

▼ Breast Cancer Detection (Project)

Name - Mushtaq k Islam

Roll No. 2102025

Dept- CSE / MTech

```
#importing the essential libraries
```

```
import numpy as np  
import matplotlib.pyplot as plt  
import pandas as pd  
import seaborn as sns
```

Load The Data

```
from sklearn.datasets import load_breast_cancer  
#loading the Breast cancer data in a variable  
cancer = load_breast_cancer()
```

cancer

```
0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0,
1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1,
1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0,
0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0,
1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1],
'target_names': array(['malignant', 'benign'], dtype='<U9')}
```

```
cancer.keys()
```

```
#Printing the keys of the Dictionary, to get enough details about the data.
```

```
dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_names', 'filen
```

```
print(cancer['DESCR'])
```

```
#Will basically give you broad description of Breast Cancer Dataset
```

	mean	std
smoothness (worst):	0.071	0.223
compactness (worst):	0.027	1.058
concavity (worst):	0.0	1.252
concave points (worst):	0.0	0.291
symmetry (worst):	0.156	0.664
fractal dimension (worst):	0.055	0.208

```
:Missing Attribute Values: None
```

```
:Class Distribution: 212 - Malignant, 357 - Benign
```

```
:Creator: Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian
```

```
:Donor: Nick Street
```

```
:Date: November, 1995
```

This is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.

<https://goo.gl/U2Uwz2>

Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image.

Separating plane described above was obtained using Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree Construction Via Linear Programming." Proceedings of the 4th Midwest Artificial Intelligence and Cognitive Science Society, pp. 97-101, 1992], a classification method which uses linear programming to construct a decision tree. Relevant features were selected using an exhaustive search in the space of 1-4 features and 1-3 separating planes.

The actual linear program used to obtain the separating plane in the 3-dimensional space is that described in:

[K. P. Bennett and O. L. Mangasarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server:

```
ftp ftp.cs.wisc.edu  
cd math-prog/cpo-dataset/machine-learn/WDBC/
```

