

# Breast Cancer Prediction using Principal Component Analysis with Logistic Regression

Mohammad Kaosain Akbar

*Concordia Institute for Information Systems Engineering (CIISE)  
Concordia University, Montreal, Quebec, Canada.*

Date of Submission: 15-10-2022

Date of Acceptance: 31-10-2022

**ABSTRACT:** Breast Cancer is considered as the major and most common form of cancer among women. Breast cancer is the second source of death causes by cancer, with first being the lung cancer. To tackle this cancer, rigorous efforts are constantly being given by scientists across the globe. Besides, the field of Machine Learning and Data Mining has made significant progress over the years for extracting and gathering valuable information, even from the most complex data sources. Based on the information extracted from data, the Machine Learning model is also capable of performing certain degree of prediction, classification, and clustering. In this paper, we have explored the relation between diagnosis of breast cancer with multiple number of attributes of a dataset. We have used a supervised learning classification algorithm called Logistic Regression for predicting the existence of breast cancer, based on five different attributes of X-ray images dataset. The dataset was obtained from the University of Wisconsin Hospitals, Madison from Dr. William H. Wolberg. Prior to applying the Machine Learning algorithm, we performed a dimensionality reduction technique over the data, known as Principal Component Analysis (PCA). Then, we analysed the performance of the Machine Learning algorithm based on the accuracy, precision, recall, f1-score, and support. Additionally, we also analysed the above scores generated by the model without the use of the PCA analysis.

**KEYWORDS:** Logistic Regression, Machine Learning, Principal Component Analysis, Breast Cancer, Prediction.

## I. INTRODUCTION

Cancer is a horrible disease which occurs when there is abnormal growth of cells within our

body. These cells can spread across the body, destroying normal body cells and can form a mass commonly known as tumour [1]. Breast Cancer is developed from abnormal cells growth in breast and signs mostly include change in breast shape and skin and lump in the breast. Although this kind of cancer is invasive cancer in women, breast cancer can also occur in men [2]. Based on the situation of cancer, Doctors arrange treatment plans, but any misdiagnoses can undoubtedly lead patient to loss significant amount of curing time which in turn may end up in dire consequences. Therefore, it is very important to diagnose and predict breast cancer at early stage

Our study chooses Logistic Regression in order to work predict diagnosis of breast cancer because this algorithm is generally used to predict the probability of a target value. Often the target value is dichotomous, which means that the target value will have two outcomes. In our case that would be either the targeted patient is diagnosed with cancer, or the patient is not diagnosed with cancer. We used Python programming language as the essential tool for developing the Machine Learning model and predicting patients being diagnosed with breast cancer or not based on the dataset generated from X-ray images. The Breast Cancer dataset which we used to train and test our model was obtained from a reputed data repository platform called "Kaggle".

We first performed PCA analysis on the breast cancer dataset. PCA analysis reduced the dimension of the data and then the data was used to train the Machine Learning model which then we calculated scores based on the predictions based by the model. Moreover, we also trained the Machine Learning model without reducing the data using PCA analysis. Then we compared the two sets of scores, generated by two Logistic Regression

Models. results show dynamic behaviour under various operating and environmental conditions and demonstrate advantages of adaptive control over the non-adaptive type.

## II. DATASET DESCRIPTION

The Breast Cancer dataset was obtained from a well known online data repository named “Kaggle”. This dataset is a collection of data from different patients who performed X-rays to check whether they were infected with Breast Cancer. The dataset has six attributes; five attributes represent different values obtained from the X-ray images such as mean of radius, mean of the texture of the X-ray, etcetera. The final attribute indicates whether the patient is diagnosed with Breast Cancer or not. The provided information by the breast cancer datasets are as follows:

- mean\_radius : Mean radius of the lump
  - mean\_texture: Mean texture of the X-ray image.
  - mean\_perimeter: Mean perimeter of the lump
  - mean\_area: Mean area of the lump
  - mean\_smoothness: Mean of smoothness of the image
  - diagnosis: Whether patient is diagnosed with Cancer or not.
- 0 represents patients not being infected with breast cancer.  
1 represents patients who were diagnosed with breast cancer.

