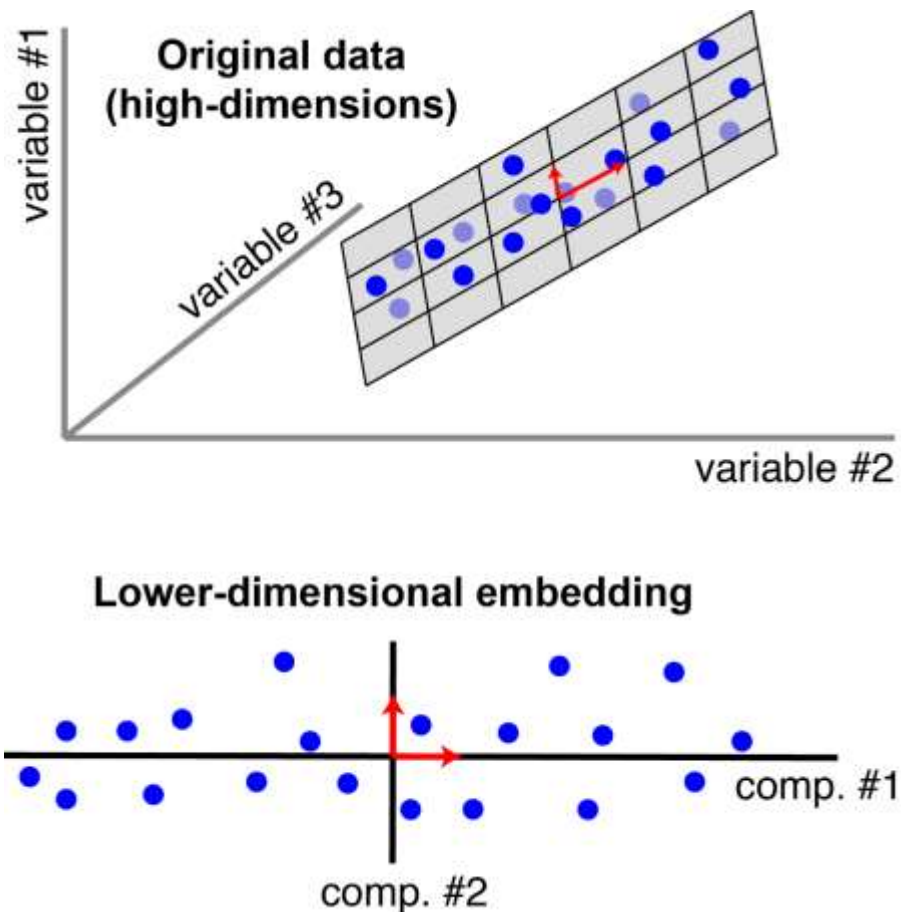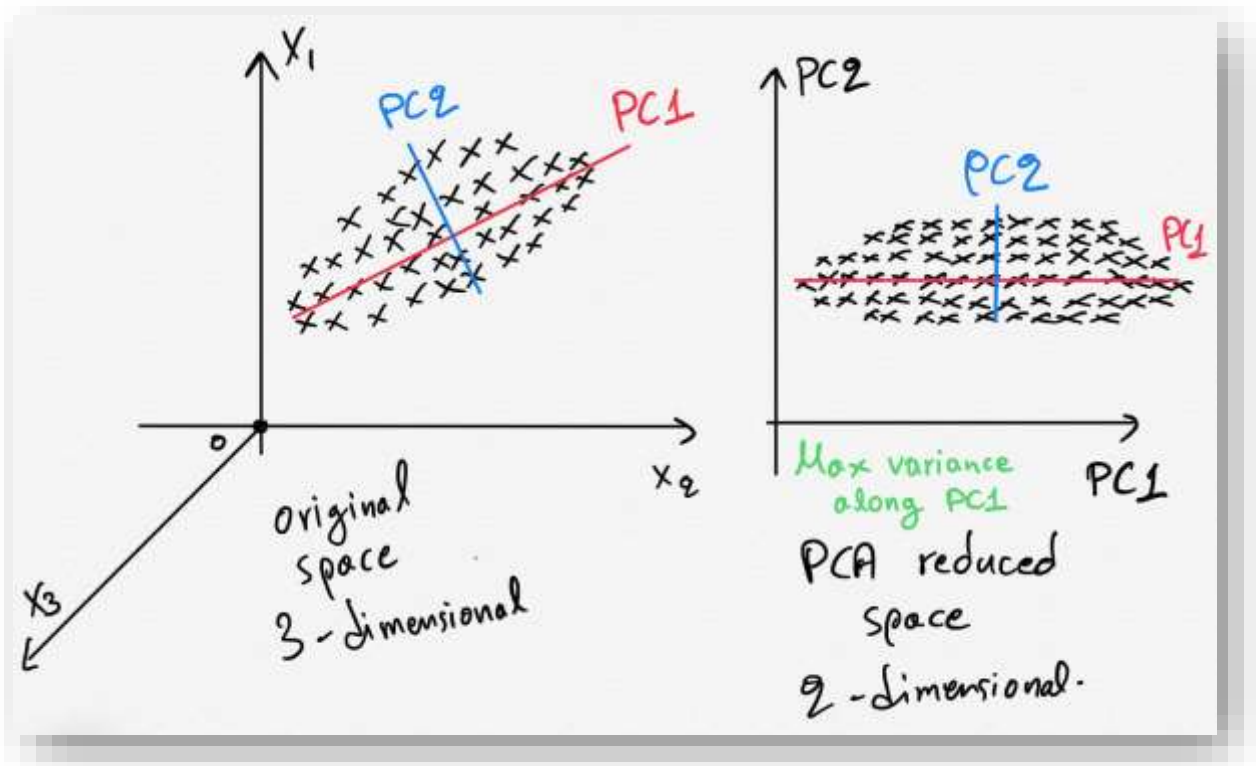