

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import seaborn as sns
        4 import matplotlib.pyplot as plt
```

```
In [2]: 1 df = pd.read_csv('data.csv')
        2 print("The size of initial dataset: ", df.shape)
```

The size of initial dataset: (569, 33)

```
In [3]: 1 df.isna().sum()
```

```
Out[3]: id                                0
        diagnosis                          0
        radius_mean                        0
        texture_mean                       0
        perimeter_mean                     0
        area_mean                          0
        smoothness_mean                    0
        compactness_mean                   0
        concavity_mean                     0
        concave points_mean                 0
        symmetry_mean                      0
        fractal_dimension_mean              0
        radius_se                           0
        texture_se                          0
        perimeter_se                        0
        area_se                             0
        smoothness_se                       0
        compactness_se                      0
        concavity_se                        0
        concave points_se                   0
        symmetry_se                         0
        fractal_dimension_se                0
        radius_worst                       0
        texture_worst                       0
        perimeter_worst                     0
        area_worst                          0
        smoothness_worst                    0
        compactness_worst                   0
        concavity_worst                     0
        concave points_worst                0
        symmetry_worst                      0
        fractal_dimension_worst              0
        Unnamed: 32                          569
        dtype: int64
```

```
In [4]: 1 df_1 = df.dropna(axis=1)
        2 print("The size of final dataset: ", df_1.shape)
```

The size of final dataset: (569, 32)

```
In [5]: 1 y = df_1.diagnosis #feature labels
        2 na_list = ['id','diagnosis']
        3 X = df_1.drop(na_list,axis = 1 ) #data matrix
        4 X.head()
```

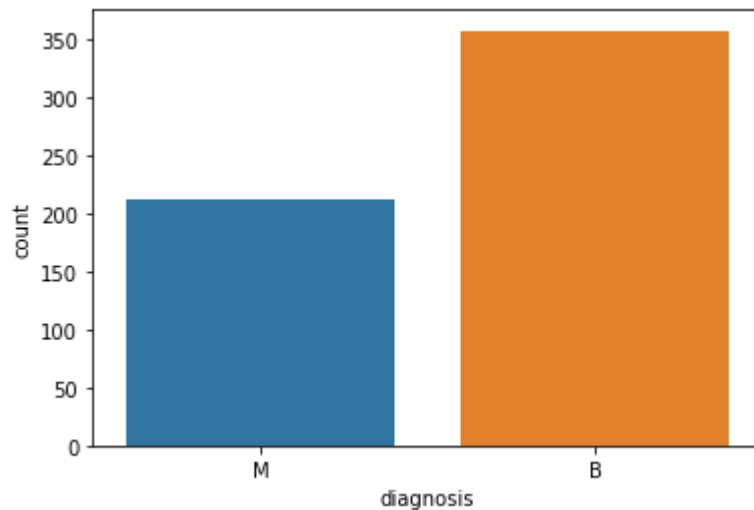
Out[5]:

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean
0	17.99	10.38	122.80	1001.0	0.11840	0.27760
1	20.57	17.77	132.90	1326.0	0.08474	0.07864
2	19.69	21.25	130.00	1203.0	0.10960	0.15990
3	11.42	20.38	77.58	386.1	0.14250	0.28390
4	20.29	14.34	135.10	1297.0	0.10030	0.13280

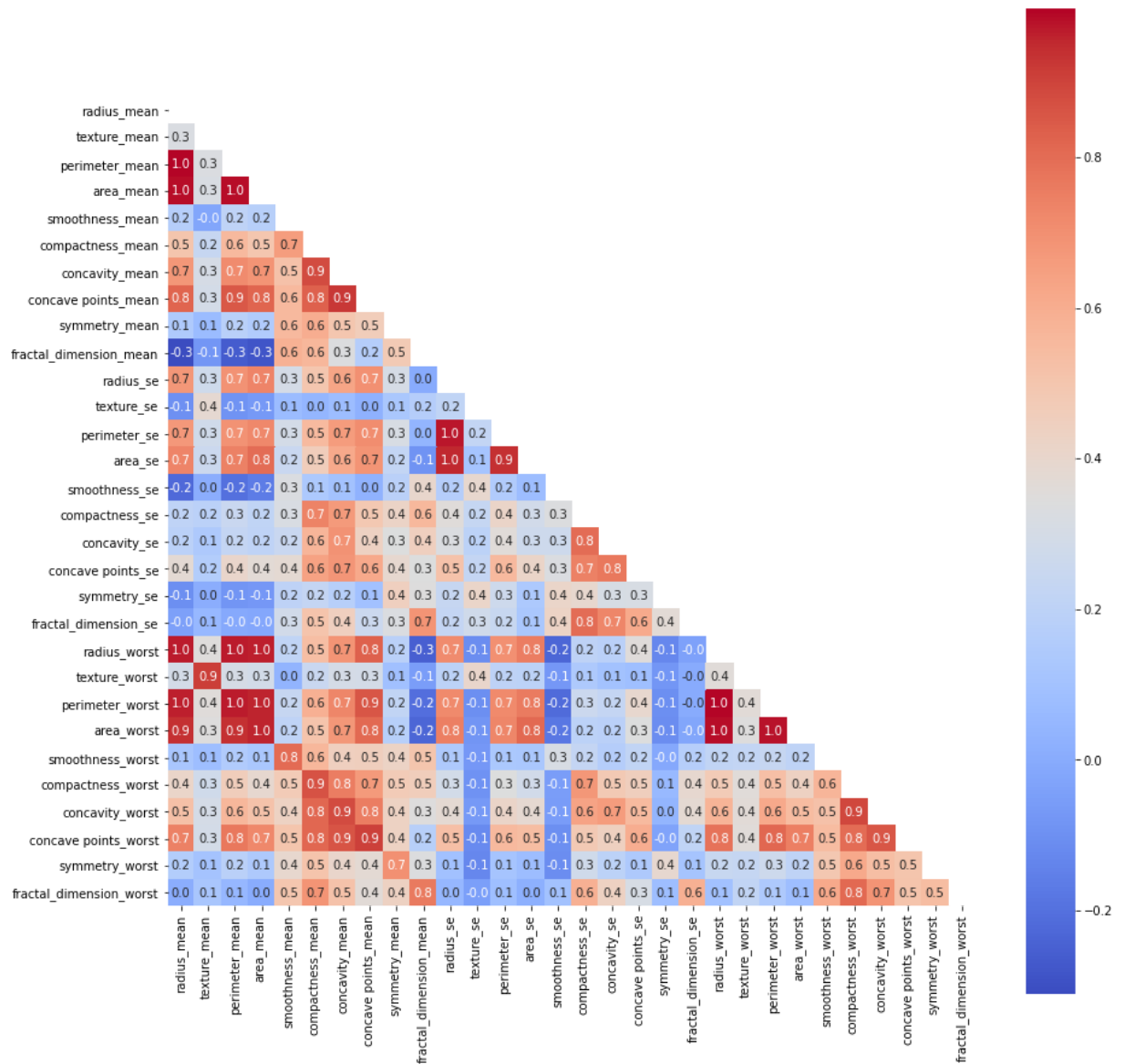
5 rows × 30 columns