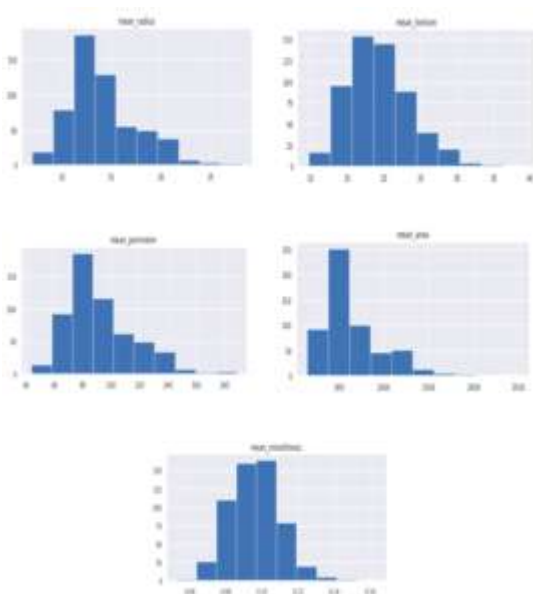


Figure 1: Histogram Representation of the attributes which is used for training the Machine Learning model.

Apart from the diagnosis attribute, we will use all the remaining attributes to train our machine learning models based on Logistic Regression. The data was gathered from the University of Wisconsin Hospitals. The data was a public dataset, and no additional charges were made while downloading the data. Figure 1 represents the histogram representation of the attributes. The figure shows the different values of the attributes and their frequency of occurrence within the breast

## III. PRINCIPAL COMPONENT ANALYSIS

### A. Theoretical Description

Principal Component Analysis is a statistical technique which is used which is greatly used to reduce the dimensions of the data. In addition to reduce dimensions, PCA also emphasizes on variation and highlights strong patterns in a dataset so data can be easily explored and visualized. Although, PCA reduces the dimension of the data, PCA tries to retain as much information possible from the original data. PCA generates Principal Components, and the first Principal component represents the larger variance of the data which means that this component accounts for most of the variability in the data. Likewise, the second Principal component represents the second most significant variance and so on. In the following, we describe the steps of performing the PCA analysis:

- First, we try to center the original data by computing the centered data matrix  $Y = X - \bar{X}$  by subtracting off-column means.
- Next, for the second step, using the centered data, we compute the covariance matrix  $S$ , represented by  $p \times p$ .

$$s = \frac{1}{n-1} Y'Y.$$

- Then, for the covariance matrix  $S$ , we compute the eigenvectors and eigenvalues using eigendecomposition

$$s = \Lambda \Lambda' = \sum_{i=1}^p \lambda_j a_j' a_j.$$

Here:

-  $\Lambda$  is a  $p \times p$  orthogonal matrix (i.e.  $\Lambda' \Lambda = I$ ) whose columns  $a_j = (a_{j1}; a_{j2}; \dots; a_{jp})$  are the eigenvectors of  $S$ .

-  $\Lambda = \text{diag}(\lambda_1; \lambda_2; \dots; \lambda_p)$  is a  $p \times p$  diagonal matrix whose elements are the eigenvalues of  $S$  arranged in decreasing order.

For the last step, we compute the transformed data matrix which is  $Z = Y A$  and this matrix is of size  $n \times p$ .

$$Z = (z'_{11}, z'_{12}, \dots, z'_{1p}, \dots, z'_{n1}, z'_{n2}, \dots, z'_{np}) = \begin{bmatrix} z'_{11} & z'_{12} & \dots & z'_{1p} \\ z'_{21} & z'_{22} & \dots & z'_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ z'_{n1} & z'_{n2} & \dots & z'_{np} \end{bmatrix}$$

It transformed data matrix contains the coordinates of the original data in the new coordinate system defined by the Principal Components. The rows of  $Z$  matrix correspond to observations  $Z_i = A' (x_i - \bar{x})$ , while its columns correspond to PC scores [3].

