prepared as function to be called later

Key Numbers



```
2]: 1 %run predict_CHD.ipynb # calling the LogisticRegression
```

```
3]: 1 predict_80_20(heart1, 'TenYearCHD')
```

```
Accuracy score: 0.8524788391777509  
Confusion matrix:  
[[697  4]  
 [118  8]]
```

Drop education column

```
1 heart_ed = heart1.drop(['education'], axis = 1)
```

```
1 predict_80_20(heart_ed, 'TenYearCHD')
```

```
Accuracy score: 0.8512696493349455  
Confusion matrix:  
[[697  4]  
 [119  7]]
```

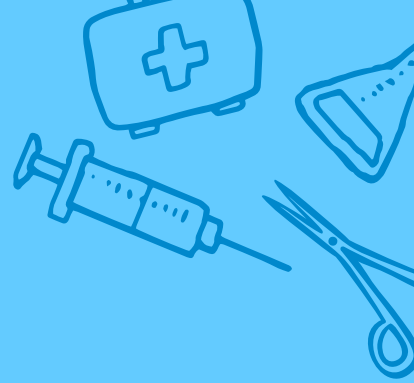




Future work

- Enhancing the model with larger dataset.
- Using some columns to make model for another disease





Thanks