```
In [6]: 1 #histograms showing the count of M and B
        2 sns.countplot(x = df_1['diagnosis'], label="Count")
```

Out[6]: <AxesSubplot:xlabel='diagnosis', ylabel='count'>



```
In [7]: 1 plt.figure(figsize=(15,15))
2 upper_traing_matrix = np.triu(X.corr())
3 sns.heatmap(X.corr(), annot=True, square=True, fmt= '.1f', cmap='coolwarm',
4 plt.show())
```



```
In [8]: 1 # features which have correlation greater than 0.95 can be dropped
2 correlation_matrix = X.corr().abs()
3 mask = np.triu(np.ones(correlation_matrix.shape), k=1)
4 upper_triangular_matrix = correlation_matrix.where(mask.astype(np.bool))
5 na_features = [c for c in upper_triangular_matrix.columns if any(upper_trian
6 print("The features which can be dropped because of high correlation are")
7 print(na_features)]
```

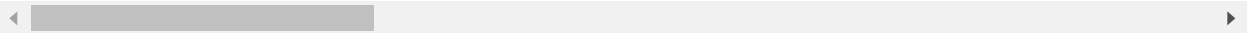
The features which can be dropped because of high correlation are
 ['perimeter_mean', 'area_mean', 'perimeter_se', 'area_se', 'radius_worst', 'perimeter_worst', 'area_worst']