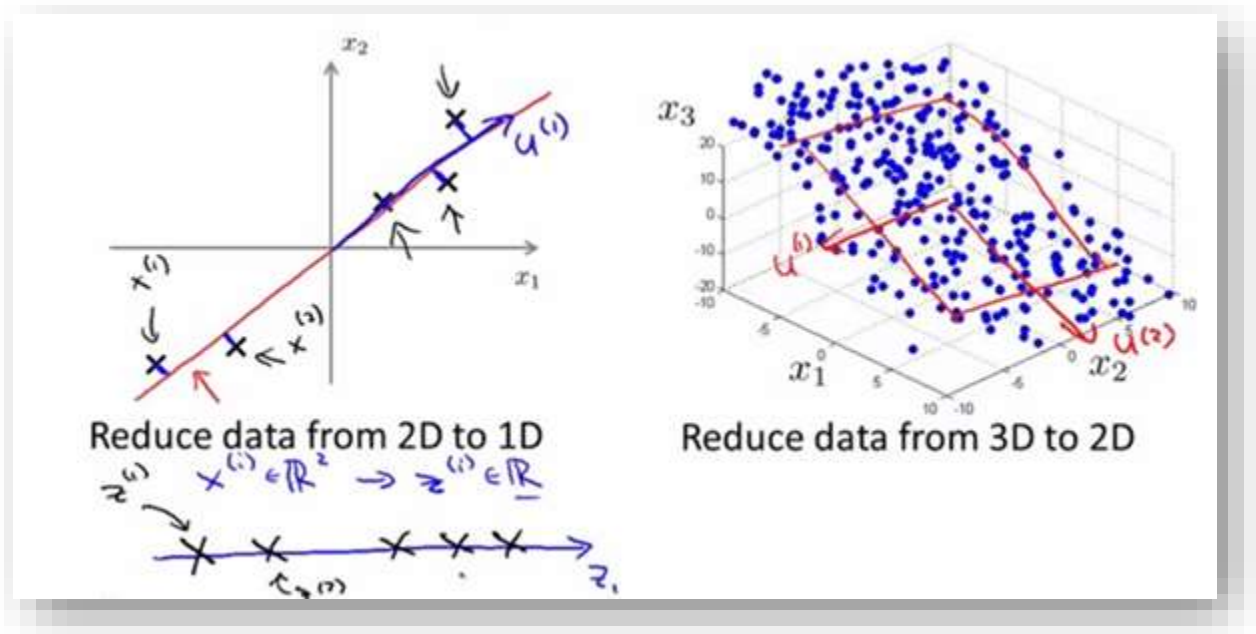# Principal Component Analysis (PCA)

