In [1]: #importing all the necessary libaries import pandas as pd import numpy as np import matplotlib.pyplot as plt from sklearn.preprocessing import LabelEncoder from sklearn.linear\_model import LogisticRegression from sklearn.model\_selection import train\_test\_split from sklearn.svm import SVC from sklearn.naive\_bayes import GaussianNB from sklearn.ensemble import RandomForestClassifier from statsmodels.stats.outliers\_influence import variance\_inflation\_factor from sklearn.metrics import accuracy\_score, confusion\_matrix, roc\_curve, roc\_auc\_score import seaborn as sns sns.set() import os In [3]: #Reading the data data = pd.read\_csv('C:\\Users\\nakkl\\OneDrive\\Desktop\\data.csv') In [4]: data Out[4]: concave ... radius\_worst texture\_worst perimeter\_worst area\_worst smoothness id diagnosis radius\_mean texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean concavity\_mean points\_mean 842302 122.80 1001.0 0.14710 ... M 17.99 10.38 0.11840 0.27760 0.30010 25.380 17.33 184.60 2019.0 842517 M 20.57 17.77 132.90 1326.0 0.08474 0.07864 0.08690 0.07017 ... 24.990 23.41 158.80 1956.0 2 84300903 M 19.69 21.25 130.00 1203.0 0.10960 0.15990 0.19740 0.12790 ... 23.570 25.53 152.50 1709.0 0.14250 **3** 84348301 11.42 20.38 77.58 386.1 0.28390 0.24140 0.10520 ... 14.910 26.50 98.87 567.7 **4** 84358402 M 20.29 14.34 135.10 1297.0 0.10030 0.13280 0.19800 0.10430 ... 22.540 16.67 152.20 1575.0 926424 142.00 1479.0 0.11100 0.11590 0.24390 0.13890 ... 166.10 2027.0 564 M 21.56 22.39 25.450 26.40 1261.0 0.09780 0.10340 1731.0 565 926682 M 20.13 28.25 131.20 0.14400 0.09791 ... 23.690 38.25 155.00 108.30 858.1 0.08455 0.10230 0.09251 34.12 1124.0 566 926954 M 16.60 28.08 0.05302 ... 18.980 126.70 20.60 140.10 567 927241 M 29.33 1265.0 0.11780 0.27700 0.35140 0.15200 ... 25.740 39.42 184.60 1821.0 568 92751 В 7.76 24.54 47.92 181.0 0.05263 0.04362 0.00000 0.00000 ... 9.456 30.37 59.16 268.6 569 rows × 32 columns #print the top 5 coloums data.head(5) Out[6]: concave ... radius\_worst texture\_worst perimeter\_worst area\_worst smoothness\_v id diagnosis radius\_mean texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean concavity\_mean points\_mean 842302 Μ 17.99 10.38 122.80 1001.0 0.11840 0.27760 0.3001 0.14710 ... 17.33 184.60 2019.0 0. 25.38 842517 20.57 17.77 132.90 1326.0 0.08474 0.07864 0.0869 0.07017 ... 24.99 23.41 158.80 1956.0 130.00 0.1974 152.50 0. **2** 84300903 19.69 21.25 1203.0 0.10960 0.15990 0.12790 ... 23.57 25.53 1709.0 0. **3** 84348301 11.42 77.58 386.1 0.14250 0.28390 0.2414 0.10520 ... 14.91 26.50 98.87 567.7 **4** 84358402 20.29 14.34 135.10 1297.0 0.10030 0.13280 0.1980 0.10430 ... 22.54 16.67 152.20 1575.0 0. M 5 rows × 32 columns In [7]: #getting the basic information data.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 569 entries, 0 to 568 Data columns (total 32 columns): # Column Non-Null Count Dtype --- ----569 non-null diagnosis 569 non-null object radius\_mean 569 non-null float64 3 texture\_mean 569 non-null float64 perimeter\_mean 569 non-null float64 4 5 area\_mean 569 non-null float64 569 non-null float64 6 smoothness\_mean 569 non-null float64 7 compactness\_mean 569 non-null float64 8 concavity\_mean 9 concave points\_mean 569 non-null float64 569 non-null float64 10 symmetry\_mean 569 non-null float64 fractal\_dimension\_mean 11 float64 12 radius\_se 569 non-null 569 non-null float64 13 texture\_se 569 non-null float64 14 perimeter\_se 15 area\_se 569 non-null float64 float64 16 smoothness\_se 569 non-null 17 compactness\_se 569 non-null float64 float64 18 concavity\_se 569 non-null concave points\_se float64 19 569 non-null 20 569 non-null float64 symmetry\_se 21 fractal\_dimension\_se 569 non-null float64 22 569 non-null float64 radius\_worst 23 569 non-null float64 