# Accuracy and Confusion Matrix

In [1]:
```python
import pandas as pd
import numpy as np
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report,roc_auc_score,confusion_matrix
```

In [2]:
```python
from  sklearn.datasets import load_breast_cancer
```

In [3]:
```python
Cancer = load_breast_cancer()
```

```
In [4]:    1  print(load_breast_cancer()['DESCR'])
```

.. _breast_cancer_dataset:

Breast cancer wisconsin (diagnostic) dataset
--------------------------------------------

**Data Set Characteristics:**

    :Number of Instances: 569

    :Number of Attributes: 30 numeric, predictive attributes and the class

    :Attribute Information:
        - radius (mean of distances from center to points on the perimeter)
        - texture (standard deviation of gray-scale values)
        - perimeter
        - area
        - smoothness (local variation in radius lengths)
        - compactness (perimeter^2 / area - 1.0)
        - concavity (severity of concave portions of the contour)
        - concave points (number of concave portions of the contour)
        - symmetry
        - fractal dimension ("coastline approximation" - 1)

        The mean, standard error, and "worst" or largest (mean of the three
        worst/largest values) of these features were computed for each image,
        resulting in 30 features.  For instance, field 0 is Mean Radius, field
        10 is Radius SE, field 20 is Worst Radius.

        - class:
                - WDBC-Malignant
                - WDBC-Benign

    :Summary Statistics:

    ===================================== ====== ======
                                           Min    Max
    ===================================== ====== ======
    radius (mean):                         6.981  28.11
    texture (mean):                        9.71   39.28
```

```
        perimeter (mean):                        43.79  188.5
        area (mean):                             143.5  2501.0
        smoothness (mean):                       0.053  0.163
        compactness (mean):                      0.019  0.345
        concavity (mean):                        0.0    0.427
        concave points (mean):                   0.0    0.201
        symmetry (mean):                         0.106  0.304
        fractal dimension (mean):                0.05   0.097
        radius (standard error):                 0.112  2.873
        texture (standard error):                0.36   4.885
        perimeter (standard error):              0.757  21.98
        area (standard error):                   6.802  542.2
        smoothness (standard error):             0.002  0.031
        compactness (standard error):            0.002  0.135
        concavity (standard error):              0.0    0.396
        concave points (standard error):         0.0    0.053
        symmetry (standard error):               0.008  0.079
        fractal dimension (standard error):      0.001  0.03
        radius (worst):                          7.93   36.04
        texture (worst):                         12.02  49.54
        perimeter (worst):                       50.41  251.2
        area (worst):                            185.2  4254.0
        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 (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.

```
In [5]:    1  y = Cancer.target
```

In [6]:
```python
1  X = pd.DataFrame(Cancer.data, columns=Cancer.feature_names)
2  X.head()
```

Out[6]:

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | mean fractal dimension | ... | worst radius | worst texture | worst perimeter | worst area |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.3001 | 0.14710 | 0.2419 | 0.07871 | ... | 25.38 | 17.33 | 184.60 | 2019.0 |
| 1 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.0869 | 0.07017 | 0.1812 | 0.05667 | ... | 24.99 | 23.41 | 158.80 | 1956.0 |
| 2 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.1974 | 0.12790 | 0.2069 | 0.05999 | ... | 23.57 | 25.53 | 152.50 | 1709.0 |
| 3 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.2414 | 0.10520 | 0.2597 | 0.09744 | ... | 14.91 | 26.50 | 98.87 | 567.7 |
| 4 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.1980 | 0.10430 | 0.1809 | 0.05883 | ... | 22.54 | 16.67 | 152.20 | 1575.0 |

5 rows × 30 columns

In [7]:
```python
1  from sklearn.metrics import confusion_matrix,accuracy_score
```

In [8]:
```python
1  clf = KNeighborsClassifier()
2  clf.fit(X,y)
```

Out[8]: KNeighborsClassifier()

In [9]:    1  y

Out[9]: array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
        0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0,
        1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0,
        1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1,
        1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0,
        0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1,
        1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0,
        0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0,
        1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1,
        1, 1, 0, 1, 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, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0,
        0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0,
        0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 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, 0, 1, 1,
        1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1,
        1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
        1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
        1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1,
        1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 0, 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, 0, 0, 0, 1])

In [14]:
```python
1  print(accuracy_score(y,clf.predict(X))*100)
2
3  print(confusion_matrix(y,clf.predict(X)))
4  #94.72% -> Accurately classified
5  #5.28% -> Misclassified
6
7
```

```
94.72759226713534
[[191  21]
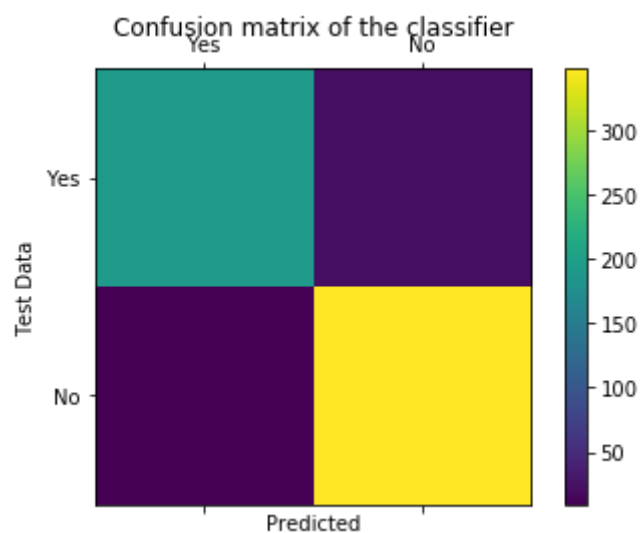 [  9 348]]
```