.. topic:: References

- W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on Electronic Imaging: Science and Technology, volume 1905, pages 861-870, San Jose, CA, 1993.
 - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and prognosis via linear programming. Operations Research, 43(4), pages 570-577, July-August 1995.
 - W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994) 163-171.

```
print(cancer['target'])  
#Printing the Target column, which is either 0=Malignenet or 1=Benign
```

```
#printing the Target names
```

```
'mean smoothness' 'mean compactness' 'mean concavity'  
'mean concave points' 'mean symmetry' 'mean fractal dimension'  
'radius error' 'texture error' 'perimeter error' 'area error'  
'smoothness error' 'compactness error' 'concavity error'  
'concave points error' 'symmetry error' 'fractal dimension error'  
'worst radius' 'worst texture' 'worst perimeter' 'worst area'
```

```
'worst smoothness' 'worst compactness' 'worst concavity'  
'worst concave points' 'worst symmetry' 'worst fractal dimension']
```

```
cancer['data'].shape
```

```
#Shape of the Data which is 569 Rows and 30 Columns
```

```
(569, 30)
```

```
df_cancer = pd.DataFrame(np.c_[cancer['data'], cancer['target']], columns= np.append(cancer
```

```
#Creating a Pandas DataFrame to deal with the data in a much easier way.
```

```
#While creating this DataFrame I added the additional column which is 'target' to the df_cancer
```

```
df_cancer.head()
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.28300	0.28300
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.10520	0.10520
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.20500	0.20500
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.14710	0.14710
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.14710	0.14710

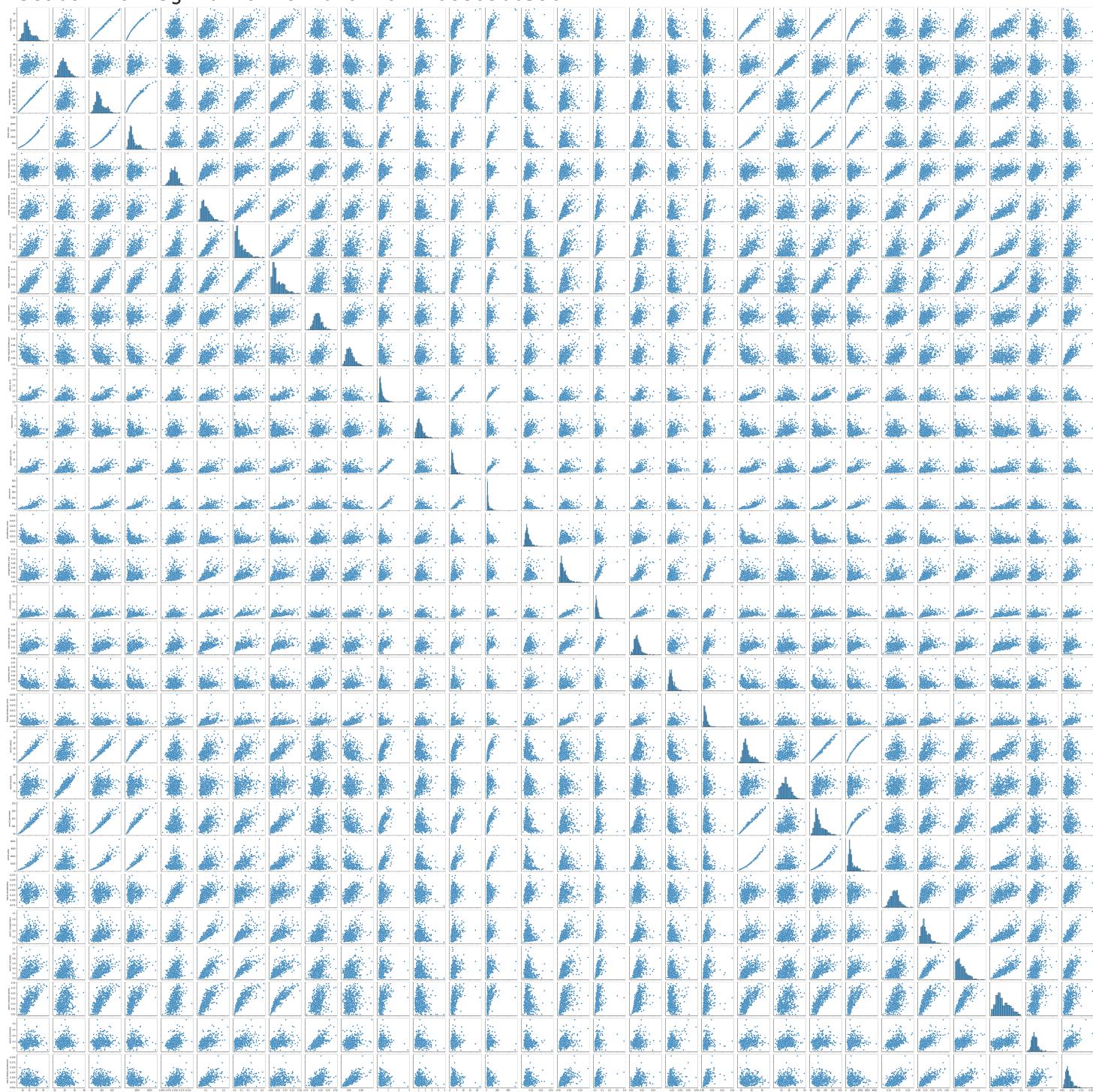
```
5 rows × 31 columns
```



▼ Visualising the data