Principal Component Analysis (PCA) is a dimensionality reduction technique commonly used in data analysis, machine learning, and statistics. It is used to simplify complex datasets by reducing the number of variables while preserving the essential information. PCA achieves this by transforming the original variables into a new set of variables called principal components, which are linear combinations of the original variables.

Principal Component Analysis (PCA) is a dimensionality reduction technique invented by Karl Pearson in 1901, which is used for identification of a smaller number of uncorrelated variables known as Principal Components from a larger set of data. The following figures illustrates the same:

Reduce data from 2D to 1D

$x^{(i)} \in \mathbb{R}^2 \rightarrow z^{(i)} \in \mathbb{R}$

Reduce data from 3D to 2D



original space
3-dimensional

Max variance along PC1

PCA reduced space
2-dimensional.

**Dimension reduction offers several benefits such as:**

- Compresses the data and thus reduces the storage space requirements
- Reduces the time required for computation since less dimensions require less computation
- Eliminates the redundant features
- Improves the model performance

**Properties of Principal Component Analysis:**

- Principal Component Analysis is a well-known dimension reduction technique
- It transforms the variables into a new set of variables called as principal components
- These principal components are linear combination of original variables and are orthogonal
- The first principal component accounts for most of the possible variation of original data
- The second principal component does its best to capture the variance in the data
- There can be only two principal components for a two-dimensional data set

## Step-By-Step Explanation of PCA (Principal Component Analysis)

Principal Component Analysis (PCA) is a mathematical technique used for dimensionality reduction and data compression while preserving the most critical information in the data. It helps to reduce the complexity of high-dimensional data by transforming it into a lower-dimensional representation. Here's a step-by-step explanation of PCA:

# Step 1: Data Preprocessing

**Data Collection**: Gather the dataset you want to analyse. Each row typically represents an observation, while each column corresponds to a variable (feature).

**Standardization**: If the variables are measured on different scales, standardize the data by subtracting the mean and dividing by the standard deviation for each variable.

Mathematically, this can be done by subtracting the mean and dividing by the standard deviation for each value of each variable. This is often called as Z-score.

$$z = \frac{value - mean}{standard\ deviation}$$

Once the standardization is done, all the variables will be transformed to the same scale.

## Step 2: Compute the Covariance Matrix

**Covariance Matrix**: Calculate the covariance matrix of the standardized data. The covariance matrix, denoted as Σ (Sigma), describes the relationships and variances between variables. The diagonal elements of Σ represent the variances of individual variables, while off-diagonal elements represent covariances between pairs of variables.

$$\Sigma = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(x_i - \bar{x})^T$$

where n is the number of observations, $x_i$ is a data point and $\bar{x}$ is the mean vector of the data.

## Step 3: Compute Eigenvalues and Eigenvectors

**Eigenvalue Decomposition**: Perform eigenvalue decomposition on the covariance matrix Σ to find its eigenvalues λ and corresponding eigenvectors v

$$\Sigma v = \lambda v$$

Where λ represents the eigenvalues and $v$ represents the eigenvectors.

## Step 4: Sort Eigenvalues

**Sort Eigenvalues**: Sort the eigenvalues in descending order. The eigenvalues indicate the amount of variance captured by each principal component (PC). By sorting them, you identify the most significant PCs.

## Step 5: Select Principal Components

**Select PCs**: Decide how many principal components (PCs) to retain. You can choose a specific number of PCs or determine it based on the cumulative explained variance. A common approach is to retain enough PCs to explain a significant portion (e.g., 95%) of the total variance.

## Step 6: Projection

**Projection**: Project the original data onto the selected principal components to obtain a lower-dimensional representation. This is done by multiplying the original data matrix ($X$) by the matrix of selected eigenvectors.

$$Y = X \cdot V_k$$

where $Y$ is the reduced-dimensional data matrix, $X$ is the original data matrix and $V_k$ is a matrix containing the top $k$ eigenvectors.

## Step 7: Interpretation

**Interpretation**: The reduced-dimensional data in matrix $Y$ now contains the information from the most important principal components. You can interpret these components to understand which variables or patterns contribute most to the variance in the data.

## Step 8: Reconstruct Data (Optional)

**Reconstruction (Optional)**: If needed, you can reverse the dimensionality reduction by multiplying $Y$ by the transpose of $Vk$ and adding back the mean (if standardization was performed). This gives you an approximation of the original data.

$$Xapprox = Y \cdot Vk^T + \bar{x}$$

PCA is a valuable tool for reducing dimensionality, visualizing data and simplifying complex datasets. It is widely used in various fields, including data science, machine learning and signal processing.

