

Team members

- Bishal Sainju
- Supratik Chanda
- Victor Lee

Introduction

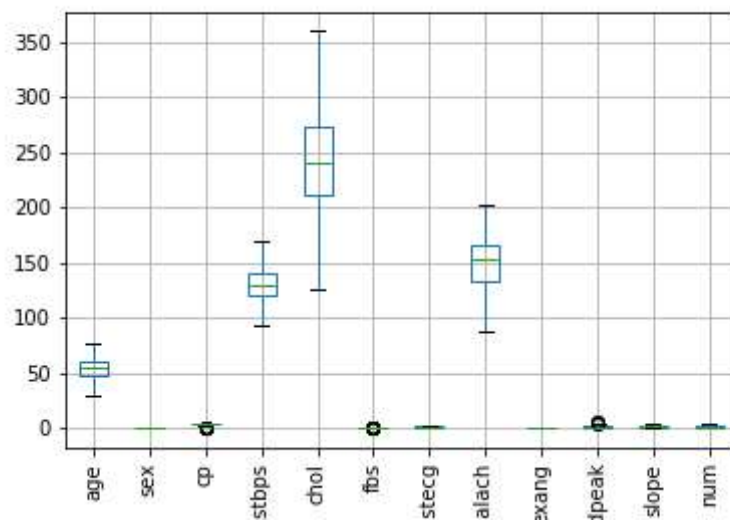
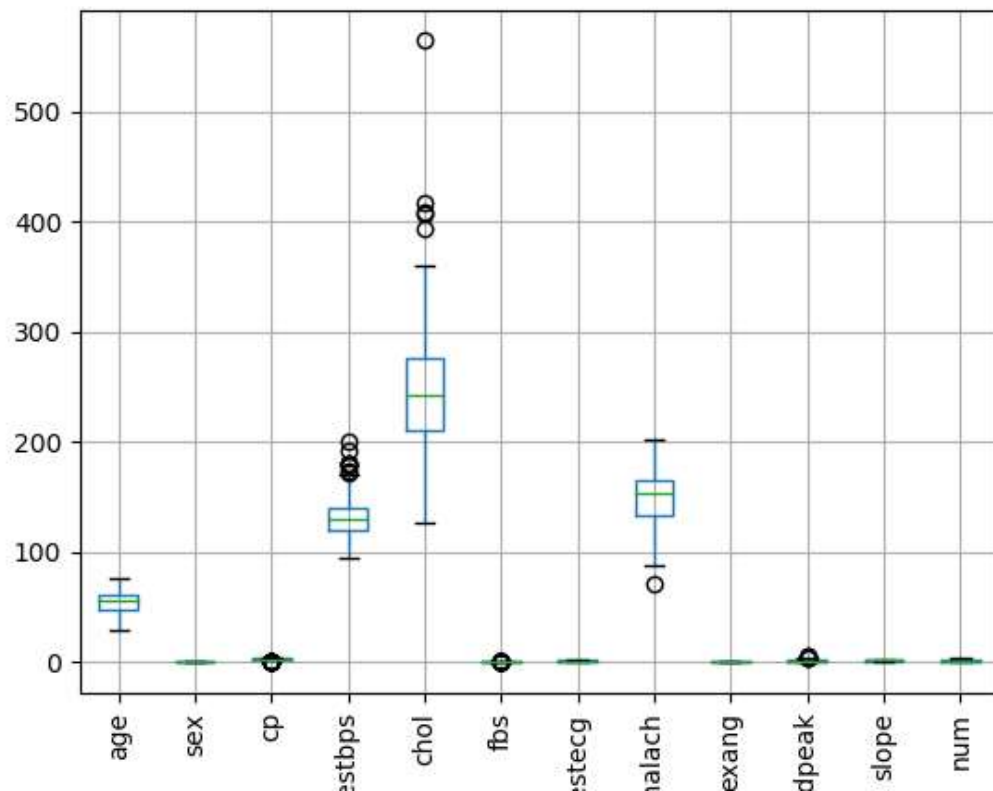
We tried to predict, if a person has a heart disease or not. We expect to create the model that would generalize well in the real world scenario.