In [15]:
```python
1  (191+348)/X.shape[0]
```

Out[15]: 0.9472759226713533

In [16]:
```python
import matplotlib.pyplot as plt
plt.figure(figsize=(25, 25))
labels = ['Yes','No']
cm = confusion_matrix(y, clf.predict(X))
print(cm)
fig = plt.figure()
ax = fig.add_subplot(111)
cax = ax.matshow(cm)
plt.title('Confusion matrix of the classifier')
fig.colorbar(cax)
ax.set_xticklabels([''] + labels)
ax.set_yticklabels([''] + labels)
plt.xlabel('Predicted')
plt.ylabel('Test Data')
```

```
[[191  21]
 [  9 348]]
```

Out[16]:  Text(0,0.5,'Test Data')

```
<Figure size 1800x1800 with 0 Axes>
```



# RoC_AUC

Roc - > Reciever Operator Char.

Auc - > Area Under the Curve

**True Positivity rate(Sensitivity) = Tp/(Fn+Tp)**

**True Negative rate(Specifity) = Tn/(Fp+Tn)**

**False Positive rate = Fp/(Tn+Fp)**

In [18]:
```python
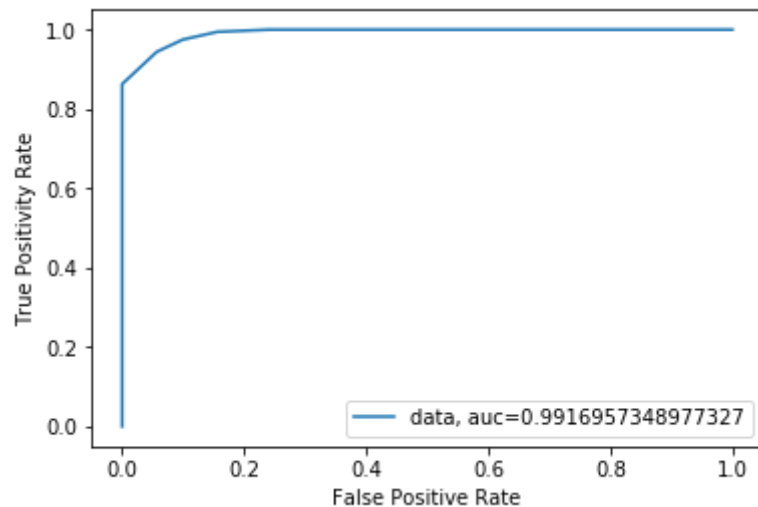from sklearn import metrics
```

In [20]:
```python
y_pred_proba = clf.predict_proba(X)[::,1]
fpr, tpr, _ = metrics.roc_curve(y,  y_pred_proba)
auc = metrics.roc_auc_score(y, y_pred_proba)
```

```
In [62]:    1  y_pred_proba = clf.predict_proba(X)[::,1]
            2  fpr, tpr, _ = metrics.roc_curve(y,  y_pred_proba)
            3  auc = metrics.roc_auc_score(y, y_pred_proba)
            4  plt.plot(fpr,tpr,label="data, auc="+str(auc))
            5  plt.legend(loc=4)
            6  plt.xlabel('False Positive Rate')
            7  plt.ylabel('True Positivity Rate')
            8  plt.show()
```



```
In [63]:    1  from sklearn.metrics import roc_auc_score
```

```
In [64]:    1  roc_auc_score(clf.predict(X),y)
```

Out[64]:  0.9490447154471544

## Classification Report => Accuracy,Recall,Precision,F1 Score(Combination of precision and Recall)

```
In [65]:    1  from sklearn.metrics import classification_report
```

```
In [66]:    1 print(classification_report(clf.predict(X),y))
```

```
              precision    recall  f1-score   support

           0       0.90      0.95      0.93       200
           1       0.97      0.94      0.96       369

    accuracy                           0.95       569
   macro avg       0.94      0.95      0.94       569
weighted avg       0.95      0.95      0.95       569
```

```
In [ ]:    1
```

# Precision

**Precision = Tp/Tp+Fp**

*Ratio of correctly predicted positives to total predicted positives*

*What proportion of predicted positive values are actually positive*

*Model Effectiveness on the captured variation*

Example :

Business use case is (We want to decrease credit limit - > Target[Whether to decrease limit or not])) and our priority is customer satisfaction so we want to avoid all cases where we decrease limit of a customer that can afford the higher limit.

Classification model: 1 -> Can decrease limit 0 -> Cannot decrease limit

False Positive Case -> We decrease cr limit of a person who can afford it (Predicted : 1 actual : 0)

**Advantage of high precision:**

To decrease FP cases we increase threshold : 0.9 :: Higher Threshold of positive classification == Higher precision value

**Disadvantage of chasing a high precision with higher thresholds:**

1.) Since we have a high threshold of prob for default to decrease the credit limit (i.e. 90%) we give higher credit limits to people with 80% chance of default that means we are giving riskier loans and hence there is a higher chance of losing money.

# Recall

**recall = Tp/(Tp+Fn)**

**Ratio of all correctly predicted positive values to all the positive values**

**What proportion of actual positives are predicted positives**

**Of how well your data is engineered and also a combination of how well you've selected your model on the basis of your data(Recall depends a lot on your data)**

# F-1 Score

**Harmonic Mean of Precision and Recall**

F1 = 2*(Precision X recall)/(Precision+recall)

**Weighted Average of acc, Precision and Recall**