```
In [9]: 1 # feature selection
2 X = X.drop(X[na_features], axis=1)
3 X.head()
```

Out[9]:

	radius_mean	texture_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean
0	17.99	10.38	0.11840	0.27760	0.3001	0.147
1	20.57	17.77	0.08474	0.07864	0.0869	0.070
2	19.69	21.25	0.10960	0.15990	0.1974	0.127
3	11.42	20.38	0.14250	0.28390	0.2414	0.105
4	20.29	14.34	0.10030	0.13280	0.1980	0.104

5 rows × 23 columns



```
In [10]: 1 print("The original values for y are: ")
          2 print(y)
          3 from sklearn.preprocessing import LabelEncoder
          4 Y_new = LabelEncoder()
          5 y= Y_new.fit_transform(y)
          6 print("The new values for y are: ")
          7 print(y)
```

The original values for y are:

0	M
1	M
2	M
3	M
4	M
	...
564	M
565	M
566	M
567	M
568	B

Name: diagnosis, Length: 569, dtype: object

The new values for y are:

[illegible]

```
In [11]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.40,
3 print("The shape of X_train is ", X_train.shape)
4 print("The shape of y_train is ", y_train.shape)
5 print("The shape of X_test is ", X_test.shape)
6 print("The shape of y_test is ", y_test.shape)
```

```
The shape of X_train is (341, 23)
The shape of y_train is (341,)
The shape of X_test is (228, 23)
The shape of y_test is (228,)
```

```
In [12]: 1 from sklearn.preprocessing import RobustScaler
          2 scaler = RobustScaler()
          3 X_train = scaler.fit_transform(X_train)
          4 X_test = scaler.transform(X_test)
```

```
In [13]: 1 from sklearn.metrics import accuracy_score
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.naive_bayes import GaussianNB
5 from sklearn.model_selection import cross_val_score
6 #Decision Tree
7 clf1 = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
8 scores_1 = cross_val_score(clf1, X, y, cv=5)
9 clf1.fit(X_train, y_train)
10 print('Training Accuracy for decision tree is ', clf1.score(X_train, y_train))
11 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_1.mean(), scores_1.std()))
12 #Random Forest
13 clf2 = RandomForestClassifier(n_estimators = 8, criterion = 'entropy', random_state = 0)
14 scores_2 = cross_val_score(clf2, X, y, cv=5)
15 clf2.fit(X_train, y_train)
16 print('Training Accuracy for random forest is ', clf2.score(X_train, y_train))
17 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_2.mean(), scores_2.std()))
18 #Naive Bayes
19 clf3 = GaussianNB()
20 scores_3 = cross_val_score(clf3, X, y, cv=5)
21 clf3.fit(X_train, y_train)
22 print('Training Accuracy for naive bayes classifier is ', clf3.score(X_train, y_train))
23 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_3.mean(), scores_3.std()))
```

Training Accuracy for decision tree is 1.0
 CV accuracy is 0.93 and standard deviation is 0.01
 Training Accuracy for random forest is 0.9970674486803519
 CV accuracy is 0.95 and standard deviation is 0.01
 Training Accuracy for naive bayes classifier is 0.9237536656891495
 CV accuracy is 0.92 and standard deviation is 0.01

```
In [14]: 1 from sklearn.metrics import confusion_matrix
2 from sklearn.model_selection import cross_val_score
3 y_pred_1 = clf1.predict(X_test)
4 conf_matrix = confusion_matrix(y_test, y_pred_1)
5 print("The confusion matrix for Decision Tree is: ")
6 print(conf_matrix)
7 print("Testing Accuracy for decision tree is ", accuracy_score(y_test, y_pred_1))
```

The confusion matrix for Decision Tree is:
 [[140 8]
 [10 70]]
 Testing Accuracy for decision tree is 0.9210526315789473

```
In [15]: 1 y_pred_2 = clf2.predict(X_test)
2 conf_matrix = confusion_matrix(y_test, y_pred_2)
3 print("The confusion matrix for Random Forest is: ")
4 print(conf_matrix)
5 print("Testing Accuracy for random forest is ", accuracy_score(y_test, y_pred_2))
```

The confusion matrix for Random Forest is:
 [[143 5]
 [11 69]]
 Testing Accuracy for random forest is 0.9298245614035088