Dataset

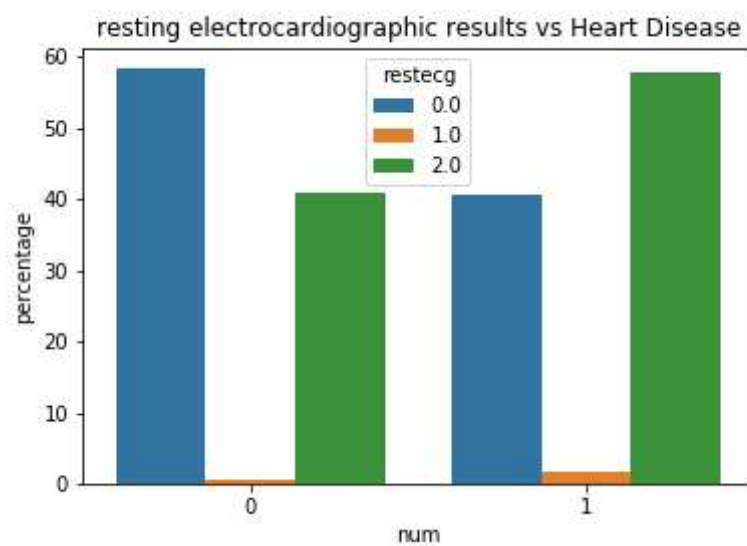
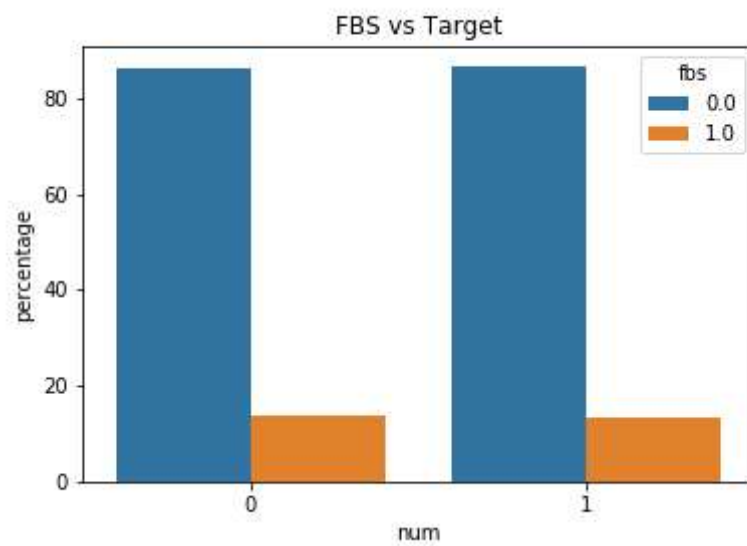
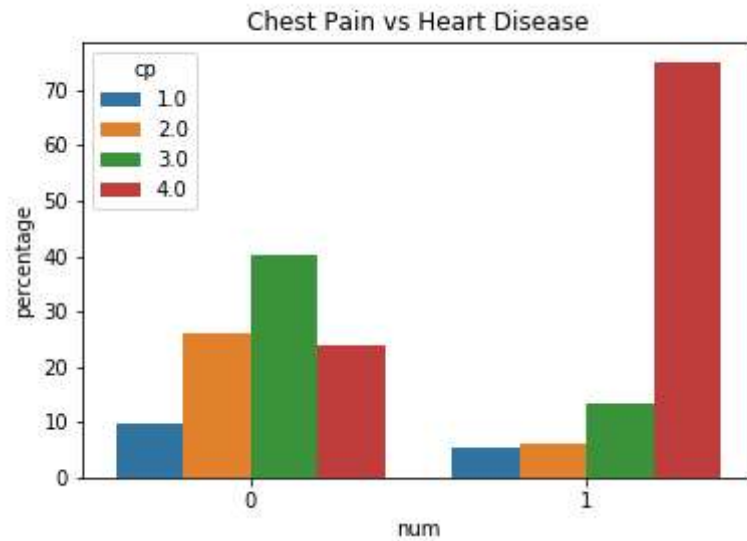
The dataset that we used is Cleveland heart disease dataset. First we cleaned the data, deleted any record having values like '?'. Then we checked for outliers and removed them. Before feeding our data to the training model, we normalized it using the standard scaler.

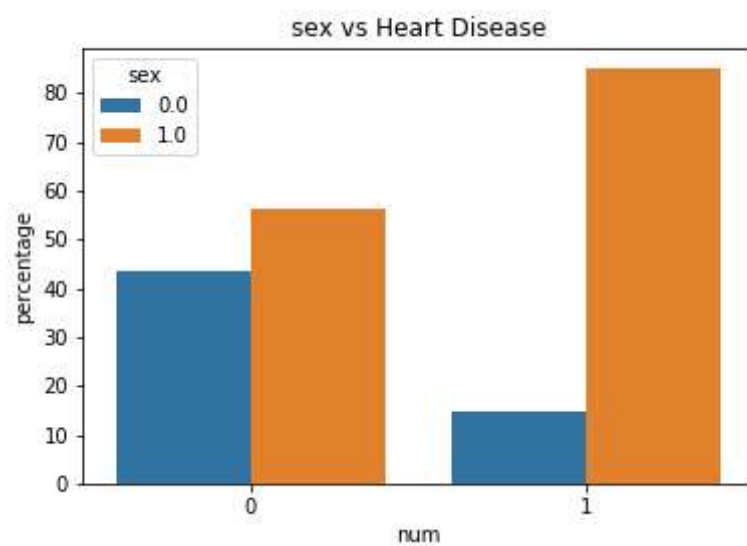
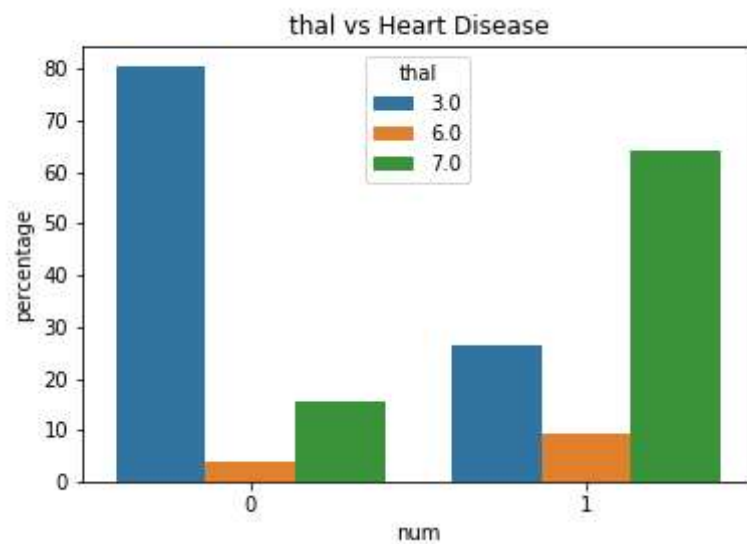
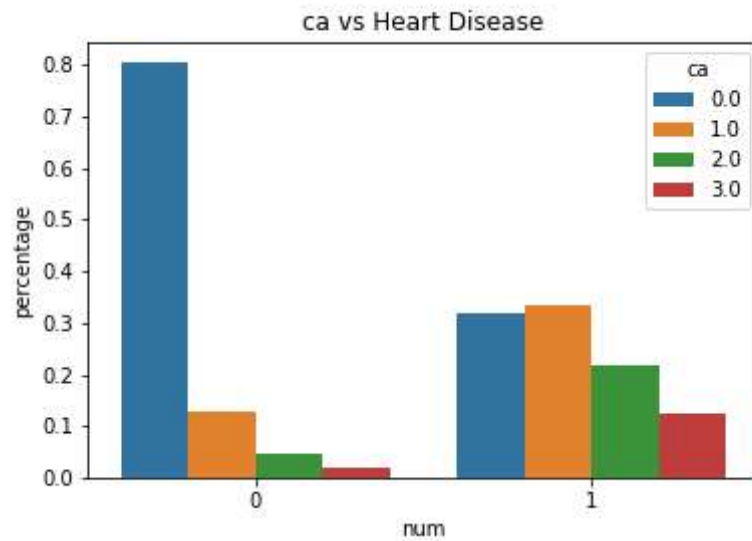
Analysis Technique

