

Principal Component Analysis

Nisha Akole, G V V Sharma*

CONTENTS

| | | |
|---|------------------------------|---|
| 1 | Objective | 1 |
| 2 | Load Dataset | 1 |
| 3 | About PCA | 1 |
| 4 | Pre-processing | 1 |
| 5 | Compute Covariance Matrix | 1 |
| 6 | Compute Eigenvectors | 2 |
| 7 | Principal Component and plot | 2 |
| 8 | Figures | 2 |

Abstract—This manual provides a brief description on how to implement Principal Component Analysis from scratch and use it for dimensionality reduction of data.

1. OBJECTIVE

Our objective is to implement PCA and reduce dimensionality of data which gives better understanding of data visually.

2. LOAD DATASET

The dataset used for PCA is available at the following link. Download all the data file in the folder where you want to write code for PCA.

<https://github.com/prabhatrai111/Commensal-Radar>

```
import numpy as np
import scipy.io as sio
from sklearn.utils.extmath import randomized_svd

mat_contents = sio.loadmat('data_all.mat')
X_data = mat_contents['data_all']
```

3. ABOUT PCA

Principal component analysis (PCA) is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components.

$$\text{cov}[PC1, PC2] = E[PC1PC2] - E[PC1]E[PC2]=0$$

Where, PC - principal components There are four main operations in the PCA:

- 1) Pre-processing
- 2) Co-variance Matrix
- 3) Eigen Vectors
- 4) Feature Vector and Plot

4. PRE-PROCESSING

Feature scaling is an important task in PCA to get an optimized output. The functions StandardScaler or MinMaxScaler will perform mean normalization by making mean = 0 and variance = 1. Standard-Scaler Function :

$$z = (X - \mu)/\sigma$$

where, μ =Mean of data X
 σ =Standard Deviation of data X

```
X = np.matrix(X_data)
μ = X.mean(0)
σ = X.std(0)
X_std = (X - μ)/σ
```

5. COMPUTE COVARIANCE MATRIX

$$X_{cov} = \begin{bmatrix} var_{PC1} & cov_{PC1,PC2} \\ cov_{PC1,PC2} & var_{PC2} \end{bmatrix} = \begin{bmatrix} \sigma_{PC1}^2 & \sigma_{PC1,PC2} \\ \sigma_{PC1,PC2} & \sigma_{PC2}^2 \end{bmatrix}$$

If variables are highly correlated, they contain redundant information.

$$X_{cov}^T = X_{cov}$$

```
X_cov = np.cov(X_std)
```

*The authors are with the Department of Electrical Engineering, Indian Institute of Technology, Hyderabad 502285 India e-mail: gadepall@iith.ac.in.

6. COMPUTE EIGENVECTORS

Eigenvalues and eigenvectors can be calculated by Singular Value Decomposition(SVD) of a covariance matrix.

$$X_{cov} = U \Sigma V^T$$

Eigenvalues are arranged in decreasing order in sigma matrix and its corresponding eigenvector is arranged in U.

$$\Sigma^T = \Sigma$$

```
U, Σ, VT = randomized_svd(X_cov,
    n_components=2, n_iter= 5, random_state=
    none)
```

n_components will consider first 2 largest eigenvalues. Hence, first 2 eigenvectors will be used for plotting the data.

7. PRINCIPAL COMPONENT AND PLOT

To visualize the data, only 2 or 3 principal components are needed. To decide K principal components out of n to retain maximum information:

$$\frac{\sum_{i=1}^K \text{FirstKDiagonalElements}}{\sum \text{AllDiagonalElements}} \geq 0.99$$

Link for PCA without function is given below:

<https://github.com/NishaAkole/AI-and-ML/blob/master/codes/pca/PCAwwithoutFun.py>

The above whole program can be combined except preprocessing and a single command(i.e using inbuilt function) can do all the

<https://github.com/NishaAkole/AI-and-ML/blob/master/codes/pca/PCAFun.py>

8. FIGURES

From fig. 5, we are not able to reduce the dimensionality below 8000. Hence, PCA is not a good approach for dimensionality reduction for this data.

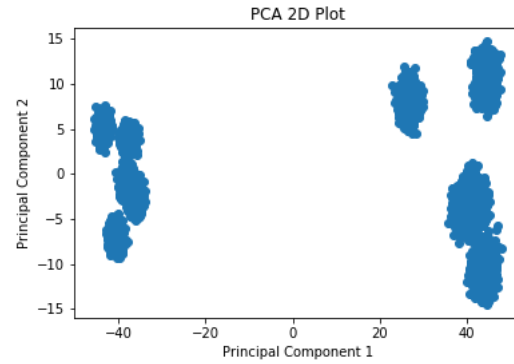


Fig. 1: PCA with Function 2D

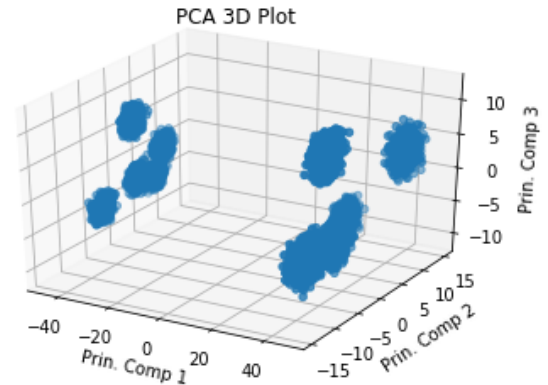


Fig. 2: PCA with Function 3D

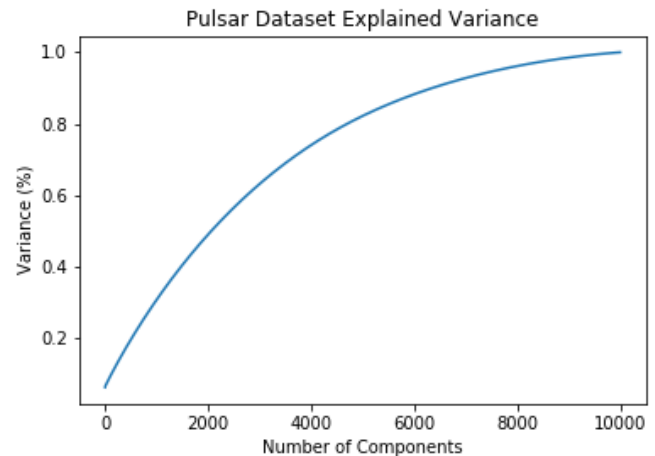


Fig. 3: Information vs No of Dimensions