# Lab2 Report

**Mohamed Elmaghraby Mohamed**
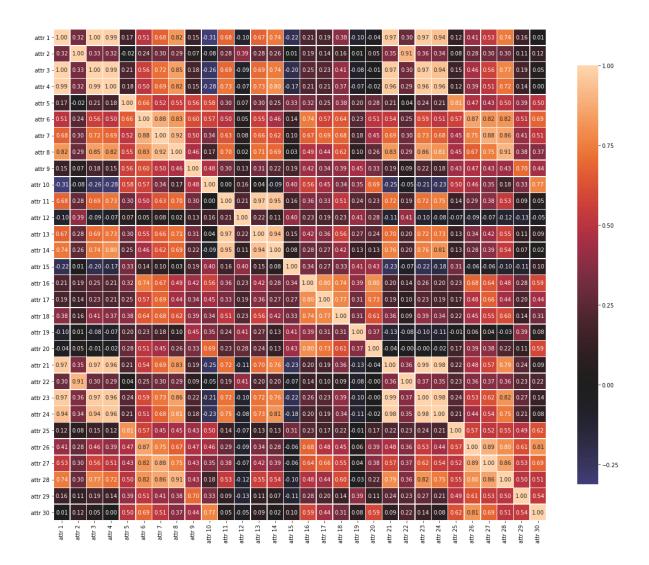
**55**

# Prepare data:

- Prepare column names.
- Import breast-cancer-Wisconsin data.
- Merge both data sets.
- Explore data :
    - Number of instances = 569
    - number of attributes = 32
    - class names :  ['M',  'B']
    - number of classes:  2

# Data Exploration:

- **Plot Box plot for the data**
    - **Conclusion :**
        - most of the values of the attributes are around zero
        - attr 4 and attr 24: has bigger and spread values
        - attr 14: has many outliers
- **Correlation matrix**
    - Use `.corr(method = 'Pearson')` to get the Pearson's correlation between each 2 attributes.
    - **Conclusion**
        - some of the attributes have a high positive correlation
        - most of the attributes have a medium positive correlation
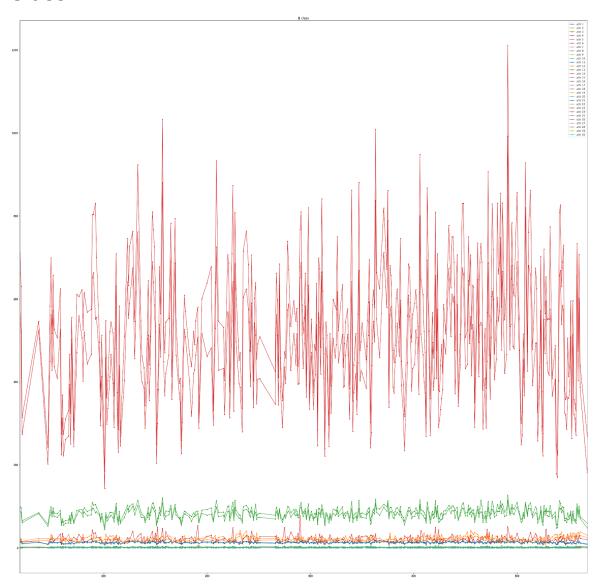        - a small set of attributes has a low positive correlation

- ■ Examples :
  - ● attrs (21, 23, 24) vs attrs (1, 3, 4)  has very strong correlation
  - ● attrs( 25, 26, 27, 28, 29, 30) has very strong correlation
  - ● attrs (15, 16, …, 20) vs attrs (1, 2 , 3 , 4) has very weak correlation
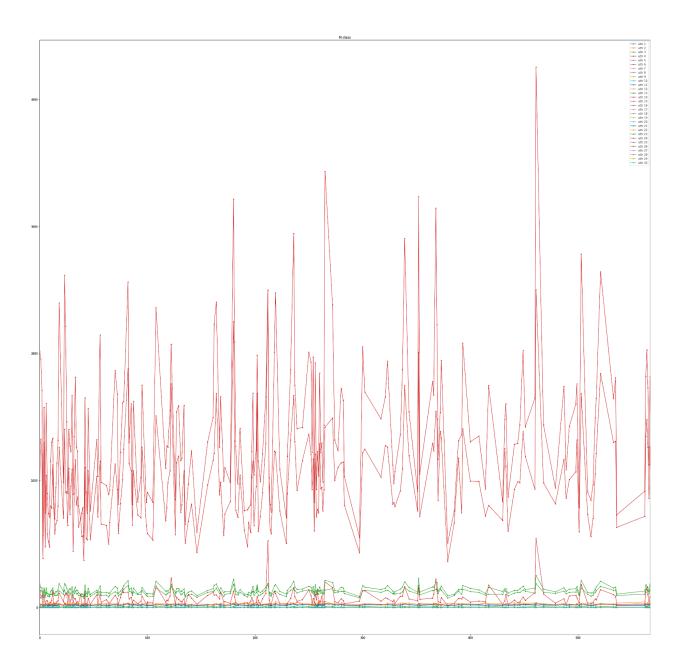- ○ Draw the correlation using heatmap.