1. Data Cleaning: First, we dropped any abnormal values and then plotted the boxplot to figure out if there was any outliers in the data. There was some harsh outliers in these attributes: ['trestbps', 'chol', 'thalach'] So, we removed those outliers, by removing all the data that was not in the Interquartile range in those attributes.



2. Data Analysis: Then, we plotted the bar graph to get some idea of the various attributes.





By, this analysis we can see that, cp_4.0 must be important over other cp values.

Similarly, we can see that, 'fbs' attribute is not much important in determining if a person has a heart disease or not.

Similarly, 'restecg_0' and 'restecg_2.0' are important predictors whereas restecg_1.0 is not.

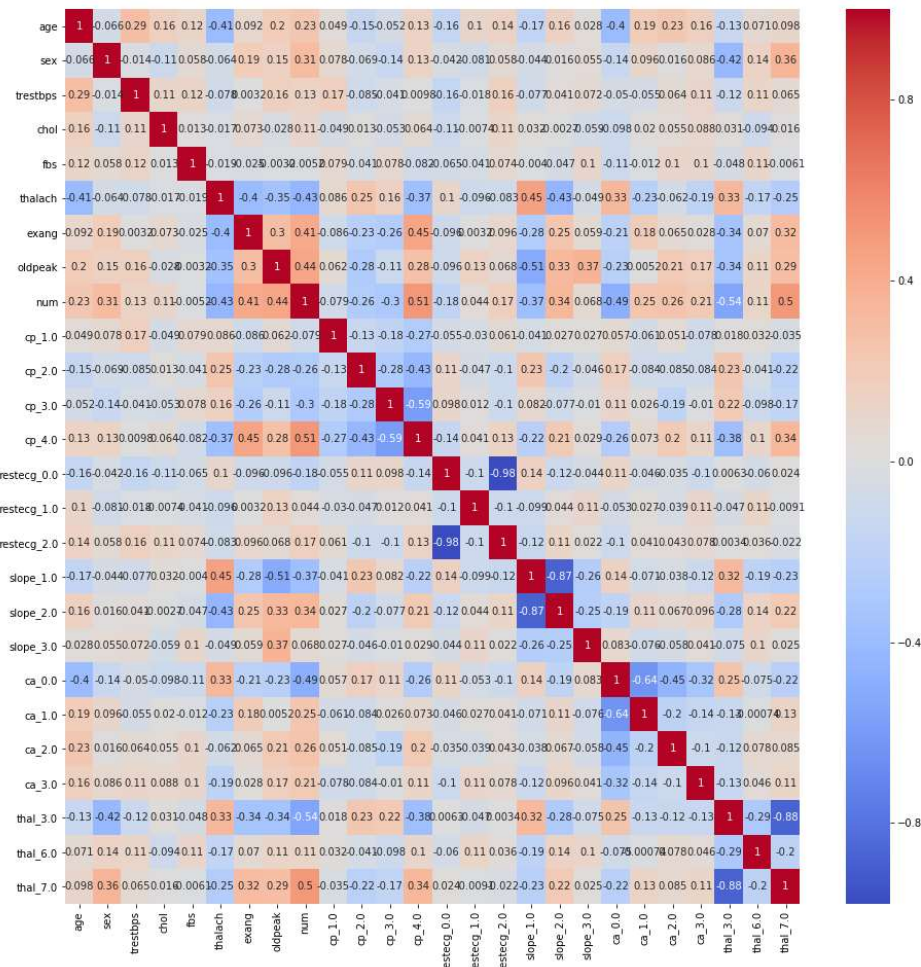
Similarly, 'ca_0' is a much more important characteristics than other values.

Similarly, thal_3.0 and thal_7.0 are important characteristics.

- Since, we have many categorical attributes, we encoded them using get_dummies.

We applied get_dummies on these attributes: ['cp', 'restecg', 'slope', 'ca', 'thal']

- Next, we wanted to know which attributes were the most important ones. So, we first standardized the data and used correlation matrix for that.



We can see that attributes like ['thal_3.0', 'cp_4.0', 'thal_7.0', 'ca_0.0', 'oldpeak', 'thalach'] are more important because they have higher correlation with the 'num'.

So, we used SelectKBest using f_regression to figure out important features and ordered them in the order of their importance.

['thal_3.0', 'cp_4.0', 'thal_7.0', 'ca_0.0', 'oldpeak', 'thalach', 'exang', 'slope_1.0', 'slope_2.0', 'sex', 'cp_3.0', 'ca_2.0', 'cp_2.0', 'ca_1.0', 'age', 'ca_3.0', 'restecg_0.0', 'restecg_2.0', 'trestbps', 'thal_6.0', 'chol', 'cp_1.0', 'slope_3.0', 'restecg_1.0', 'fbs']

We got what was expected.

5. Next, we built the custom KNN model and cross-validation method, and using for loops determined which model had the best f1-score.
6. Later, we also used the GridSearchCV method to find out the best parameters for the best estimates.

Results

1. From the custom KNN model that we built, we obtained a maximum f1-score of 86.32, when we select 19 best features and 6 nearest neighbors as a parameter to our KNN model.