texture\_worst 24 569 non-null float64 perimeter\_worst 25 area\_worst 569 non-null float64 26 569 non-null float64 smoothness\_worst 569 non-null 27 float64 compactness\_worst 28 569 non-null float64 concavity\_worst 29 569 non-null float64 concave points\_worst 569 non-null 30 float64 symmetry\_worst 31 fractal\_dimension\_worst 569 non-null float64 dtypes: float64(30), int64(1), object(1) memory usage: 142.4+ KB In [8]: #getting the statistical data data.describe() Out[8]: concave id radius\_mean texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean concavity\_mean symmetry\_mean ... radius\_worst texture\_worst perimeter\_worst area\_wors points\_mean **count** 5.690000e+02 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 ... 569.000000 569.000000 569.000000 569.00000 mean 3.037183e+07 14.127292 19.289649 91.969033 654.889104 0.096360 0.104341 0.088799 0.048919 0.181162 ... 16.269190 107.261213 880.58312 25.677223 std 1.250206e+08 3.524049 4.301036 24.298981 351.914129 0.014064 0.052813 0.079720 0.038803 0.027414 ... 4.833242 6.146258 33.602542 569.35699 min 8.670000e+03 9.710000 43.790000 143.500000 0.052630 0.019380 0.000000 0.000000 0.106000 ... 7.930000 12.020000 50.410000 185.20000 6.981000 0.029560 0.020310 8.692180e+05 11.700000 16.170000 75.170000 420.300000 0.086370 0.064920 0.161900 ... 13.010000 21.080000 84.110000 515.30000 **50%** 9.060240e+05 13.370000 18.840000 86.240000 551.100000 0.095870 0.092630 0.061540 0.033500 0.179200 ... 14.970000 25.410000 97.660000 686.50000 0.105300 0.130400 0.130700 8.813129e+06 15.780000 21.800000 104.100000 782.700000 0.074000 0.195700 ... 18.790000 29.720000 125.400000 1084.00000 max 9.113205e+08 28.110000 39.280000 188.500000 2501.000000 0.163400 0.345400 0.426800 0.201200 0.304000 ... 36.040000 49.540000 251.200000 4254.00000 8 rows × 31 columns #count the number of rows and coloums in data In [10]: data.shape (569, 32)Out[10]: In [11]: #get the empty values(NaN, NAN, na) in each coloum data.isna().sum() 0 id Out[11]: diagnosis 0 radius\_mean 0 texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean concavity\_mean concave points\_mean symmetry\_mean 0 fractal\_dimension\_mean radius\_se texture\_se perimeter\_se area\_se smoothness\_se compactness\_se concavity\_se concave points\_se symmetry\_se fractal\_dimension\_se radius\_worst texture\_worst perimeter\_worst area\_worst smoothness\_worst compactness\_worst concavity\_worst concave points\_worst symmetry\_worst fractal\_dimension\_worst dtype: int64 since there are no empty values so no need of filling the data we can use it as it is In [12]: #Get the count of Maligant(M) or Benign(B) cells data['diagnosis'].value\_counts() 357 Out[12]: 212 Name: diagnosis, dtype: int64 #visualise the count sns.countplot(data['diagnosis'],label='count') C:\Users\nakkl\OneDrive\Documents\Python Scripts\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn( <AxesSubplot:xlabel='diagnosis', ylabel='count'> Out[13]: 350 300 250 200 200 150 100 50 0 M В diagnosis In [14]: #looking at the data which are to be encoded(replaced) data.dtypes int64 id Out[14]: diagnosis object radius\_mean float64 float64 texture\_mean float64 perimeter\_mean area\_mean float64 float64 smoothness\_mean compactness\_mean float64 concavity\_mean float64 concave points\_mean float64 float64 symmetry\_mean fractal\_dimension\_mean float64 float64 radius\_se texture\_se float64 float64 perimeter\_se float64 area\_se float64 smoothness\_se compactness\_se float64 float64 concavity\_se float64 concave points\_se float64 symmetry\_se fractal\_dimension\_se float64 float64 radius\_worst float64 texture\_worst perimeter\_worst float64 area\_worst float64 float64 smoothness\_worst float64 compactness\_worst float64 concavity\_worst concave points\_worst float64 