```
In [16]: 1 y_pred_3 = clf3.predict(X_test)
2 conf_matrix = confusion_matrix(y_test, y_pred_3)
3 print("The confusion matrix for naive bayes classifier is: ")
4 print(conf_matrix)
5 print("Testing Accuracy for naive bayes classifier is ", accuracy_score(y_te
```

The confusion matrix for naive bayes classifier is:

```
[[141  7]
 [ 10 70]]
```

Testing Accuracy for naive bayes classifier is 0.9254385964912281

In [17]:

```

1  from sklearn.model_selection import train_test_split
2  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25,
3  print("The shape of X_train is ", X_train.shape)
4  print("The shape of y_train is ", y_train.shape)
5  print("The shape of X_test is ", X_test.shape)
6  print("The shape of y_test is ", y_test.shape)
7  from sklearn.preprocessing import RobustScaler
8  scaler = RobustScaler()
9  X_train = scaler.fit_transform(X_train)
10 X_test = scaler.transform(X_test)
11 from sklearn.metrics import accuracy_score
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.ensemble import RandomForestClassifier
14 from sklearn.naive_bayes import GaussianNB
15 #Decision Tree
16 clf1 = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
17 scores_1 = cross_val_score(clf1, X, y, cv=5)
18 clf1.fit(X_train, y_train)
19 print('Training Accuracy for decision tree is ', clf1.score(X_train, y_train)
20 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_1.mea
21 #Random Forest
22 clf2 = RandomForestClassifier(n_estimators = 8, criterion = 'entropy', rando
23 scores_2 = cross_val_score(clf2, X, y, cv=5)
24 clf2.fit(X_train, y_train)
25 print('Training Accuracy for random forest is ', clf2.score(X_train, y_train)
26 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_2.me
27 #Naive Bayes
28 clf3 = GaussianNB()
29 scores_3 = cross_val_score(clf3, X, y, cv=5)
30 clf3.fit(X_train, y_train)
31 print('Training Accuracy for naive bayes classifier is ', clf3.score(X_train
32 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_3.me
33 from sklearn.metrics import confusion_matrix
34 from sklearn.model_selection import cross_val_score
35 y_pred_1 = clf1.predict(X_test)
36 conf_matrix = confusion_matrix(y_test, y_pred_1)
37 print("The confusion matrix for Decision Tree is: ")
38 print(conf_matrix)
39 print("Testing Accuracy for decision tree is ", accuracy_score(y_test, y_pre
40 y_pred_2 = clf2.predict(X_test)
41 conf_matrix = confusion_matrix(y_test, y_pred_2)
42 print("The confusion matrix for Random Forest is: ")
43 print(conf_matrix)
44 print("Testing Accuracy for random forest is ", accuracy_score(y_test, y_pre
45 y_pred_3 = clf3.predict(X_test)
46 conf_matrix = confusion_matrix(y_test, y_pred_3)
47 print("The confusion matrix for naive bayes classifier is: ")
48 print(conf_matrix)
49 print("Testing Accuracy for naive bayes classifier is ", accuracy_score(y_te

```

The shape of X_train is (426, 23)

The shape of y_train is (426,)

The shape of X_test is (143, 23)

The shape of y_test is (143,)

Training Accuracy for decision tree is 1.0

CV accuracy is 0.93 and standard deviation is 0.01

Training Accuracy for random forest is 0.9929577464788732


```
CV accuracy is 0.95 and standard deviation is 0.01
Training Accuracy for naive bayes classifier is 0.9131455399061033
CV accuracy is 0.92 and standard deviation is 0.01
The confusion matrix for Decision Tree is:
[[84  5]
 [ 4 50]]
Testing Accuracy for decision tree is 0.9370629370629371
The confusion matrix for Random Forest is:
[[86  3]
 [ 4 50]]
Testing Accuracy for random forest is 0.951048951048951
The confusion matrix for naive bayes classifier is:
[[85  4]
 [ 4 50]]
Testing Accuracy for naive bayes classifier is 0.9440559440559441
```

In [18]:

```

1  from sklearn.model_selection import train_test_split
2  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.80,
3  print("The shape of X_train is ", X_train.shape)
4  print("The shape of y_train is ", y_train.shape)
5  print("The shape of X_test is ", X_test.shape)
6  print("The shape of y_test is ", y_test.shape)
7  from sklearn.preprocessing import RobustScaler
8  scaler = RobustScaler()
9  X_train = scaler.fit_transform(X_train)
10 X_test = scaler.transform(X_test)
11 from sklearn.metrics import accuracy_score
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.ensemble import RandomForestClassifier
14 from sklearn.naive_bayes import GaussianNB
15 #Decision Tree
16 clf1 = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
17 scores_1 = cross_val_score(clf1, X, y, cv=5)
18 clf1.fit(X_train, y_train)
19 print('Training Accuracy for decision tree is ', clf1.score(X_train, y_train)
20 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_1.me
21 #Random Forest
22 clf2 = RandomForestClassifier(n_estimators = 8, criterion = 'entropy', rando
23 scores_2 = cross_val_score(clf2, X, y, cv=5)
24 clf2.fit(X_train, y_train)
25 print('Training Accuracy for random forest is ', clf2.score(X_train, y_train)
26 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_2.me
27 #Naive Bayes
28 clf3 = GaussianNB()
29 scores_3 = cross_val_score(clf3, X, y, cv=5)
30 clf3.fit(X_train, y_train)
31 print('Training Accuracy for naive bayes classifier is ', clf3.score(X_train
32 print("CV accuracy is %0.2f and standard deviation is %0.2f" % (scores_3.me
33 from sklearn.metrics import confusion_matrix
34 from sklearn.model_selection import cross_val_score
35 y_pred_1 = clf1.predict(X_test)
36 conf_matrix = confusion_matrix(y_test, y_pred_1)
37 print("The confusion matrix for Decision Tree is: ")
38 print(conf_matrix)
39 print("Testing Accuracy for decision tree is ", accuracy_score(y_test, y_pre
40 y_pred_2 = clf2.predict(X_test)
41 conf_matrix = confusion_matrix(y_test, y_pred_2)
42 print("The confusion matrix for Random Forest is: ")
43 print(conf_matrix)
44 print("Testing Accuracy for random forest is ", accuracy_score(y_test, y_pre
45 y_pred_3 = clf3.predict(X_test)
46 conf_matrix = confusion_matrix(y_test, y_pred_3)
47 print("The confusion matrix for naive bayes classifier is: ")
48 print(conf_matrix)
49 print("Testing Accuracy for naive bayes classifier is ", accuracy_score(y_te

```

The shape of X_train is (113, 23)

The shape of y_train is (113,)

The shape of X_test is (456, 23)

The shape of y_test is (456,)

Training Accuracy for decision tree is 1.0

CV accuracy is 0.93 and standard deviation is 0.01

```

Training Accuracy for random forest is 0.9911504424778761
CV accuracy is 0.95 and standard deviation is 0.01
Training Accuracy for naive bayes classifier is 0.911504424778761
CV accuracy is 0.92 and standard deviation is 0.01
The confusion matrix for Decision Tree is:
[[255  35]
 [ 18 148]]
Testing Accuracy for decision tree is 0.8837719298245614
The confusion matrix for Random Forest is:
[[274  16]
 [ 18 148]]
Testing Accuracy for random forest is 0.9254385964912281
The confusion matrix for naive bayes classifier is:
[[268  22]
 [ 17 149]]
Testing Accuracy for naive bayes classifier is 0.9144736842105263

```

```

In [19]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25,
3 print("For the 2nd experiment, the shape of X_train is ", X_train.shape)
4 print("For the 2nd experiment, the shape of y_train is ", y_train.shape)
5 print("For the 2nd experiment, the shape of X_test is ", X_test.shape)
6 print("For the 2nd experiment, the shape of y_test is ", y_test.shape)
7 from sklearn.model_selection import train_test_split
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.80,
9 print("For the 3rd experiment, the shape of X_train is ", X_train.shape)
10 print("For the 3rd experiment, the shape of y_train is ", y_train.shape)
11 print("For the 3rd experiment, the shape of X_test is ", X_test.shape)
12 print("For the 3rd experiment, the shape of y_test is ", y_test.shape)

```

```

For the 2nd experiment, the shape of X_train is (426, 23)
For the 2nd experiment, the shape of y_train is (426,)
For the 2nd experiment, the shape of X_test is (143, 23)
For the 2nd experiment, the shape of y_test is (143,)
For the 3rd experiment, the shape of X_train is (113, 23)
For the 3rd experiment, the shape of y_train is (113,)
For the 3rd experiment, the shape of X_test is (456, 23)
For the 3rd experiment, the shape of y_test is (456,)

```