#### B. Application of PCA in reduction of Breast Cancer data dimensions

The data was obtained from Kaggle and was created for "AI for Social Good: Women Coders' Bootcamp". The data do not have time parameters. It only consists of X-ray image features and output of patient being diagnosed with Cancer or not

#	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

Figure 2: Datatype of the individual attributes along with the non-null count.

From figure 2, we can see that all attributes of Breast Cancer dataset have datatype of float64 and the dataset do not have any null values.

#### 1) Box Plot

Initially, we normalized the data so that values of numeric columns of the dataset gets into same scale.

$$\frac{\text{value} - \mu}{\sigma}$$

Then, we constructed the side-by-side box plots of the five attributes of the dataset.

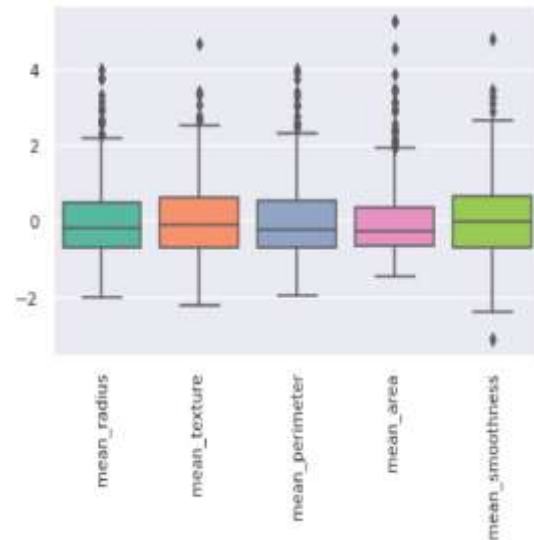


Figure 3: Box plot of the dataset

Figure 3 shows the box plot of the dataset. From the boxplot, we can see that all the attributes have certain numbers of outliers. These outliers occurred because health and medical data often contains abnormal patient condition, instrumentation errors or even recording errors [4].

#### 2) Covariance Matrix

After the data is centred, we compute the Covariance Matrix. Covariance Matrix helps us to measure of how variables change with respect to each other. It is positive when variables tend to show similar behaviour and negative otherwise.

From figure 4, we can see the covariance matrix of the dataset. It can be observed that the values at diagonals of the matrix are 1. since along the diagonal, variables are same with respect to one another. It is also observed that all the values of the covariance matrix are positive, lying between 0 and 1.

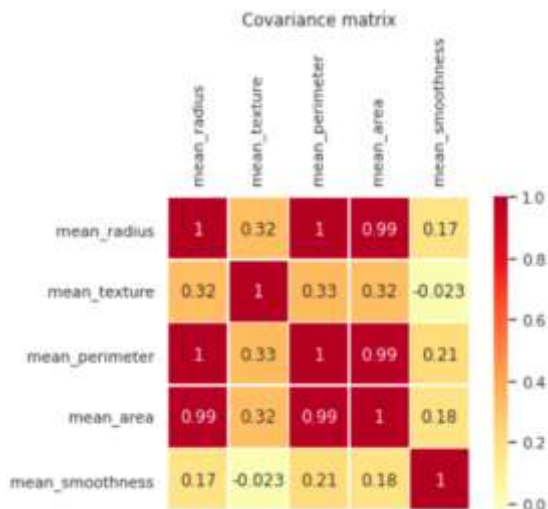


Figure 4: Covariance Matrix of the dataset

### 3) Pair-plot

Pair-plots usually plots a pairwise relationships in adataset [5]. Pair-plot represents the distribution of singlevariables as well as correlation between two variables. Fromfigure 5, we can see that certain attributes shows linearrelationships with other attributes of the dataset which meansthey are attributes are highlycorrelated. Additionally, fromthe diagonals of the Pair-plot we can see the histograms anddetermine the skewness.