```
sns.pairplot(df_cancer , vars =['mean radius', 'mean texture', 'mean perimeter', 'mean area',  
'mean smoothness', 'mean compactness' , 'mean concavity',  
'mean concave points', 'mean symmetry', 'mean fractal dimension',  
'radius error', 'texture error', 'perimeter error', 'area error',  
'smoothness error', 'compactness error', 'concavity error',  
'concave points error', 'symmetry error', 'fractal dimension error',  
'worst radius', 'worst texture', 'worst perimeter', 'worst area',  
'worst smoothness', 'worst compactness', 'worst concavity',  
'worst concave points', 'worst symmetry', 'worst fractal dimension'])
```

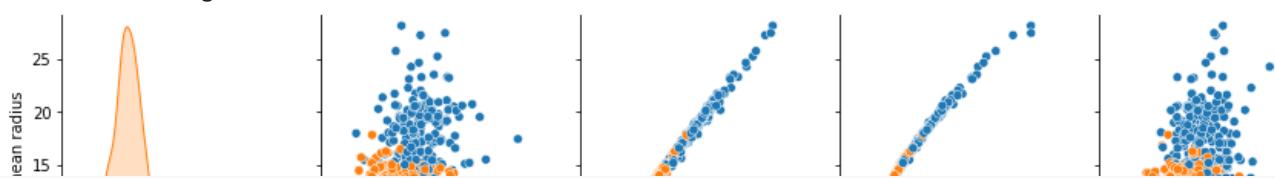
<seaborn.axisgrid.PairGrid at 0x7f06809bb590>



```
#in above plotting we are not able to differentiate much,so we use 'hue' on target column, v
sns.pairplot(df_cancer, hue = 'target', vars = ['mean radius', 'mean texture', 'mean perime
 'mean smoothness'])

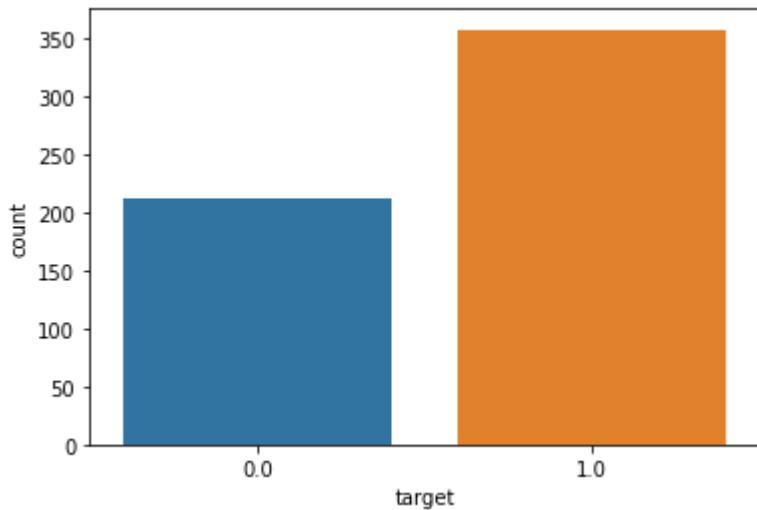
#blue points are malignant case which are severe cases or life threatnning cases.
#orange points are not very severe or life threatnning
```

