

Principal Component Analysis on the Breast Cancer – Python

Subash Kumar

Bachelor of Computer Science Engineering, Anna University

To Cite this Article

Subash Kumar, "Principal Component Analysis on the Breast Cancer – Python", *International Journal for Modern Trends in Science and Technology*, 6(10): 134-136, 2020.

Article Info

Received on 02-October-2020, Revised on 16-October-2020, Accepted on 18-October-2020, Published on 26-October-2020.

ABSTRACT

Timely diagnosis of any disease is critical in medical field, with increasing population of breast cancer patients, this paper is dedicated to all medical professionals who are trying to save many lives, let us discuss the importance of principal component analysis on the breast cancer in this journal.

KEYWORDS: Machine Learning, Deep Learning, Artificial Intelligence, Medicine, Breast Cancer, Cancer prognosis and prediction, Data Science, Principal Component Analysis, Python

INTRODUCTION

Principal Component Analysis is an unsupervised statistical technique used to examine the interrelations among a set of variables in order to identify the underlying structure of those variables, it is also known as factor analysis, when regression determines a line of best fit to a dataset, a factor analysis determines several orthogonal lines of the best fit to the dataset. The term "Orthogonal" means "at right angles", actually the lines are perpendicular to each other in n-dimensional space, n-dimensional space is the variable sample space, there are many dimensions as there are many variables, so for example in a dataset with 4 variables the sample space is 4 dimensional.

Principal Component Analysis is just a transformation of the data and it attempts to find out what features explain the most variance in the data, for example if we have an original dataset of two components, we try to get rid of the components that don't explain the much of the variance in the data.

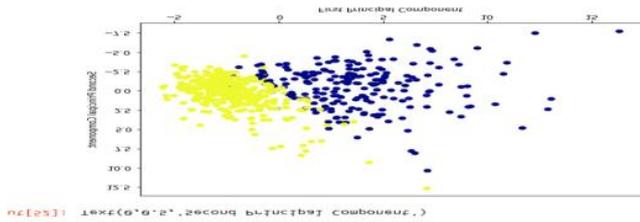
In this paper we will use 12 attributes, out of which 10 real valued features from each cell nucleus

obtained from the Wisconsin diagnostic breast cancer dataset that explains about the stage of breast cancer M (Malignant) and B (Benign).

Features explanations:

1. ID: Patient id
2. Diagnosis (M = Malignant, B = Benign)
3. Radius(mean of distances from center to points on the perimeter) (worst). Worst texture. Texture (standard deviation of gray-scale values) (worst). Worst perimeter. perimeter (worst)
4. Texture (Breast cancer can cause changes and inflammation in skin cells that can lead to texture changes), here we can take the standard deviation of grey scaled values
5. Perimeter: Size of the core tumor
6. Area: Area of the core tumor
7. Smoothness: Local variation in radius length
8. Compactness: $(\text{perimeter}^2 / \text{area} - 1.0)$
9. Concavity (severity of concave portions of the contour)
10. Concave points (Number of concave portions of the contour)
11. Symmetry

Output:



11. We can see visually what these principal components represents through heat map, the below heat map shows the relationship between the correlation of the various features and the principal components

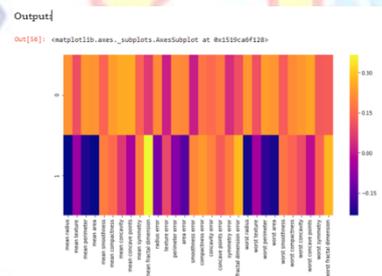
Input:

```
df_comp = pd.DataFrame(pca.components_, columns=cancer['feature_names'])
```

Output:

	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
0	0.218902	0.103725	0.227537	0.220995	0.142590	0.239285	0.258400	0.260854	0.138167	0.064363	0.227997	0.104469	0.236640
1	-0.233857	-0.059708	-0.215101	-0.231077	0.186113	0.151892	0.090165	-0.034768	0.190349	0.386575	-0.219866	-0.045407	-0.199878

```
plt.figure(figsize=(12,6))
sns.heatmap(df_comp, cmap='plasma')
```



RESULTS

The Wisconsin Breast cancer data with 30 features was analyzed using PCA component analysis basically each principal component shows here as a row and the higher the number are the dark looking colors or it is yellow which is more correlated to a specific feature in the columns and that is really the best explanation for what these principal components represents their combinations of all these features and we can see what features were specifically important to one of the principal components versus the other

REFERENCES

- [1] Author: Ian Joliffe Reference work entry First Online: 02 December 2014 DOI: https://doi.org/10.1007/978-3-642-04898-2_455
- [2] Karl Pearson F.R.S. , 1901. LIII. On lines and planes of closest fit to systems of points in space. The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science, 2(11), pp.559–572.
- [3] DOI:10.1016/0169-7439(87)80084-9 Corpus ID: 54979942 Principal component analysis S. Wold, K. Esbensen, P. Geladi Published 1987 Computer Science Chemometrics and Intelligent Laboratory Systems
- [4] Principal components analysis. Groth D, Hartmann S, Klie S, Selbig J Author information Methods in Molecular Biology (Clifton, N.J.), 31 Dec 2012, 930:527-547 DOI: 10.1007/978-1-62703-059-5_22 PMID: 23086856