

Importing Libraries

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

import warnings
warnings.filterwarnings('ignore')
```

Import Dataset

```
In [2]: df = pd.read_csv(r"E:\FSDS With GEN AI_NIT\Breast_cancer_data.csv")
df.head()
```

```
Out[2]:
```

	mean_radius	mean_texture	mean_perimeter	mean_area	mean_smoothness	diagnosis
0	17.99	10.38	122.80	1001.0	0.11840	0
1	20.57	17.77	132.90	1326.0	0.08474	0
2	19.69	21.25	130.00	1203.0	0.10960	0
3	11.42	20.38	77.58	386.1	0.14250	0
4	20.29	14.34	135.10	1297.0	0.10030	0

```
In [3]: df.shape
```

```
Out[3]: (569, 6)
```

View Summary of the dataset

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   mean_radius      569 non-null   float64
1   mean_texture     569 non-null   float64
2   mean_perimeter   569 non-null   float64
3   mean_area        569 non-null   float64
4   mean_smoothness  569 non-null   float64
5   diagnosis        569 non-null   int64
dtypes: float64(5), int64(1)
memory usage: 26.8 KB
```

Check the Summary of the Tagert variable

```
In [5]: df.count()
```

```
Out[5]: mean_radius      569
mean_texture      569
mean_perimeter    569
mean_area         569
mean_smoothness   569
diagnosis         569
dtype: int64
```

```
In [6]: df['diagnosis'].value_counts()
```

```
Out[6]: diagnosis
1      357
0      212
Name: count, dtype: int64
```

- The target variable is diagnosis. It contains 2 values - 0 and 1.
-

0 is for Negative prediction and 1 for Positive predictio.+

We can see that the problem is binary classification task.

Declare Feature Vector and Target Variable

```
In [8]: x = df[['mean_radius','mean_texture','mean_perimeter','mean_area','mean_smoothness']]  
y = df['diagnosis']
```

```
In [11]: x
```

```
Out[11]:
```

	mean_radius	mean_texture	mean_perimeter	mean_area	mean_smoothness
0	17.99	10.38	122.80	1001.0	0.11840
1	20.57	17.77	132.90	1326.0	0.08474
2	19.69	21.25	130.00	1203.0	0.10960
3	11.42	20.38	77.58	386.1	0.14250
4	20.29	14.34	135.10	1297.0	0.10030
...
564	21.56	22.39	142.00	1479.0	0.11100
565	20.13	28.25	131.20	1261.0	0.09780
566	16.60	28.08	108.30	858.1	0.08455
567	20.60	29.33	140.10	1265.0	0.11780
568	7.76	24.54	47.92	181.0	0.05263

569 rows × 5 columns

```
In [12]: y
```

```
Out[12]: 0      0
          1      0
          2      0
          3      0
          4      0
          ..
        564      0
        565      0
        566      0
        567      0
        568      1
Name: diagnosis, Length: 569, dtype: int64
```

Split the data into train and test split

```
In [9]: from sklearn.model_selection import train_test_split

X_train,X_test,y_train,y_test = train_test_split(X,y,test_size = 0.3, random_state = 0)
```

Built a LGBM classifier Model

```
In [11]: import lightgbm as lgb

cls = lgb.LGBMClassifier()
cls.fit(X_train,y_train)
```

[illegible]

localhost:8888/doc/tree/FSDS 2nd Jan 2025/Class Folder/4. April/Breast Cancer prediction Analysis Using LGBM.ipynb

```
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
```

```
Out[11]: ▾ LGBMClassifier
          LGBMClassifier()
```

Model Prediction based on test data

```
In [12]: y_pred = cls.predict(X_test)
```

Compute Accuracy

```
In [13]: from sklearn.metrics import accuracy_score
          accuracy=accuracy_score(y_pred, y_test)
          print('LightGBM Model accuracy score: {0:0.4f}'.format(accuracy_score(y_test, y_pred)))
```

LightGBM Model accuracy score: 0.9298

Compare the train and test score accuracy

```
In [16]: y_pred_train = cls.predict(X_train)

print('Traing set accuracy is: {0:0.4f}'.format(accuracy_score(y_train,y_pred_train)))
```

Traing set accuracy is: 1.0000

Check for Overfitting

```
In [18]: print('Training set accuracy (Bias): {0:0.4f}'.format(cls.score(X_train, y_train)))
print('Testing set accuracy (Variance): {0:0.4f}'.format(cls.score(X_test, y_test)))
```

Training set accuracy (Bias): 1.0000

Testing set accuracy (Variance): 0.9298

The training and test set accuracy are quite comparable. So, we cannot say there is overfitting.

Compute Confusion Matrix

```
In [21]: from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test,y_pred)
print('Confusion Matrix\n\n', cm)
print('\n True Positive (TP):', cm[0,0])
print('\n True Negative (TN):', cm[1,1])
print('\n False Positive (FP):', cm[0,1])
print('\n False Negative (FN):', cm[1,0])
```


Confusion Matrix

```
[[ 55   8]
 [  4 104]]
```

True Positive (TP): 55

True Negative (TN): 104

False Positive (FP): 8

False Negative (FN): 4

Classification Metrics

```
In [24]: from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.93	0.87	0.90	63
1	0.93	0.96	0.95	108
accuracy			0.93	171
macro avg	0.93	0.92	0.92	171
weighted avg	0.93	0.93	0.93	171

visualize confusion matrix with seaborn heatmap

```
In [22]: cm_matrix = pd.DataFrame(data=cm, columns=['Actual Positive:1', 'Actual Negative:0'],
                                index=['Predict Positive:1', 'Predict Negative:0'])

sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='YlGnBu')
```

Out[22]: <Axes: >