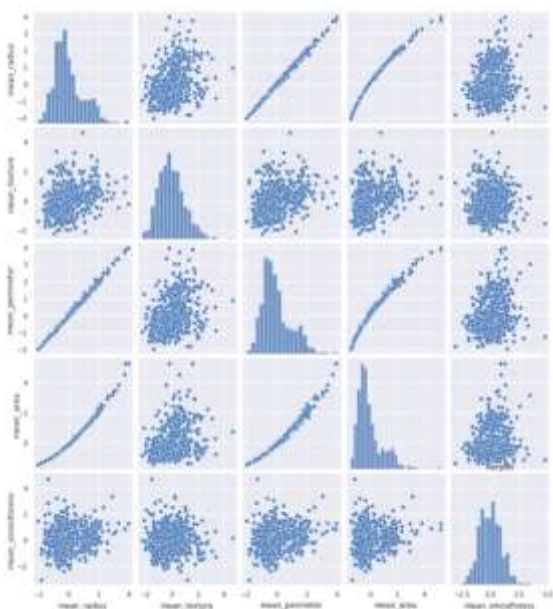


Figure 5: Pair-plot of the breast cancer dataset

### 4) Scree plot and Pareto Chart

Next, we applied Principal Component analysis on thedataset. In our Python programming

language, we used thebuilt-in framework library called the scikit-learn framework.Then we obtain the Principal Components.

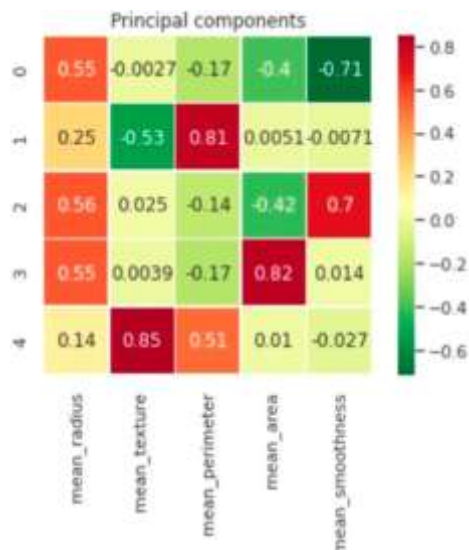


Figure 6: Principal Components generated after performing principal component analysis on the dataset.

In order to decide how many Principal Components shouldwe select; we first need to plot the Scree plot and Pareto Chart.The Scree plots provides us with explained variance ofindividual Principal Components that is which is thepercentage of variance counted for by the  $j^{th}$  PrincipalComponent.

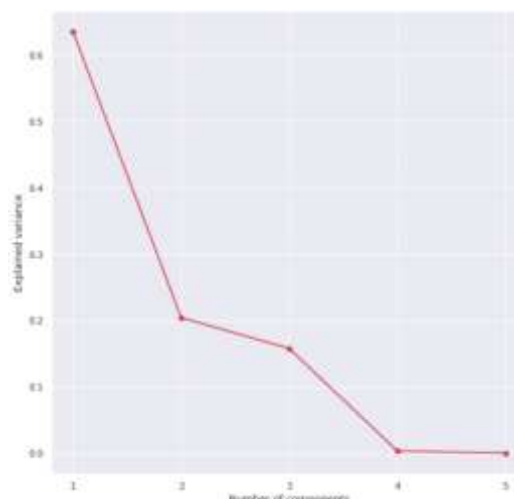


Figure 7: Scree Plot



Component. Explained variance is given by:

$$l_j = \frac{\lambda_j}{\sum_{j=1}^p \lambda_j} * 100\%, \forall j \in [1, p]$$

From our scree plot, we can see that we can select the value of “r” as 2, based on the elbow generated from the beginning of horizontal point 2.