The 19 best features that it selected were:

N = 19: ['thal_3.0', 'cp_4.0', 'thal_7.0', 'ca_0.0', 'oldpeak', 'thalach', 'exang', 'slope_1.0', 'slope_2.0', 'sex', 'cp_3.0', 'ca_2.0', 'cp_2.0', 'ca_1.0', 'age', 'ca_3.0', 'restecg_0.0', 'restecg_2.0', 'trestbps']

The scores are as follows:

accuracy: 0.883, f1-score: 0.8632, precision: 0.872, recall: 0.8568

2. However, we wanted to implement the pipeline and grid search library to do our task much efficiently. So we run our algorithm and searched using gridsearchCV, and found out that the best parameters were the following:

```
{'classify__algorithm': 'auto', 'classify__n_neighbors': 11, 'classify__weights': 'uniform',
'reduce_dim': SelectKBest(k=18, score_func=<function f_regression>), 'reduce_dim__k': 18,
'reduce_dim__score_func': }
```

The best estimate of the f1-score, we obtained was 0.8493

accuracy: 86.88, precision: 88.01, recall: 82.75

This shows that our custom knn had best f1-score than the built in, it might be because of the difference in the implementaion of cross-validation score and knn classifier.

The built-in cross validation uses stratifying parameter, whereas we have not implemented that. We found from the gridsearch that best parameter was not weighted by distance, so we did not implement weightier distance in out implementation of KNN.

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 from sklearn.model_selection import train_test_split
        6
        7
        8
        9 hd = pd.read_csv('cleveland.csv')
       10 hd.dropna(inplace=True)
       11 hd = hd[hd.ca != '?']
       12 hd = hd[hd.thal != '?']
       13 print(hd.shape)
       14 hd.boxplot()
       15 plt.xticks(rotation = 90)
       16 plt.savefig('outlier.png')
       17 plt.show()
       18
```

(297, 14)

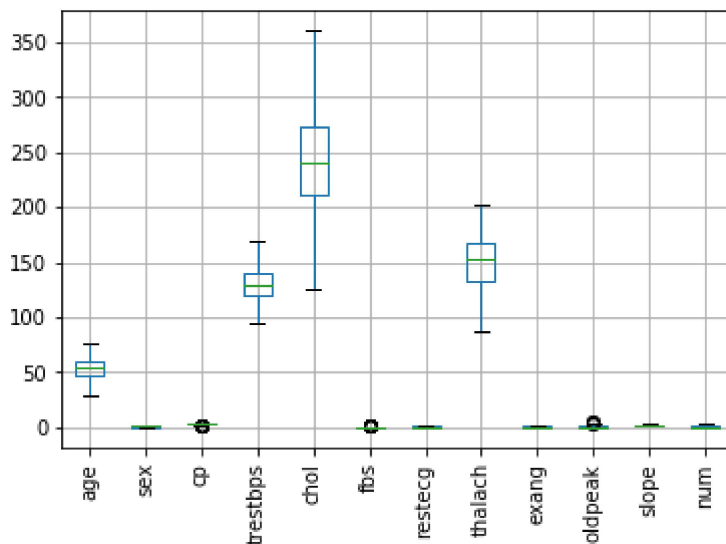
<Figure size 640x480 with 1 Axes>

```

In [2]: 1 def remove_outliers(hd, cols):
2         for col in cols:
3             Q3 = hd[col].quantile(.75)
4             Q1 = hd[col].quantile(.25)
5             IQR = Q3 - Q1
6             hd = hd[~((hd[col] < Q1 - 1.5 * IQR) | (hd[col] > Q3 + 1.5 * IQR))]
7         return hd
8
9     hd = remove_outliers(hd, ['trestbps', 'chol', 'thalach'])
10    print(hd.shape)
11    hd.boxplot()
12    plt.xticks(rotation = 90)
13    plt.savefig('outlier_rem.png')
14    plt.show()

```

(282, 14)



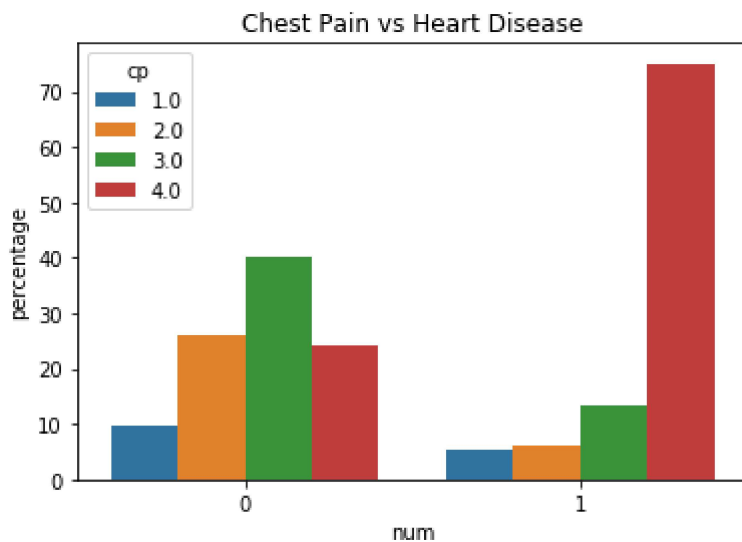
```

In [3]: 1 #num: 1, 2, 3 = 1
2     hd['num'][hd['num'] != 0] = 1
3     # hd['num']