```
<seaborn.axisgrid.PairGrid at 0x7f06694b87d0>
```



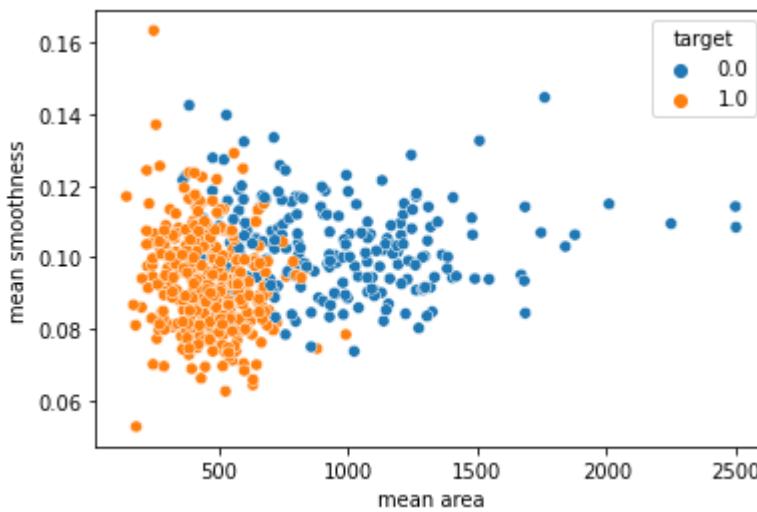
```
#will simply tell you how many Malignant and Benign cases we have.  
#Malignant= 200~ and Benign = 350~ approx.  
sns.countplot(df_cancer['target'])
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass  
FutureWarning  
<matplotlib.axes._subplots.AxesSubplot at 0x7f0664a1c7d0>
```

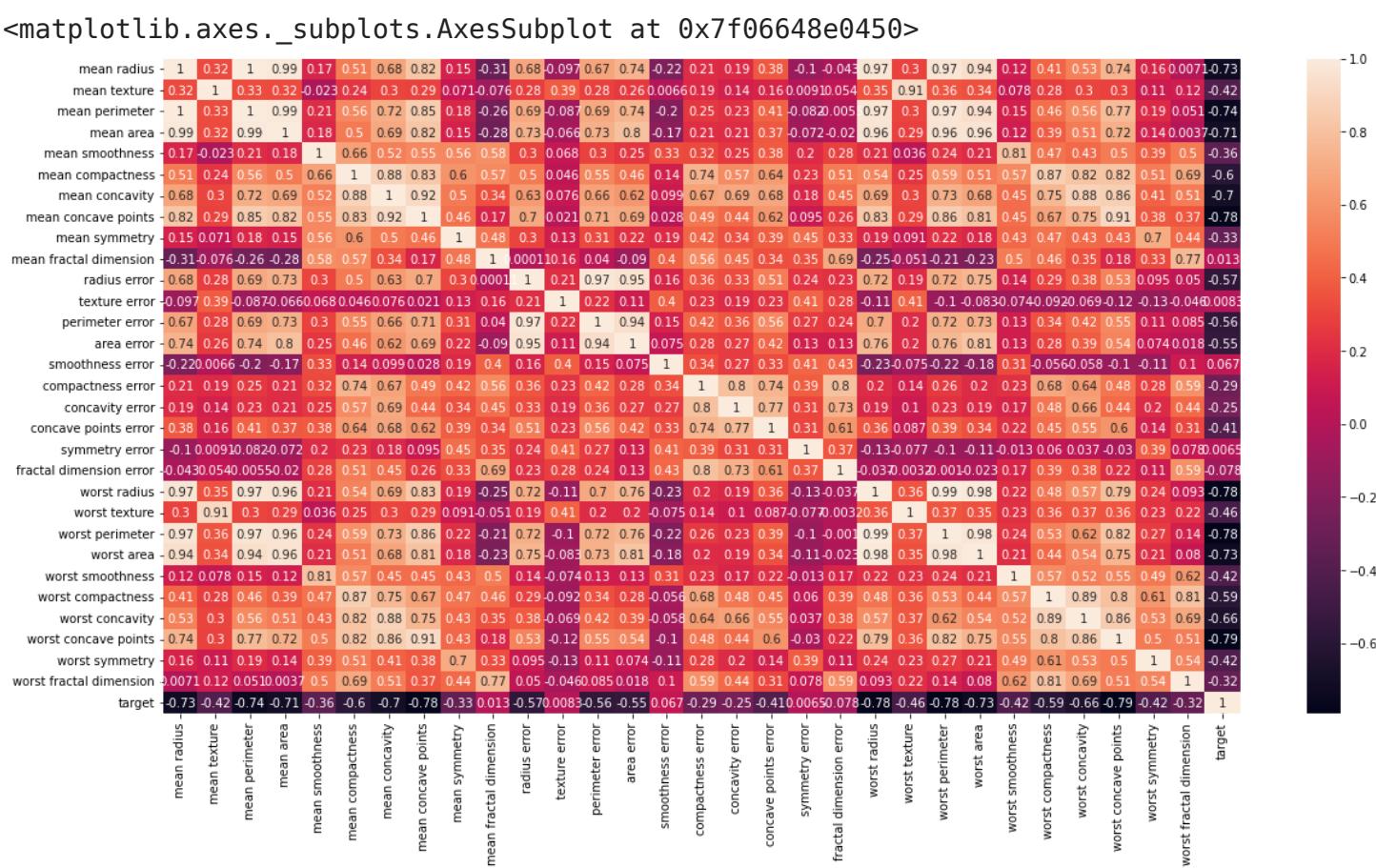


```
#plotting a scatter plot diagram for mean area anf mean smoothness, you can plot any feature  
sns.scatterplot(x='mean area',y='mean smoothness',hue='target',data =df_cancer)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f0664890350>
```



```
#here we made a heatmap figure of correlation of all the columns  
plt.figure(figsize =(20,10))  
sns.heatmap(df_cancer.corr(), annot =True)
```



SPLITTING THE DATASET

```
x = df_cancer.drop(['target'], axis = 1)
#train test split
```

X

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	sym
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08690	0.07017	
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.12790	
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	
...