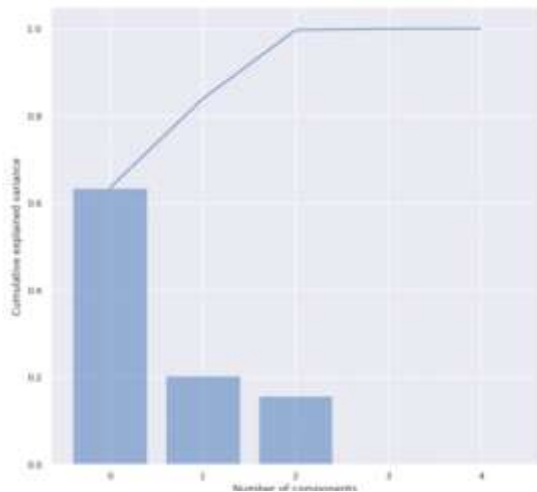


Figure 8: Pareto Chart

Generally, the rule is to be considered retain several components that represent at least 80% of the data variance. From our Pareto chart, we can see that almost 82% of the total variance can be explained by two of the initial Principal Components that was generated right after the Principal Component Analysis.

#### 5) Scatter Plot

In Figure 9, we can see the scatter plot from which we can find out which attributes have similar contribution as that of the computed Principal Components. Considering the parameters mean\_smoothness which is situated on the top left corner of the plot, we can deduce that the attribute of mean smoothness has high contribution on Principal Component 2 (A2).

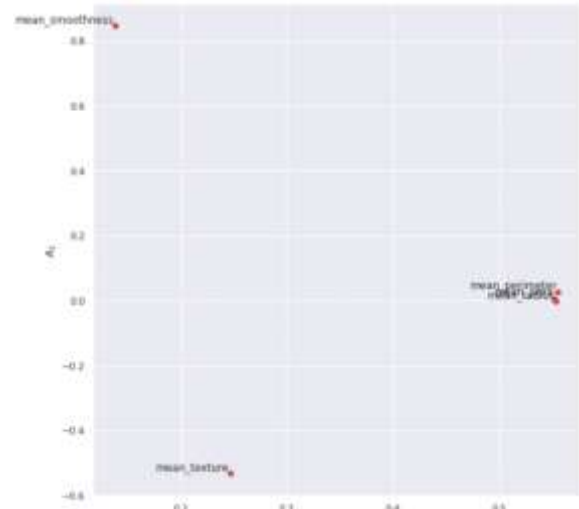


Figure 9: Scatter Plot

#### 6) Biplot

From biplot, for each observation, we can see principal component coefficient and principal component scores together. Points on the biplot represent observations whereas the lines on the biplot are the attributes of the dataset. From the length and direction of the vectors, we can deduce how each attribute contributes towards the individual PCs. Considering mean\_texture, this attribute contributes more to PC2 than mean\_smoothness attribute, because mean\_texture is closer to the PC2 than mean\_smoothness.

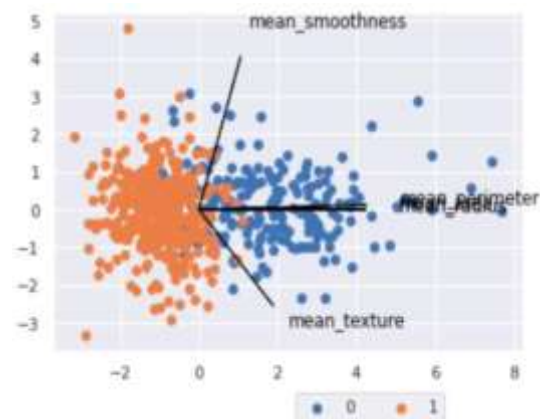


Figure 10: Biplot

### IV. MACHINE LEARNING ALGORITHM: LOGISTIC REGRESSION

Dealing with Cancer detection data is a very sensitive issue as, we need to properly identify the cancer patient efficiently. Since our dataset deals with classification problem, we find that for solving our problem, Logistic Regression would definitely be an ideal solution.

