# Assignment - 2

Name: Nusrat Jannat
ID: 20121 76145
Session: 2019-20
Year: 4th year even semester
Topic: Principal Component Analysis
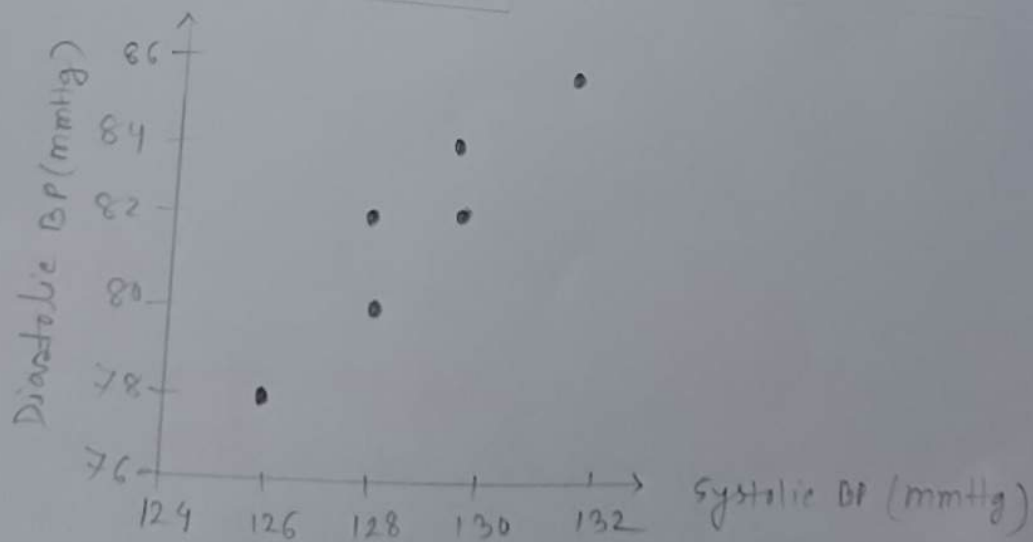Course: Machine Learning

Principal Component Analysis (PCA): It is a popular unsupervised dimensionality reduction technique in machine learning used to transform high dimensional data into a lower dimensional representation.

To explain how the PCA works, we will use the example data to combine the two blood pressure variables into just one variable based on data from six individuals.

The datasets are as following:

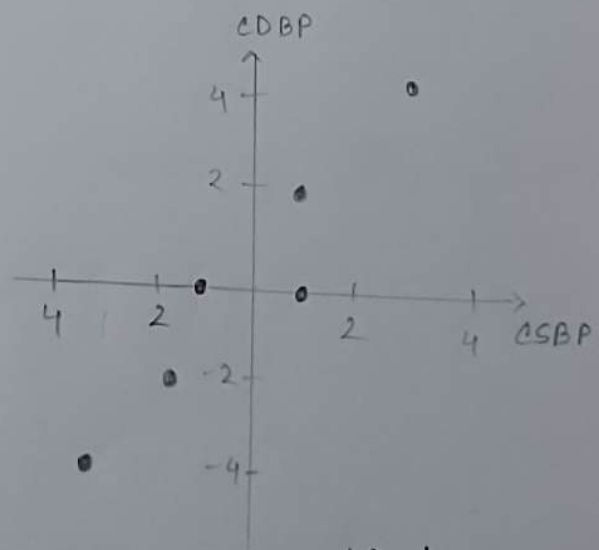| Systolic BP | Diastolic BP |
|-------------|--------------|
| 126 | 78 |
| 128 | 80 |
| 128 | 82 |
| 130 | 82 |
| 130 | 84 |
| 132 | 86 |

To compute PCA, we need to perform following steps:

1. Center the data.
2. Calculate the covariance matrix.
3. Calculate eigenvalues of the covariance matrix.
4. Calculate eigenvectors of the covariance matrix.
5. Order the eigenvectors.
6. Calculate the principle components.

Step 1: Center the data: In this case, we will only center the data which means that we subtract all the values for each variable by its corresponding mean. The means by which we subtract systolic blood pressure from each observation is 129 and diastolic BP is 82.

| Centered SBP | Centered DBP |
|---|---|
| -3 | -4 |
| -1 | -2 |
| -1 | 0 |
| 1 | 0 |
| 1 | 2 |
| 3 | 4 |



When we center the data, it means that we center the data points around the origin. It will help us when we will rotate the data.