```
y= df_cancer['target']
```

```
y
```

```
0      0.0
1      0.0
2      0.0
3      0.0
4      0.0
      ...
564    0.0
565    0.0
566    0.0
567    0.0
568    1.0
```

```
Name: target, Length: 569, dtype: float64
```

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=5)
```

```
x_train
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean sym
--	----------------	-----------------	-------------------	--------------	--------------------	---------------------	-------------------	---------------------------	-------------

y_train

```

306    1.0
410    1.0
197    0.0
376    1.0
244    0.0
...
8      0.0
73     0.0
400    0.0
118    0.0
206    1.0
Name: target, Length: 455, dtype: float64

```

y_test

```

28     0.0
163    1.0
123    1.0
361    1.0
549    1.0
...
414    0.0
515    1.0
186    0.0
3      0.0
261    0.0
Name: target, Length: 114, dtype: float64

```

▼ TRAINING THE MODEL USING LOGISTIC REGRESSION

```

from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix

logistic_model = LogisticRegression(random_state = 0)
logistic_model.fit(x_train, y_train)

/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: ConvergenceWarning: 
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
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
LogisticRegression(random_state=0)

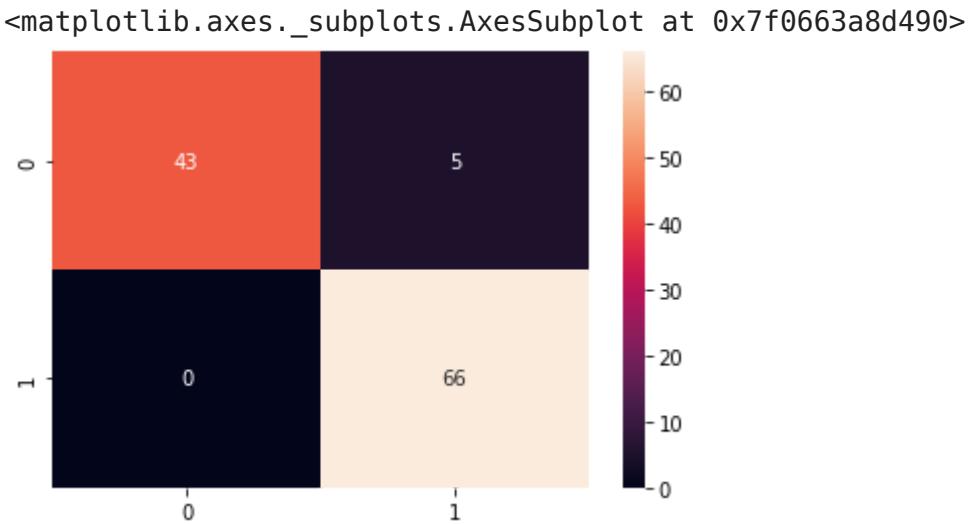
y_predict =logistic_model.predict(x_test)
y_predict

```

```
array([0., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1.,  
      1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 0., 0., 0., 0., 0., 0.,  
      1., 1., 0., 1., 1., 0., 1., 1., 1., 1., 0., 1., 1., 1., 0., 0., 1., 1., 1., 1., 1., 1.,  
      1., 1., 1., 1., 0., 0., 0., 1., 0., 0., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,  
      1., 0., 1., 0., 1., 1., 1., 1., 1., 0., 1., 0., 1., 1., 0., 1., 0., 1., 1., 0., 0., 0., 0.,  
      1., 0., 1., 0., 0., 0., 0., 1., 1., 0., 0., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0.,  
      1., 1., 0., 0., 1., 0., 1., 0., 1., 0., 1., 0., 1., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 0.])
```