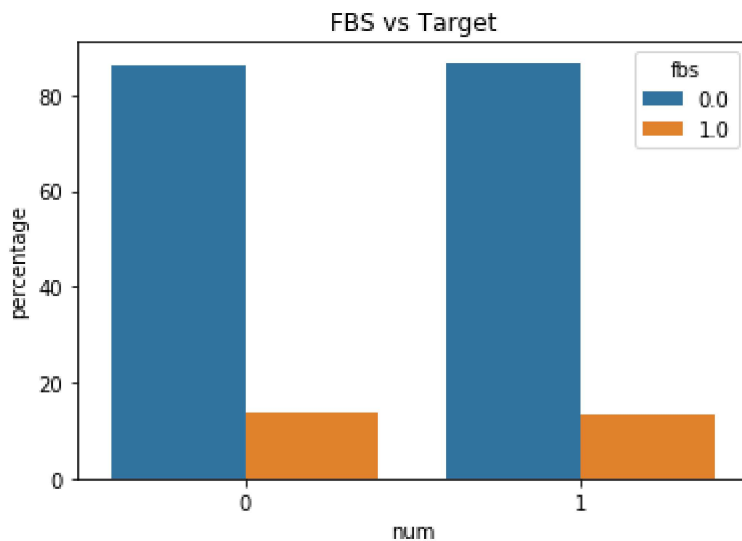
```



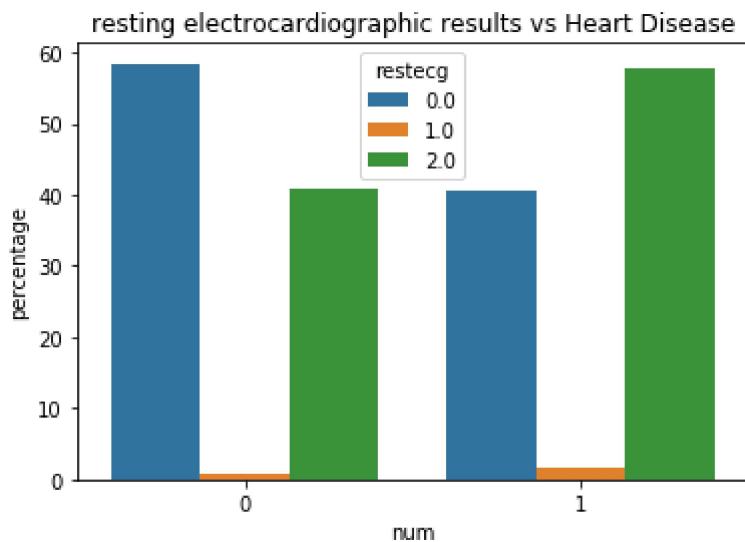
```
In [4]: 1 temp = (hd.groupby(['num']))['cp'].value_counts(normalize=True)\
2         .mul(100).reset_index(name = "percentage")
3         sns.barplot(x = "num", y = "percentage", hue = "cp", data = temp)\
4         .set_title("Chest Pain vs Heart Disease")
5         plt.savefig('1.png')
```



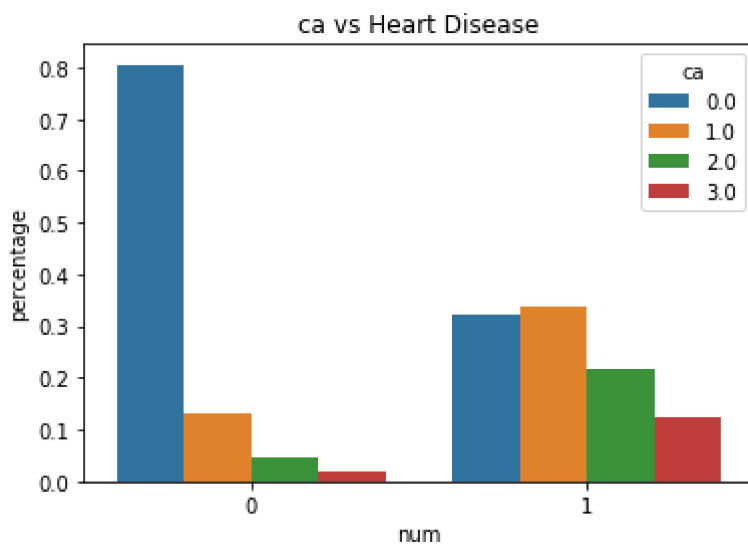
```
In [5]: 1 temp = (hd.groupby(['num']))['fbs'].value_counts(normalize=True)\
2         .mul(100).reset_index(name = "percentage")
3         sns.barplot(x = "num", y = "percentage", hue = "fbs", data = temp).set_title(
4         plt.savefig('2.png'))
```



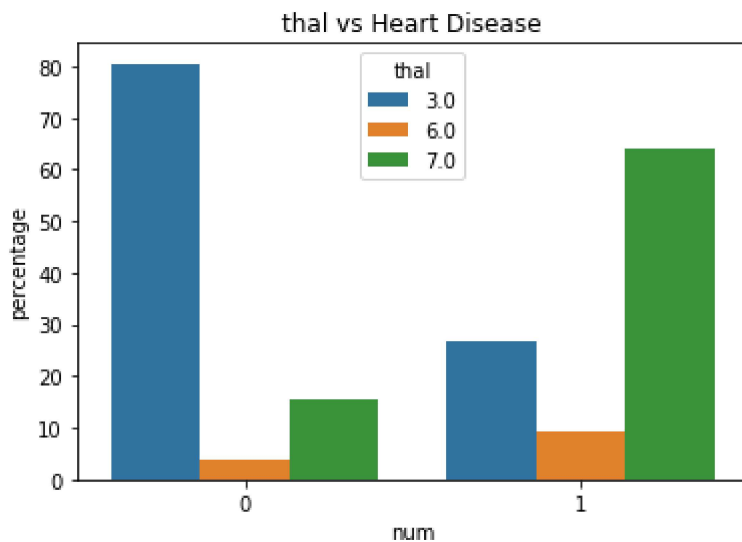
```
In [6]: 1 temp = (hd.groupby(['num']))['restecg'].value_counts(normalize=True)\
2         .mul(100).reset_index(name = "percentage")
3         sns.barplot(x = "num", y = "percentage", hue = "restecg", data = temp)\
4         .set_title("resting electrocardiographic results vs Heart Disease")
5         plt.savefig('3.png')
```



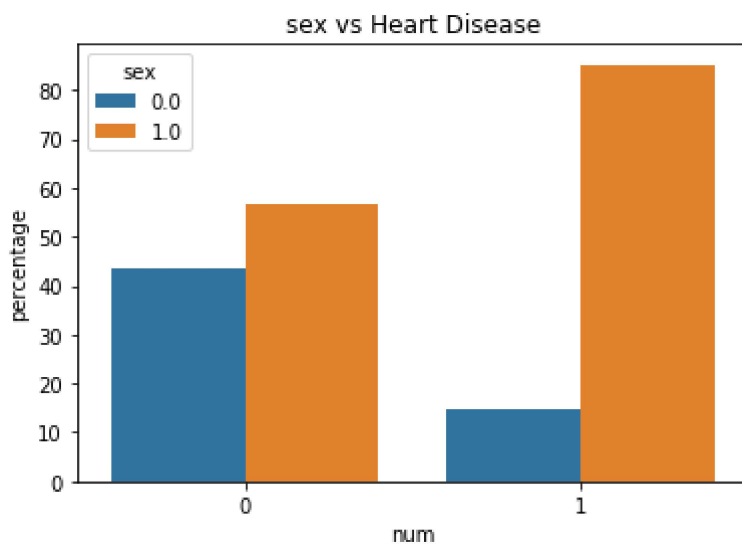
```
In [7]: 1 temp = (hd.groupby(['num']))['ca'].value_counts(normalize=True)\
2         .reset_index(name = "percentage")
3         sns.barplot(x = "num", y = "percentage", hue = "ca", data = temp)\
4         .set_title("ca vs Heart Disease")
5         plt.savefig('4.png')
```