Step2: Calculate the covariance matrix: We calculate the covariance matrix based on the centered data.

$$\text{var}\,(\text{CSBP}) = \frac{1}{n-1} \sum_{j=1}^{n} \left(\text{CSBP}_i - \overline{\text{CSBP}}\right)^2$$

$$= \left((-3)^2 + (-1)^2 + (-1)^2 + 1^2 + 1^2 + 3^2\right)/(6-1) = 22/5 = 4.4$$

$$\text{var}\,(\text{CDBP}) = \frac{1}{n-1} \sum_{j=1}^{n} \left(\text{CDBP}_i - \overline{\text{CDBP}}\right)^2$$

$$= \left((-4)^2 + (-2)^2 + 0^2 + 6^2 + 2^2 + 4^2\right)/(6-1) = 40/5 = 8$$

$$\text{cov}\,(\text{CSBP}, \text{CDBP}) = \frac{1}{n-1} \sum \left(\text{CSBP}_i - \overline{\text{CSBP}}\right)\cdot\left(\text{CDBP}_i - \overline{\text{CDBP}}\right)$$

$$= \left((-3)\cdot(-4) + (-1)\cdot(-2) + (-1)\cdot 0 + 1\cdot 0 + 1\cdot 2 + 3\cdot 4\right)/(6-1)$$

$$= 28/5 = 5.6$$

|      | SBP | DBP |
|------|-----|-----|
| SBP  | 4.4 | 5.6 |
| DBP  | 5.6 | 8.0 |

covariance matrix

$$A = \begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{bmatrix}$$

Step3: Calculate the eigenvalues of the covariance matrix

$$\det |A - \lambda I| = 0$$

$$\det \left| \begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \right| = 0$$

$$\det \left| \begin{bmatrix} (4.4 - \lambda) & 5.6 \\ 5.6 & (8.0 - \lambda) \end{bmatrix} \right| = 0$$

Calculating the determinstic of the matrix

$(4.4 - \lambda)(8.0 - \lambda) - 5.6 \cdot 5.6 = 0$

$\Rightarrow 3.84 - 12.4\lambda + \lambda^2 = 0$

So $\lambda_1 = 0.32$ , $\lambda_2 = 12.08$

Step 4: Calculate eigenvectors of the covariance matrix.

while $\lambda_2 = 12.08$

$A \cdot V = \lambda \cdot v$

$\Rightarrow \begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = 12.08 \begin{bmatrix} x \\ y \end{bmatrix}$

So we get

$4.4x + 5.6y = 12.08x$ 　　　　$5.6x + 8.0y = 12.08y$

$\Rightarrow 5.6y = 7.68x$ 　　　　　$\Rightarrow 5.6x = 4.08y$

$\therefore y = 1.37x$ 　　　　　　　$\therefore y = 1.37x$

So eigenvector, $v_2 = \begin{bmatrix} 1 \\ 1.37x \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1.37 \end{bmatrix}$

normalisation $v_2 = \begin{bmatrix} \frac{1}{\sqrt{1^2 + 1.37^2}} \\ \frac{1.37}{\sqrt{1^2 + 1.37^2}} \end{bmatrix} = \begin{bmatrix} 0.59 \\ 0.81 \end{bmatrix}$

Now, while $\lambda_1 = 0.32$

$A \cdot v = \lambda \cdot v$

$\Rightarrow \begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = 0.32 \begin{bmatrix} x \\ y \end{bmatrix}$

Here eigen vectors, $v_1 = \begin{bmatrix} -1 \\ 0.72 \end{bmatrix}$

we get after normalization,

$$V_1 = \begin{bmatrix} \dfrac{-1}{\sqrt{(1)^2 + (0.72)^2}} \\ \dfrac{0.72}{\sqrt{(1)^2 + (0.72)^2}} \end{bmatrix} = \begin{bmatrix} -0.81 \\ 0.59 \end{bmatrix}$$

Step 5: Order the eigenvectors: We order the eigenvectors based on their corresponding eigenvalues, where the eigenvector with the largest eigenvalue becomes our first eigenvector.