```
cm = confusion_matrix(y_test,y_predict)
```

```
sns.heatmap(cm ,annot=True)
```



```
print(classification_report(y_test,y_predict))
```

	precision	recall	f1-score	support
0.0	1.00	0.90	0.95	48
1.0	0.93	1.00	0.96	66
accuracy			0.96	114
macro avg	0.96	0.95	0.95	114
weighted avg	0.96	0.96	0.96	114

The accuracy we received is 96% or 95.8% for exact accuracy.

Using KNeighborsClassifier Method of neighbors class to use Nearest Neighbor algorithm

```
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.metrics import classification_report , confusion matrix
```

```
knn_model = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
knn_model.fit(x_train, y_train)
```

KNeighborsClassifier()

```
y_predict =knn_model.predict(x_test)
```

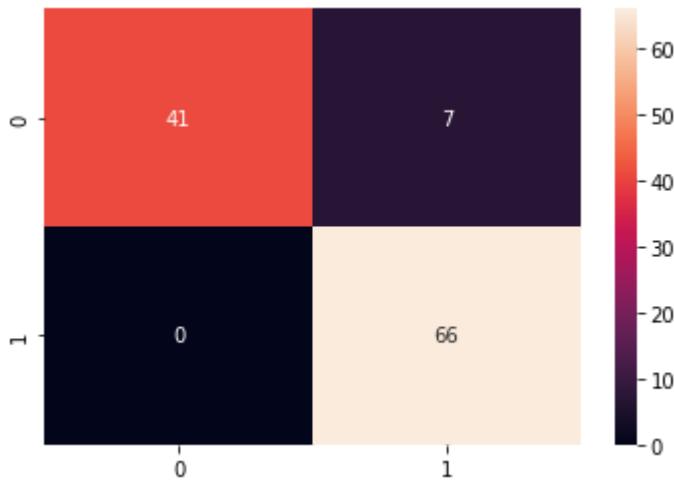
```
y_predict
```

```
array([0., 1., 1., 1., 0., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1.,
       1., 1., 1., 1., 1., 1., 1., 1., 0., 0., 1., 0., 0., 0., 0., 1., 0.,
       1., 1., 0., 1., 1., 0., 1., 1., 0., 1., 1., 0., 0., 1., 0., 1., 1.,
       1., 1., 1., 0., 1., 0., 1., 0., 0., 1., 1., 1., 1., 1., 1., 1., 1.,
       1., 0., 1., 0., 1., 1., 1., 1., 1., 0., 1., 0., 1., 1., 1., 1., 1.,
       1., 0., 1., 0., 0., 0., 1., 1., 0., 0., 1., 1., 1., 1., 1., 1., 1.,
       1., 1., 0., 0., 1., 0., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 0.,
       1., 1., 0., 0., 1., 0., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 0.])
```

```
cm = confusion_matrix(y_test,y_predict)
```

```
sns.heatmap(cm ,annot=True)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f0663807690>
```



```
from sklearn import metrics
```

```
print("Train set Accuracy: ", metrics.accuracy_score(y_train, knn_model.predict(x_train)))
```

```
print("Test set Accuracy: ", metrics.accuracy_score(y_test, y_predict))
```

```
Train set Accuracy: 0.945054945054945
```

```
Test set Accuracy: 0.9385964912280702
```

```
print(classification_report(y_test,y_predict))
```

	precision	recall	f1-score	support
0.0	1.00	0.85	0.92	48
1.0	0.90	1.00	0.95	66
accuracy			0.94	114
macro avg	0.95	0.93	0.94	114
weighted avg	0.94	0.94	0.94	114

Trying for a random k value for example 9

```
k = 9
```

```
neigh6 = KNeighborsClassifier(n_neighbors = k).fit(x_train,y_train)
```

```

yhat6 = neigh6.predict(x_test)

print("Train set Accuracy: ", metrics.accuracy_score(y_train, neigh6.predict(x_train)))
print("Test set Accuracy: ", metrics.accuracy_score(y_test, yhat6))