Logistic regression is basically a supervised classification algorithm. In a classification problem, the target variable (or output),  $y$ , can take only discrete values for given set of features (or inputs),  $X$ .

Contrary to popular belief, logistic regression is a regression model. The model builds a regression model to predict the probability that a given data entry belongs to the category numbered as "1". Just like Linear regression assumes that the data follows a linear function, Logistic regression models the data using the sigmoid function [7].

$$g(z) = \frac{1}{1 + e^{-z}}$$

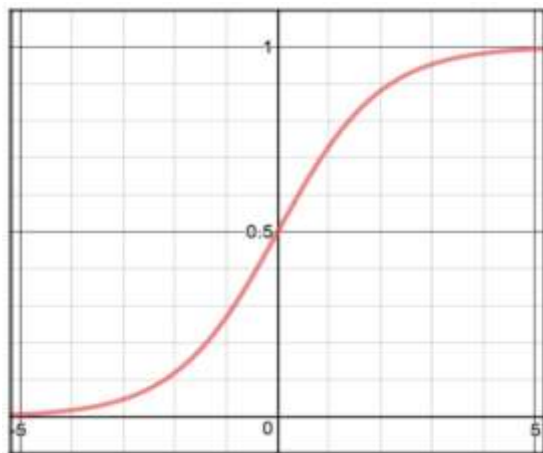


Figure 11: Graph of a sigmoid Function.

Logistic regression becomes a classification technique only when a decision threshold is brought into the picture. The setting of the threshold value is a very important aspect of Logistic regression and is dependent on the classification problem itself.

## V. EXPERIMENT WITH MACHINE LEARNING ALGORITHM

### A. Applying Logistic Regression on the original Breast Cancer Datasets

Initially, the original dataset is provided to the Logistic Regression. We use then use sklearn framework to perform Logistic Regression on the original dataset. We used train and test split operation to divide the data as training and testing set. After training the model with training set, we test our model using the testing set. Then we print the Accuracy Score, Classification report and Confusion matrix, once we are done testing our model.

	precision	recall	f1-score	support
0	0.91	0.86	0.89	36
1	0.94	0.96	0.95	78
accuracy			0.93	114
macro avg	0.92	0.91	0.92	114
weighted avg	0.93	0.93	0.93	114

Accuracy of ML Algorithms on the Full Data: 0.9298245614835088

Chart 1: Classification Report and accuracy score of the Machine Learning model on original dataset.

From the classification report, we can observe that our model generated a satisfactory output. The model had 91% precision in identifying patients who does not have breast cancer and 94% precision in identifying patients who is suffering from breast cancer. Both outcomes have an average f1-score and recall score of more than 90%. The accuracy of the training model is 93%.



Figure 12: Confusion Matrix for the ML algorithm on the Original Dataset.

From the confusion matrix generated by the model trained and tested using the original data, we can see that true negative score is 31 and true positive score is 75. Both of these scores are significantly higher than false positive and false negative scores, indicating that model was almost efficiently generating the expected output.

### B. Applying Logistic Regression on the data generated after Principal Component Analysis.

Next, we used the reduced data generated by the PCA to train our Machine Learning model based on Logistic Regression. Again, we used sklearn framework [6] to perform logistic operation but this time on the PCA generated dataset. We again used the similar train and test split operation to divide the data as training and testing set. After training the model with training set, we test our model using the testing set. Finally, we printed the Accuracy Score, Classification report and Confusion matrix, once we are done

testing our model with data generated after Principal Component Analysis.

	precision	recall	f1-score	support
0	0.91	0.95	0.93	41
1	0.97	0.95	0.96	73
accuracy			0.95	114
macro avg	0.94	0.95	0.94	114
weighted avg	0.95	0.95	0.95	114

Accuracy of PCA Data: 0.9473684210526315

Chart 2: Classification Report and accuracy score of the Machine Learning model on PCA generated dataset.