float64 symmetry\_worst float64 fractal\_dimension\_worst dtype: object In [16]: #encoding the data labelencoder\_y = LabelEncoder() data.iloc[:,1] = labelencoder\_y.fit\_transform(data.iloc[:,1].values) data.iloc[:,1] 1 Out[16]: 1 1 3 1 4 564 1 565 1 566 1 567 1 568 Name: diagnosis, Length: 569, dtype: int32 In [21]: #create a pair plot sns.pairplot(data.iloc[:,1:6], hue='diagnosis') <seaborn.axisgrid.PairGrid at 0x16dc87ddd00> 25 radius\_mean 20 15 10 40 35 texture\_mean 30 20 15 10 diagnosis 175 perimeter\_mean 150 125 100 75 50 2500 2000 area\_mean 1500 1000 500 20 30 150 200 1000 2000 3000 40 radius\_mean texture\_mean perimeter\_mean area\_mean In [23]: #print the frist 5 rows of new data data.head(5) Out[23]: concave id diagnosis radius\_mean texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean concavity\_mean radius\_worst texture\_worst perimeter\_worst area\_worst smoothness\_v points\_mean 842302 122.80 0.11840 17.33 184.60 2019.0 0. 1 17.99 10.38 1001.0 0.27760 0.3001 0.14710 ... 25.38 132.90 0.08474 0.07864 158.80 842517 20.57 17.77 1326.0 0.0869 0.07017 ... 24.99 23.41 1956.0 **2** 84300903 19.69 21.25 130.00 1203.0 0.10960 0.15990 0.1974 0.12790 ... 23.57 25.53 152.50 1709.0 77.58 0.14250 0.2414 567.7 0. **3** 84348301 11.42 20.38 386.1 0.28390 0.10520 14.91 26.50 98.87 **4** 84358402 20.29 14.34 135.10 1297.0 0.10030 0.13280 0.1980 0.10430 ... 22.54 16.67 152.20 1575.0 5 rows × 32 columns In [24]: #get the correlation of the coloums data.iloc[:,1:12].corr() Out[24]: diagnosis radius\_mean texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean concavity\_mean concave points\_mean symmetry\_mean fractal\_dimension\_mean diagnosis 1.000000 0.358560 0.596534 0.776614 0.330499 -0.012838 0.730029 0.415185 0.742636 0.708984 0.696360 0.730029 1.000000 0.323782 0.997855 0.987357 0.170581 0.506124 0.676764 0.822529 0.147741 -0.311631 radius\_mean 0.321086 -0.023389 -0.076437 texture\_mean 0.415185 0.323782 1.000000 0.329533 0.236702 0.302418 0.293464 0.071401 perimeter\_mean 0.997855 0.329533 1.000000 0.986507 0.207278 0.556936 0.716136 0.850977 0.183027 -0.261477 1.000000 0.177028 0.498502 0.685983 0.823269 0.151293 -0.283110 area\_mean 0.708984 0.987357 0.321086 0.986507 smoothness\_mean 0.358560 0.170581 -0.023389 0.207278 0.177028 1.000000 0.659123 0.521984 0.553695 0.557775 0.584792 0.498502 0.659123 0.831135 0.596534 0.506124 0.236702 0.556936 1.000000 0.883121 0.602641 0.565369 compactness\_mean 0.696360 0.676764 0.302418 0.716136 0.685983 0.521984 0.883121 1.000000 0.921391 0.500667 0.336783 concavity\_mean 0.823269 0.553695 0.831135 0.921391 1.000000 0.462497 0.166917 concave points\_mean 0.776614 0.822529 0.293464 0.850977 0.330499 0.147741 0.071401 0.183027 0.151293 0.557775 0.602641 0.500667 0.462497 1.000000 0.479921 symmetry\_mean -0.012838 -0.311631 -0.076437 -0.261477 -0.283110 0.584792 0.565369 0.336783 0.166917 0.479921 1.000000 fractal\_dimension\_mean #visualize the correlation plt.figure(figsize=(8,8)) sns.heatmap(data.iloc[:,1:12].corr()) <AxesSubplot:> Out[39]: diagnosis radius\_mean - 0.8 texture\_mean - 0.6 perimeter\_mean area\_mean - 0.4 smoothness\_mean compactness\_mean - 0.2 concavity\_mean 0.0 concave points\_mean symmetry\_mean fractal\_dimension\_mean perimeter\_mean smoothness\_mean concave points\_mean symmetry\_mean fractal\_dimension\_mean texture\_mean compactness\_mean In [40]: plt.figure(figsize=(8,8)) sns.heatmap(data.iloc[:,1:12].corr(), annot=True, fmt='.0%') #annot(annotation) is used to label the data #fmt(format) is used to show given labeled data in form of percentage <AxesSubplot:> 1.0 diagnosis 100% 73% 42% 74% 71% 36% 60% 70% 78% 33% -1% 73% 100% 32% 100% 99% 17% 51% 68% 82% 15% -31% radius\_mean - 0.8 42% 32% 100% 33% 32% -2% 24% 30% 29% texture\_mean - 0.6 