Train set Accuracy:  0.9384615384615385
Test set Accuracy:  0.9385964912280702

```

```

Ks = 10
mean_acc = np.zeros((Ks-1))
std_acc = np.zeros((Ks-1))
ConfustionMx = [];
for n in range(1,Ks):

    #Train Model and Predict
    neigh = KNeighborsClassifier(n_neighbors = n).fit(x_train,y_train)
    yhat=neigh.predict(x_test)
    mean_acc[n-1] = metrics.accuracy_score(y_test, yhat)

    std_acc[n-1]=np.std(yhat==y_test)/np.sqrt(yhat.shape[0])

mean_acc

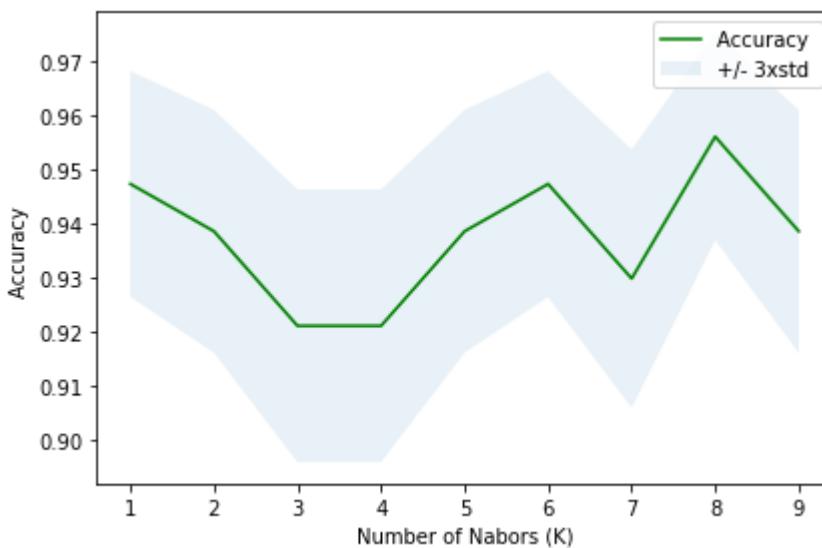
array([0.94736842, 0.93859649, 0.92105263, 0.92105263, 0.93859649,
       0.94736842, 0.92982456, 0.95614035, 0.93859649])

```

```

plt.plot(range(1,Ks),mean_acc,'g')
plt.fill_between(range(1,Ks),mean_acc - 1 * std_acc,mean_acc + 1 * std_acc, alpha=0.10)
plt.legend(['Accuracy ', '+/- 3xstd'])
plt.ylabel('Accuracy ')
plt.xlabel('Number of Nabors (K)')
plt.tight_layout()
plt.show()

```



```

print( "The best accuracy was with", mean_acc.max(), "with k=", mean_acc.argmax()+1)

The best accuracy was with 0.956140350877193 with k= 8

```

```
k = 8
```

```
neigh6 = KNeighborsClassifier(n_neighbors = k).fit(x_train,y_train)

yhat6 = neigh6.predict(x_test)

print("Train set Accuracy: ", metrics.accuracy_score(y_train, neigh6.predict(x_train)))
print("Test set Accuracy: ", metrics.accuracy_score(y_test, yhat6))

Train set Accuracy:  0.9318681318681319
Test set Accuracy:  0.956140350877193
```

Now that we have the value of k, we can try to improve the accuracy for our model

```
knn_model = KNeighborsClassifier(n_neighbors = 8, metric = 'minkowski', p = 2)
knn_model.fit(x_train, y_train)