From the classification report, we can observe that our model had 91% precision in identifying patients who does not have breast cancer and 97% precision in identifying patients who is suffering from breast cancer. Both outcomes have a recall score 95% and f1-score of more than 90%. The accuracy of the training model is 95%.

From the confusion matrix generated by the model trained and tested using the PCA generated data, we can see that true negative score is 39 and true positive score is 69. Again, both scores are significantly higher than false positive and false negative scores, indicating that model was almost efficiently generating the expected output.

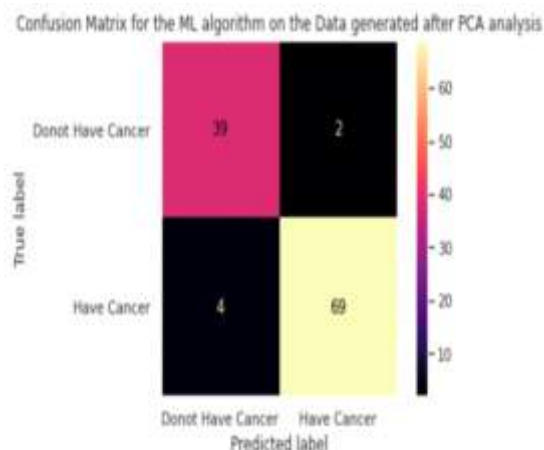


Figure 13: Confusion Matrix for the ML algorithm on the Dataset obtained after PCA.

## VI. FINDINGS

Based on the classification report, accuracy scores and confusion matrix, it can be said for that the classification report generated by the model which was trained by the dataset generated

after PCA analysis shows better accuracy and precision score in comparison to the report generated by the Machine Learning model trained by only the original dataset. We can also look at the confusion matrices and conclude that the total number of false negative and false positive generated by the model trained with PCA generated data is lower in contrast to that of the model trained by the original dataset. Considering all these findings we can say that applying PCA on our breast cancer datasets not only reduces the dimensions of our data but also create a machine learning model from which we are also getting better accuracy score, higher precision, recall and f1-scores as well as getting a better confusion matrix outcome. Although we get similar result from the model trained by the original dataset, yet model trained by PCA datasets showed better output for breast cancer detection data using Logistic Regression Algorithm.

## VII. CONCLUSION

In this paper, we have demonstrated how to perform PCA analysis on a dataset. Initially we performed Principal Component Analysis on the breast Cancer Datasets, and we showed that the attributes are correlated but act differently on the principal components based on biplot. Then we applied the original data and data generated after

PCA analysis to train two separate ML model using Logistic Regression. We compared their classification reports, accuracy scores and confusion matrix. We found that Machine Learning model trained with data generated by PCA had better accuracy, precision, f1 and recall score in comparison to model which was trained using original breast cancer datasets. In the future we hope to repeat these experiments with breast cancer datasets with more attributes so that we can obtain the exact picture about the models trained by PCA generated datasets and original datasets.

## REFERENCES

- [1]. Cancer Treatment Centers for America <https://www.cancercenter.com/>
- [2]. Mayo Clinic <https://www.mayoclinic.org/diseases-conditions/breastcancer/symptoms-causes/syc-20352470>
- [3]. A. Ben Hamza, Chapter 5 : Statistical Process and Quality Control, INSE 6220 course content, Winter 2017..K. Elissa, "Title of paper if known," unpublished.
- [4]. V. Deneshkumar, K. Senthamarai kannan, M. Manikandan, Identification of Outliers in Medical Diagnostic System Using Data Mining Techniques, International

- Journal of Statistics and Applications, Vol.  
4 No. 6, 2014, pp. 241-248.  
doi:10.5923/j.statistics.20140406.01.
- [5]. PythonTutorial  
<https://pythonbasics.org/seabornpairplot>
- [6]. Sklearn: <https://scikit-learn.org/>
- [7]. Sigmoid Function:  
[https://en.wikipedia.org/wiki/Sigmoid\\_function](https://en.wikipedia.org/wiki/Sigmoid_function)