74% 100% 33% 100% 99% 21% 56% 72% 85% 18% -26% perimeter\_mean 71% 99% 32% 99% 100% 18% 50% 69% 82% 15% -28% area\_mean - 0.4 36% 17% -2% 21% 18% 100% 66% 52% 55% 56% 58% smoothness\_mean 60% 51% 24% 56% 50% 66% 100% 88% 83% 60% 57% compactness\_mean - 0.2 70% 68% 30% 72% 69% 52% 88% 100% 92% 50% 34% concavity\_mean 0.0 78% 82% 29% 85% 82% 55% 83% 92% 100% 46% 17% concave points\_mean symmetry\_mean 33% 15% 7% 18% 15% 56% 60% 50% 46% 100% 48% - -0.2 -1% -31% -8% -26% -28% 58% 57% 34% 17% 48% 100% fractal\_dimension\_mean radius\_mean texture\_mean smoothness\_mean concave points\_mean fractal\_dimension\_mean perimeter\_mean compactness\_mean concavity\_mean symmetry\_mean In [35]: #split the data set into independent(X) and dependent(Y) data sets X = data.iloc[:,2:31].values#x has the features of cancer or not Y = data.iloc[:,1].values #y has the diagnosis In [36]: #split the data set into 75% training and 25% testing X\_train, X\_test, Y\_train, Y\_test = train\_test\_split(X, Y, test\_size = 0.25 , random\_state = 0) In [41]: #Scale the data(feature scaling) from sklearn.preprocessing import StandardScaler sc = StandardScaler() X\_test = sc.fit\_transform(X\_test) X\_train = sc.fit\_transform(X\_train) In [43]: #create a function for the models def models(X\_train, Y\_train): #Logistic Regression log = LogisticRegression(random\_state=0) log.fit(X\_train, Y\_train) **#Decision Tree** from sklearn.tree import DecisionTreeClassifier tree = DecisionTreeClassifier(criterion = 'entropy', random\_state=0 ) tree.fit(X\_train, Y\_train) #Random Forest Classifier forest = RandomForestClassifier(n\_estimators = 10, criterion = 'entropy') forest.fit(X\_train,Y\_train) #Print the models accuarcy on the training data print('[0]Logistic Regression Training Accuracy:', log.score(X\_train,Y\_train)) print('[1]Decision Tree Classifier Training Accuracy:', tree.score(X\_train,Y\_train)) print('[2]Random Forest Classifier Training Accuracy:', forest.score(X\_train,Y\_train)) return log, tree, forest #Getting all of the models model = models(X\_train, Y\_train) [0]Logistic Regression Training Accuracy: 0.9906103286384976 [1]Decision Tree Classifier Training Accuracy: 1.0 [2] Random Forest Classifier Training Accuracy: 0.9976525821596244 In [46]: #test model accuracy on test data on confusion matrix for i in range( len(model) ): print('Model' , i) cm = confusion\_matrix(Y\_test, model[i].predict(X\_test)) TP = cm[0][0]TN = cm[1][1]FP = cm[1][0]FN = cm[0][1]print(cm) print('Testing Accuracy =', (TP + TN)/ (TP + TN + FN + FP)) print() Model 0 [[86 4] [ 3 50]] Testing Accuracy = 0.951048951048951Model 1 [[83 7] [ 2 51]] Testing Accuracy = 0.9370629370629371Model 2 [[89 1] [ 6 47]] Testing Accuracy = 0.951048951048951#another way to find the accuracy In [48]: from sklearn.metrics import classification\_report for i in range( len(model) ): print('Model' , i) print( classification\_report(Y\_test,model[i].predict(X\_test))) print( accuracy\_score(Y\_test, model[0].predict(X\_test))) print() Model 0 recall f1-score support precision 0.96 90 0 0.97 0.96 1 0.93 0.94 0.93 53 0.95 143 accuracy macro avg 0.95 0.95 0.95 143 weighted avg 0.95 0.95 143 0.951048951048951 Model 1 precision recall f1-score support 0 0.98 0.92 0.95 90 1 0.88 0.96 0.92 53 0.94 143 accuracy 0.94 0.93 143 macro avg 0.93 weighted avg 0.94 0.94 0.94 143 0.951048951048951 Model 2 recall f1-score precision 0 0.94 0.99 0.96 90 1 0.98 0.89 0.93 53 0.95 143 accuracy macro avg 0.96 0.94 0.95 143 weighted avg 0.95 0.95 0.95 143 0.951048951048951 In [49]: #print the prediction of Random Forest Classifier pred = model[2].predict(X\_test) print(pred) print() print(Y\_test) 1 1 0 0 0 0 0 0 0 0 0 1 0 1 1 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 1 0 1 0 0 0 0 0 1 0 0 0 0 0 1 1 0 0 0 0 1