KNeighborsClassifier(n_neighbors=8)
```

```
y_predict =knn_model.predict(x_test)
```

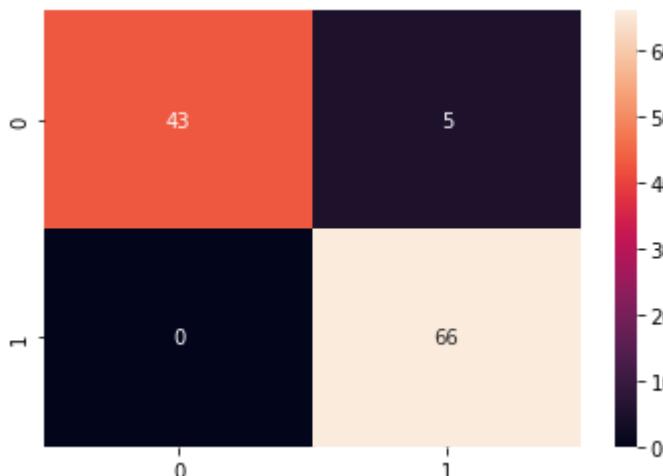
```
y_predict
```

```
array([0., 1., 1., 1., 0., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1.,
       1., 1., 1., 0., 1., 1., 1., 1., 0., 0., 1., 0., 0., 0., 0., 1., 0.,
       1., 1., 0., 1., 0., 1., 1., 0., 1., 1., 0., 0., 1., 0., 1., 0., 1.,
       1., 1., 1., 0., 0., 0., 1., 0., 0., 0., 1., 1., 1., 1., 1., 1., 1.,
       1., 0., 1., 0., 1., 1., 1., 1., 0., 1., 0., 1., 0., 1., 0., 0., 0.,
       1., 0., 1., 0., 0., 0., 1., 1., 0., 0., 1., 1., 1., 1., 1., 1., 0.,
       1., 1., 0., 0., 1., 0., 1., 0., 1., 0., 1., 0., 1., 0., 1., 0., 0., 0.])
```

```
cm = confusion_matrix(y_test,y_predict)
```

```
sns.heatmap(cm ,annot=True)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f0664016a50>
```



```
from sklearn import metrics
print("Train set Accuracy: ", metrics.accuracy_score(y_train, neigh.predict(x_train)))
print("Test set Accuracy: ", metrics.accuracy_score(y_test, yhat))

Train set Accuracy:  0.9384615384615385
```

```
Test set Accuracy: 0.9385964912280702
```

```
print(classification_report(y_test,y_predict))
```

	precision	recall	f1-score	support
0.0	1.00	0.90	0.95	48
1.0	0.93	1.00	0.96	66
accuracy			0.96	114
macro avg	0.96	0.95	0.95	114
weighted avg	0.96	0.96	0.96	114

As we can see the Test set accuracy for the model has increased from 94% to 96% when taken the best value for nearest neighbours

▼ TRAINING THE MODEL USING SVM

```
from sklearn.svm import SVC
from sklearn.metrics import classification_report , confusion_matrix
```

```
from sklearn.svm import SVC
svm_model = SVC(kernel = 'linear', random_state = 0)
svm_model.fit(x_train, y_train)
```

```
SVC(kernel='linear', random_state=0)
```

```
y_predict =svm_model.predict(x_test)
```

```
y_predict
```

```
array([0., 1., 1., 1., 0., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1.,
       1., 1., 0., 1., 1., 1., 1., 0., 0., 1., 0., 0., 0., 0., 1., 0.,
       1., 1., 0., 1., 1., 1., 1., 0., 1., 1., 1., 0., 1., 0., 1., 1.,
       1., 1., 1., 0., 0., 0., 1., 0., 0., 0., 1., 1., 1., 1., 1., 1.,
       1., 0., 1., 0., 1., 1., 1., 0., 1., 0., 0., 0., 1., 0., 0., 0.,
       1., 0., 1., 0., 0., 0., 1., 1., 0., 0., 1., 1., 1., 1., 1., 0.,
       1., 1., 0., 0., 0., 1., 0., 1., 0., 0., 1., 1., 1., 1., 1., 0.,
```

```
cm = confusion_matrix(y_test,y_predict)
sns.heatmap(cm ,annot=True)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f06635fee10>
```

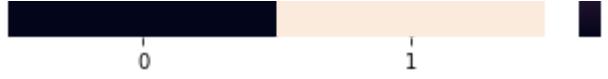


CALCULATING F1 SCORE



```
from sklearn.metrics import f1_score  
f1_score(y_test, yhat, average='weighted')
```

```
0.9377278780169241
```



✓ 0s completed at 12:46