- **Line plot**
  - I plotted a line plot for each class independently
  - **Conclusion**
    - by comparing both attr 3 and attr 23 seems to has a higher value in class 'B' than in class' (colored by green)
    - same for attrs (1, 2, 21, 22, 11, 12) ( colored by blue and orange)

## Class B



## Class M

M class

# Data Preprocessing:

- Split the data int train and test sets using StratifiedShuffleSplit
- **Normalization:**
  - **Z-score:**
    - Calculate z-score using `zscore(x) from scipy library`
    - **Conclusion**:
      - The result data spread on less scale, so values are more close to each other.
      - A value is exactly equal to the mean of all the values of the feature, it will be normalized to 0. If it is below the mean, it will be a negative number, and if it is above the mean it will be a positive number.
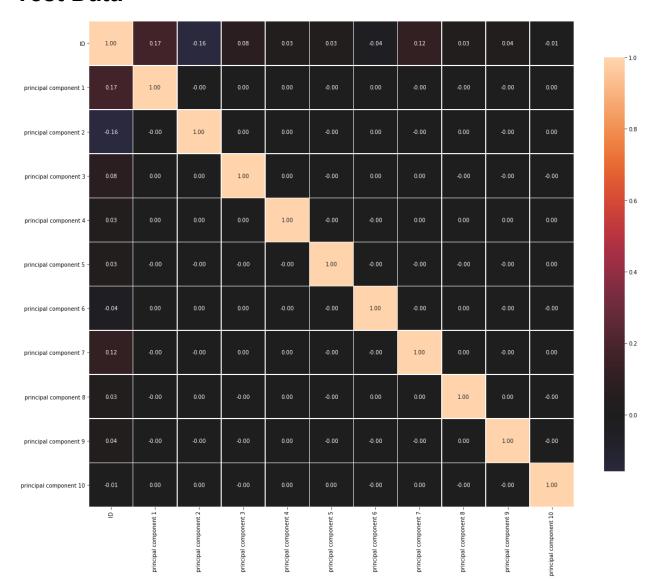
- **Dimensionality reduction**
  - Feature Projection:
    - Use `PCA from the scikit-learn` library seeking to sum of variance ratio we need of 0.95.
    - Number of PCA components = 10
    - Plot correlation matrix for result data.
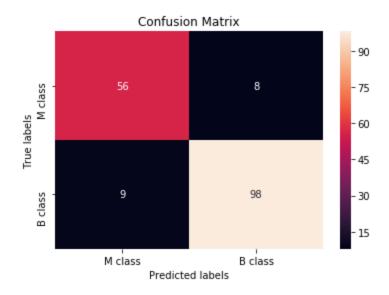
# Train Data

## Test Data



## ■ Conclusion

- Principal component analysis convert a set of observations of correlated variables into a set of values of linearly uncorrelated variables called principal components
- less than half of the attributes can describe the data with 0.95 percent.

# Classification

- **Decision tree**
  - Used GridSearchCV for tunning the depth parameter among the given values [ 5, 7, 10, 15, 20, 30] and get 5 as the best value.
  - Train the model with DecisionTreeClassifier from Sk-learn
  - Compute precision, recall, F-score using precision_recall_fscore_support from SK-learn
    - precision: 0.9193548387096774 for "M" class , 0.9357798165137615 for "B" class
    - recall: 0.890625 for "M" class , 0.9532710280373832 for "B" class
    - fscore: 0.9047619047619047 for "M" class , 0.9444444444444445 for "B" class
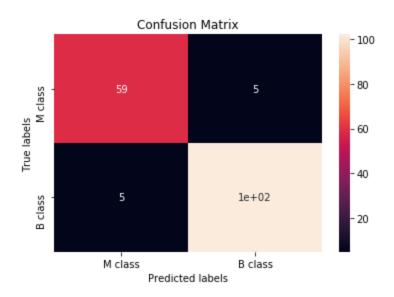  - Compute Confusion matrix using confusion_matrix from SK-learn



Confusion Matrix

**Note**

The same algorithm tested for depth = 15, 50 and less value for Evaluation matrices are obtained, so parameter tunning helped to get better accuracy.

## ● AdaBoost Classifier
  - ○ Used GridSearchCV for tunning
    - ■ n_estimators in range [10, 50, 100, 200, 500]
    - ■ learning_rate in range [0.25, 0.5, 1, 1.5, 2, 3]
    - ■ Best value of learning_rate = 1.5
    - ■ Best value of n_estimators = 200
  - ○ Train the model with AdaBoostClassifier from Sk-learn
  - ○ Compute precision, recall, F-score using precision_recall_fscore_support from SK-learn
    - ■ precision: 0.9375 for "M" class , 0.9626168224299065 for "B" class
    - ■ recall: 0.9375 for "M" class , 0.9626168224299065 for "B" class
    - ■ fscore: 0.9375 for "M" class , 0.9626168224299065 for "B" class
  - ○ Compute Confusion matrix using confusion_matrix from SK-learn



Confusion Matrix

**Note**

The same algorithm tested for learning rate = 0.5, 2..5 and n_estimtors = 50, 300 and less value for Evaluation matrices are obtained, so parameter tunning helped to get better accuracy.

- **Random Forest Classifier**
  - Used GridSearchCV for tunning
    - depth in range [5, 10, 15, 20, 50]
    - n_estimators in range [10, 50, 100, 200, 500]
    - Best value of depth = 10
    - Best value of n_estimators = 50
  - Train the model with RandomForestClassifier from Sk-learn
  - Compute precision, recall, F-score using precision_recall_fscore_support from SK-learn
    - precision: 0.9016393442622951 for "M" class , 0.9181818181818182 for "B" class
    - recall: 0.859375 for "M" class , 0.9439252336448598 for "B" class
    - fscore: 0.88 for "M" class , 0.9308755760368663 for "B" class
  - Compute Confusion matrix using confusion_matrix from SK-learn



Confusion Matrix

**Note**

The same algorithm tested for depth = 2, 20 and n_estimtors = 50, 200 and less value for Evaluation matrices are obtained, so parameter tunning helped to get better accuracy.

- ## Conclusion
  - All the models have very good values for the precession, recall, and F-score
  - The best model is the AdaBoost classifier as the number of wrongly classified samples is less = 9