```
In [8]: 1 temp = (hd.groupby(['num']))['thal'].value_counts(normalize=True)\
2         .mul(100).reset_index(name = "percentage")
3         sns.barplot(x = "num", y = "percentage", hue = "thal", data = temp)\
4         .set_title("thal vs Heart Disease")
5         plt.savefig('5.png')
```



```
In [9]: 1 temp = (hd.groupby(['num']))['sex'].value_counts(normalize=True)\
2         .mul(100).reset_index(name = "percentage")
3         sns.barplot(x = "num", y = "percentage", hue = "sex", data = temp)\
4         .set_title("sex vs Heart Disease")
5         plt.savefig('6.png')
```



```
In [10]: 1 #discretize the data (get_dummies)
2 hd = pd.get_dummies(hd, columns=['cp', 'restecg', 'slope', 'ca', 'thal'], dr
3 hd.head()
```

Out[10]:

	age	sex	trestbps	chol	fbs	thalach	exang	oldpeak	num	cp_1.0	...	slope_1.0	slope_2.0
0	63.0	1.0	145.0	233.0	1.0	150.0	0.0	2.3	0	1	...	0	0
1	67.0	1.0	160.0	286.0	0.0	108.0	1.0	1.5	1	0	...	0	1
2	67.0	1.0	120.0	229.0	0.0	129.0	1.0	2.6	1	0	...	0	1
3	37.0	1.0	130.0	250.0	0.0	187.0	0.0	3.5	0	0	...	0	0
4	41.0	0.0	130.0	204.0	0.0	172.0	0.0	1.4	0	0	...	1	0

5 rows × 26 columns

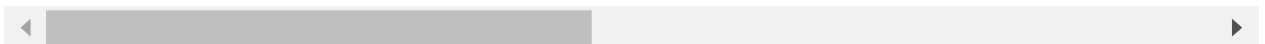


```
In [11]: 1 #standardize the data
2 from sklearn import preprocessing
3 col_to_norm = hd.columns.tolist()
4 col_to_norm.remove('num')
5 hd[col_to_norm] = preprocessing.StandardScaler().fit_transform(hd[col_to_norm])
6 hd.head()
```

Out[11]:

	age	sex	trestbps	chol	fbs	thalach	exang	oldpeak	num	cp
0	0.969368	0.662401	0.956135	-0.226835	2.533980	0.012802	-0.684653	1.123914	0	3.43
1	1.410275	0.662401	1.916813	0.947896	-0.394636	-1.836246	1.460593	0.420181	1	-0.25
2	1.410275	0.662401	-0.644994	-0.315494	-0.394636	-0.911722	1.460593	1.387814	1	-0.25
3	-1.896523	0.662401	-0.004542	0.149966	-0.394636	1.641724	-0.684653	2.179514	0	-0.25
4	-1.455616	-1.509659	-0.004542	-0.869612	-0.394636	0.981350	-0.684653	0.332214	0	-0.25

5 rows × 26 columns

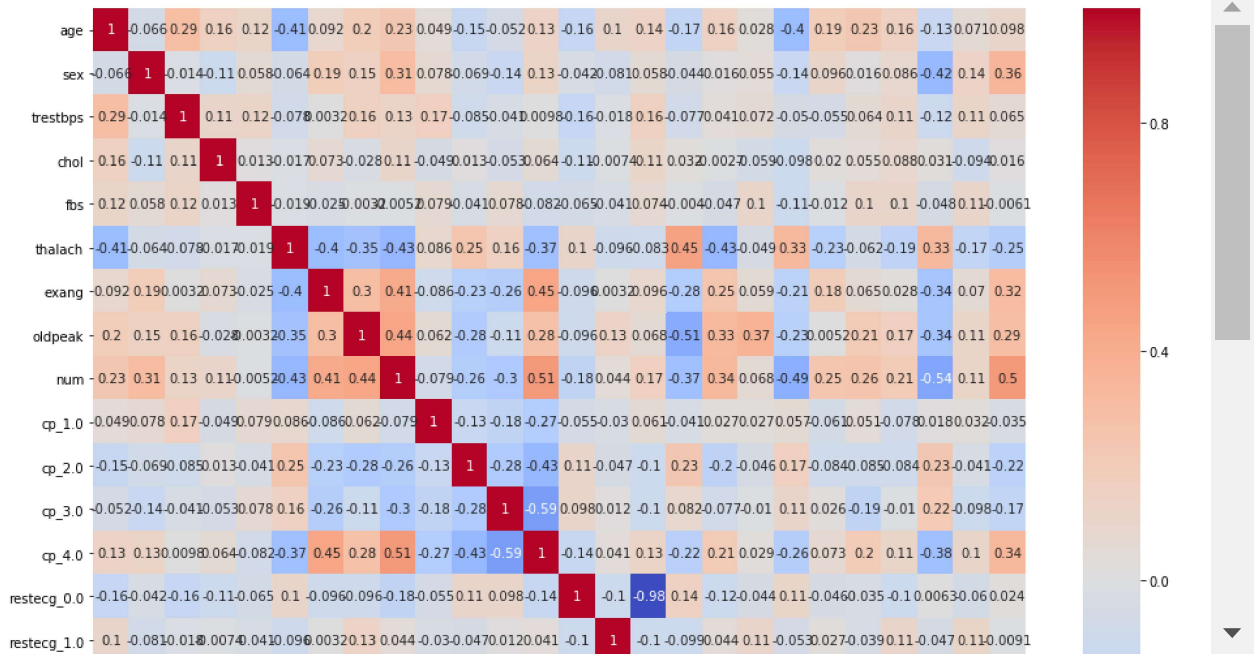


In [12]:

```

1 plt.figure(figsize=(16,16))
2 sns.heatmap(hd.corr(),annot=True,cmap='coolwarm')
3 plt.savefig('correlation.png')
4 plt.show()

```



In [13]:

```
1 print(hd.corr()['num'])
```

```

age          0.230561
sex          0.309960
trestbps     0.131340
chol         0.105463
fbs          -0.005178
thalach      -0.433597
exang        0.414825
oldpeak      0.438209
num          1.000000
cp_1.0       -0.079296
cp_2.0       -0.261296
cp_3.0       -0.299103
cp_4.0       0.508388
restecg_0.0  -0.177411
restecg_1.0  0.044314
restecg_2.0  0.168382
slope_1.0    -0.374651
slope_2.0    0.341924
slope_3.0    0.067508
ca_0.0       -0.489967
ca_1.0       0.246310
ca_2.0       0.261680
ca_3.0       0.209577
thal_3.0     -0.541220
thal_6.0     0.111589
thal_7.0     0.498312
Name: num, dtype: float64

```

```
In [14]: 1 #Which features are the most important (arrange them in descending order)
2 from sklearn.feature_selection import SelectKBest, f_regression
3
4 X = hd[col_to_norm].values.astype(np.float)
5 y = hd['num'].values.astype(np.float)
6 k_best = SelectKBest(f_regression, k=len(col_to_norm)).fit(X, y)
7 score = np.array(k_best.scores_)
8 rank = score.argsort()[-len(col_to_norm):][::-1]
9 features = []
10 for i in rank:
11     features.append(col_to_norm[i])
12 print(features)
```

```
['thal_3.0', 'cp_4.0', 'thal_7.0', 'ca_0.0', 'oldpeak', 'thalach', 'exang', 'slope_1.0', 'slope_2.0', 'sex', 'cp_3.0', 'ca_2.0', 'cp_2.0', 'ca_1.0', 'age', 'ca_3.0', 'restecg_0.0', 'restecg_2.0', 'trestbps', 'thal_6.0', 'chol', 'cp_1.0', 'slope_3.0', 'restecg_1.0', 'fbs']
```

In [15]:

```

1  #building custom model
2  import numpy as np
3  from sklearn.metrics import f1_score, accuracy_score, precision_score, recall
4      classification_report, confusion_matrix
5  from sklearn.neighbors import NearestNeighbors
6
7  def classification_report_cm(y_true, y_pred):
8      print (classification_report(y_true, y_pred))
9      print (confusion_matrix(y_true, y_pred))
10
11 def accuracy(y_true, y_pred):
12     return accuracy_score(y_true, y_pred)
13
14 def precision(y_true, y_pred):
15     return precision_score(y_true, y_pred)
16
17 def recall(y_true, y_pred):
18     return recall_score(y_true, y_pred)
19
20 def f1(y_true, y_pred):
21     return f1_score(y_true, y_pred)
22
23 def knn(X_train, y_train, X_test, k):
24     neigh = NearestNeighbors(n_neighbors=k).fit(X_train, y_train)
25     y_pred_list = neigh.kneighbors(X_test, return_distance=False)
26     y_pred = [[y_train[i] for i in indices] for indices in y_pred_list]
27     # print(y_pred)
28     y_pred = [max(y, key=y.count) for y in y_pred]
29     return np.asarray(y_pred)
30
31 def chunkIt(seq, num):
32     avg = len(seq) / float(num)
33     out = []
34     last = 0.0
35     while last < len(seq):
36         out.append(seq[int(last):int(last + avg)])
37         last += avg
38     return out
39
40 def get_score(data, k, n):
41     # print(data.shape)
42     data_list = chunkIt(data, n)
43     f1_list = []
44     acc_list = []
45     prec_list = []
46     rec_list = []
47     for i in range(n):
48         data_split = data_list.copy()
49         test = np.asarray(data_split.pop(i))
50         train = np.asarray([j for k in data_split for j in k])
51         X_train, y_train, X_test, y_test = train[:, :-1], train[:, -1], test[:, :-1], test[:, -1]
52         y_pred = knn(X_train, y_train, X_test, k)
53         print('Cross-Val: {}'.format(i))
54         classification_report_cm(y_test, y_pred)
55         f1_score = f1(y_test, y_pred)
56         acc = accuracy(y_test, y_pred)

```