So, $V_1 = \begin{bmatrix} 0.59 \\ 0.81 \end{bmatrix}$ ; $h_2 = 12.08$

$V_2 = \begin{bmatrix} -0.81 \\ 0.59 \end{bmatrix}$ ; $h_1 = 0.32$

So the matrix $V = \begin{bmatrix} 0.59 & -0.81 \\ 0.81 & 0.59 \end{bmatrix}$

Step 6: Calculate the principal components: We will now use this matrix to transform our original centered data so that the two variable are completely uncorrelated.
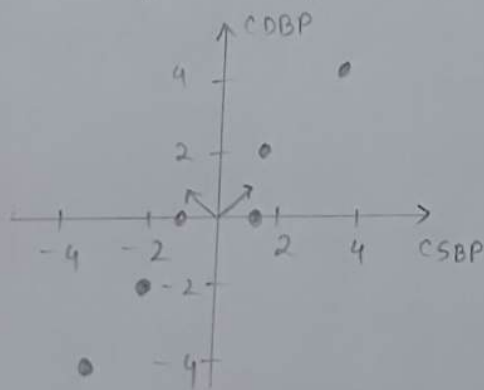
write the centered data as matrix D.

$$D = \begin{bmatrix} -3 & -4 \\ -1 & -2 \\ -1 & 0 \\ 1 & 0 \\ 1 & 2 \\ 3 & 4 \end{bmatrix}$$
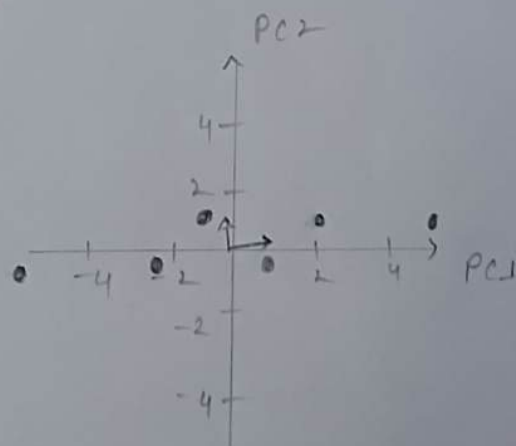
Now we will multiply our data matrix D by matrix $v$, which includes our eigenvectors as columns.to get transformed data.

$$D v = \begin{bmatrix} -3 & -4 \\ -1 & -2 \\ -1 & 0 \\ 1 & 0 \\ 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 0.59 & -0.81 \\ 0.81 & 0.59 \end{bmatrix} = \begin{matrix} PC1 & PC2 \\ \begin{bmatrix} -5.0 & 0.1 \\ -2.2 & -0.4 \\ -0.6 & 0.8 \\ 0.6 & -0.8 \\ 2.2 & 0.4 \\ 5.0 & -0.1 \end{bmatrix} \end{matrix}$$

When we go from our original data matrix to the transformed data, this can be seen like we rotate the data clockwise until the two eigenvectors point in the same direction as the x and y-axes of the plot.



Before transform

After transform

Now let's compare the centered data with the transformed data.

| Centered SBP | Centered DBP |
|---|---|
| −3 | −4 |
| −1 | −2 |
| −1 | 0 |
| 1 | 0 |
| 1 | 2 |
| 3 | 4 |
| van = 4.4 | varr = 8.0 |

Covariance matrix $= \begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{bmatrix}$

| PC1 | PC2 |
|---|---|
| −5.0 | 0.1 |
| −2.2 | −0.4 |
| −0.6 | 0.8 |
| 0.6 | −0.8 |
| 2.2 | 0.4 |
| 5.0 | −0.1 |
| Varr = 12.08 | van = 0.32 |

Covariance matrix $= \begin{bmatrix} 12.08 & 0 \\ 0 & 0.32 \end{bmatrix}$

Now we can see for

% var of PC1 $= \dfrac{12.08}{12.08 + 0.32} = 97.4\%$.

% var of PC2 $= \dfrac{0.32}{12.08 + 0.32} = 2.6\%$.

As we can see that the first principal component captures 97.4% of the total variance that means that almost all variance is kept in the first principle component. and PC1 is the 1D representation of this 2D dataset.