```
57     prec = precision(y_test, y_pred)
58     rec = recall(y_test, y_pred)
59     print('f1: {}, acc: {}'.format(f1_score, acc))
60     print("")
61     f1_list.append(f1_score)
62     acc_list.append(acc)
63     prec_list.append(prec)
64     rec_list.append(rec)
65     return np.asarray(f1_list).mean().round(4), np.asarray(acc_list).mean().r
66         np.asarray(prec_list).mean().round(4), np.asarray(rec_list).mean()
```



```

In [16]: 1 import json
          2
          3 #shuffle data
          4 hd2 = hd.sample(frac=1, random_state=5).reset_index(drop=True)
          5 # K = range(5, 15)
          6 # num_attr = range(15, 21)
          7 K = [6]
          8 num_attr = [19]
          9 score = {}
         10 for n in num_attr:
         11     score[n] = {}
         12     feat = features[:n]
         13     print('N = {}: {}'.format(n, feat))
         14     feat.append('num')
         15     hd1 = hd2[feat]
         16     hd_data = hd1.values.astype(float)
         17     for k in K:
         18         score[n][k] = {}
         19         score[n][k]['f1'], score[n][k]['acc'], score[n][k]['prec'], score[n][
         20 print(json.dumps(score, indent=4, sort_keys=True))
         21

```

N = 19: ['thal_3.0', 'cp_4.0', 'thal_7.0', 'ca_0.0', 'oldpeak', 'thalach', 'exang', 'slope_1.0', 'slope_2.0', 'sex', 'cp_3.0', 'ca_2.0', 'cp_2.0', 'ca_1.0', 'age', 'ca_3.0', 'restecg_0.0', 'restecg_2.0', 'trestbps']

Cross-Val: 0

	precision	recall	f1-score	support
0.0	0.88	0.86	0.87	35
1.0	0.77	0.81	0.79	21
avg / total	0.84	0.84	0.84	56

[[30 5]

[4 17]]

f1: 0.7906976744186046, acc: 0.8392857142857143

Cross-Val: 1

	precision	recall	f1-score	support
0.0	0.93	0.96	0.95	28
1.0	0.96	0.93	0.95	28
avg / total	0.95	0.95	0.95	56

[[27 1]

[2 26]]

f1: 0.9454545454545454, acc: 0.9464285714285714

Cross-Val: 2

	precision	recall	f1-score	support
0.0	0.83	0.92	0.87	26
1.0	0.93	0.84	0.88	31

avg / total	0.88	0.88	0.88	57
-------------	------	------	------	----

[[24 2]

[5 26]]

f1: 0.8813559322033899, acc: 0.8771929824561403

Cross-Val: 3

	precision	recall	f1-score	support
0.0	0.91	0.86	0.89	36
1.0	0.77	0.85	0.81	20

avg / total	0.86	0.86	0.86	56
-------------	------	------	------	----

[[31 5]

[3 17]]

f1: 0.8095238095238095, acc: 0.8571428571428571

Cross-Val: 4

	precision	recall	f1-score	support
0.0	0.87	0.93	0.90	29
1.0	0.92	0.86	0.89	28

avg / total	0.90	0.89	0.89	57
-------------	------	------	------	----

[[27 2]

[4 24]]

f1: 0.8888888888888889, acc: 0.8947368421052632

```
{
  "19": {
    "6": {
      "acc": 0.883,
      "f1": 0.8632,
      "prec": 0.872,
      "rec": 0.8568
    }
  }
}
```

In [17]:

```
1 X = hd[col_to_norm].values.astype(np.float)
2 y = hd['num'].values.astype(np.float)
```

```

In [18]: 1 #Using built in KNN, pipeline and gridsearch
2 from sklearn.decomposition import PCA
3 from sklearn.feature_selection import SelectKBest, chi2, f_regression, f_classif
4         mutual_info_regression, RFE
5 from sklearn.model_selection import GridSearchCV
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.pipeline import Pipeline
8
9 pipe = Pipeline([
10     ('reduce_dim', None),
11     ('classify', KNeighborsClassifier())
12 ])
13
14 # N_FEATURES_OPTIONS = range(10, 25)
15 N_FEATURES_OPTIONS = [18]
16 # K_OPTIONS = range(5, 15)
17 K_OPTIONS = [11]
18 SCORE_FUNC = [f_regression]
19 ALG = ['auto']
20 WT = ['uniform']
21 param_grid = [
22     # {
23     #     'reduce_dim': [PCA()],
24     #     'reduce_dim__n_components': N_FEATURES_OPTIONS,
25     #     'classify__n_neighbors': K_OPTIONS,
26     #     'classify__weights': WT,
27     #     'classify__algorithm': ALG
28     # },
29     {
30         'reduce_dim': [SelectKBest()],
31         'reduce_dim__k': N_FEATURES_OPTIONS,
32         'reduce_dim__score_func': SCORE_FUNC,
33         'classify__n_neighbors': K_OPTIONS,
34         'classify__weights': WT,
35         'classify__algorithm': ALG
36     }
37 ]
38 grid = GridSearchCV(pipe, cv=10, param_grid=param_grid, scoring='recall')
39 grid.fit(X, y)
40 print(grid.best_params_)
41 print(grid.best_score_)

```

```

{'classify__algorithm': 'auto', 'classify__n_neighbors': 11, 'classify__weights': 'uniform', 'reduce_dim': SelectKBest(k=18, score_func=<function f_regression at 0x000001501C4FFEA0>), 'reduce_dim__k': 18, 'reduce_dim__score_func': <function f_regression at 0x000001501C4FFEA0>}
0.8274686306601201

```

In []:

1

In []:

1

In []:

1