### Python Code to Find PCA

To perform Principal Component Analysis (PCA) for the provided dataset in Python, you can follow these steps:

1. Create a NumPy array to represent the data.

2. Standardize the data.

3. Compute the covariance matrix.

4. Find the eigenvalues and eigenvectors.

5. Sort the eigenvalues and eigenvectors.

6. Choose the number of principal components to retain.

7. Project the data onto the selected principal components.

8. Visualize the results (optional).

# Example

**1. Find the Principal Component Analysis (PCA) for the following dataset:**

| $x_1$ | 4 | 8 | 13 | 7 |
|---|---|---|---|---|
| $x_2$ | 11 | 4 | 5 | 14 |

## Solution:

First, we need to standardize the data by subtracting the mean and dividing by the standard deviation for each variable.

Mean for $x_1$ and $x_2$ = [8, 8.5]

Standard deviation for $x_1$ and $x_2$ = [3.24037035, 4.15331193]

Standardized data is

| $x_1$_std | -1.2344268 | 0.0 | 1.5430335 | -0.3086067 |
|---|---|---|---|---|
| $x_2$_std | 0.60192927 | -1.08347268 | -0.84270097 | 1.32424438 |

Next, we calculate the covariance matrix ($\Sigma$) for the standardized data:

Covariance Matrix is

$$\Sigma = \begin{bmatrix} 1.33333333 & -0.81734138 \\ -0.81734138 & 1.33333333 \end{bmatrix}$$

Now, we find the eigenvalues ($\lambda$) and corresponding eigenvectors ($v$) of the covariance matrix ($\Sigma$).

The eigenvalues are $\lambda_1 = 2.15067471$ and $\lambda_2 = 0.51599195$

The corresponding eigenvectors are

$$\begin{bmatrix} 0.70710678 & 0.70710678 \\ -0.70710678 & 0.70710678 \end{bmatrix}$$

Sort the eigenvalues in descending order and arrange the corresponding eigenvectors accordingly. In this case, $\lambda_1 > \lambda_2$, so:

Principal Component 1 (PC$_1$) corresponds to $\lambda_1$ and $v_1$.

Principal Component 2 (PC$_2$) corresponds to $\lambda_2$ and $v_2$.

Now, we need to choose the number of principal components to retain.

Let the number of principal components = 1

Now, Project the data onto the selected principal components.

This can be done by multiplying the original data matrix by the matrix of selected eigenvectors.

$$Y = X \cdot V_1 = \begin{bmatrix} -1.2344268 & 0.60192927 \\ 0. & -1.08347268 \\ 1.5430335 & -0.84270097 \\ -0.3086067 & 1.32424438 \end{bmatrix} \begin{bmatrix} 0.70710678 \\ -0.70710678 \end{bmatrix}$$

Projected data is $\begin{bmatrix} -1.29849983 \\ 0.76613088 \\ 1.68696902 \\ -1.15460007 \end{bmatrix}$

Explained variance ratio for the chosen component is

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} \times 100 = \frac{2.15067471}{2.15067471 + 0.51599195} \times 100 = 80.65\%$$

## Above example using Python Code

```python
1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  # Define the dataset
5  data = np.array([[4,11],[8,4],[13,5],[7,14]])
6
7  # Step 2: Standardize the data
8  mean = np.mean(data, axis=0)
9  std_dev = np.std(data, axis=0)
10 data_std = (data - mean) / std_dev
11 print('Mean = ', mean)
12 print('Standard deviation =', std_dev)
13 print('Standardized data is\n', data_std)
14
15 # Step 3: Compute the covariance matrix
16 cov_matrix = np.cov(data_std, rowvar=False)
17 print('\nCovariance Matrix is \n',cov_matrix)
18
19 # Step 4: Find the eigenvalues and eigenvectors
20 eigenvalues, eigenvectors = np.linalg.eig(cov_matrix)
21 print('\nEigenvalues =',eigenvalues)
22 print('\nCorresponding Eigenvectors are\n ',eigenvectors)
23
24 # Step 5: Sort the eigenvalues and eigenvectors
25 sorted_indices = np.argsort(eigenvalues)[::-1]
26 eigenvalues = eigenvalues[sorted_indices]
27 eigenvectors = eigenvectors[:, sorted_indices]
28 print('\nSorted Eigenvalues are',eigenvalues)
29 print('Sorted Eigenvectors are\n',eigenvectors)
30
31 # Step 6: Choose the number of principal components to retain
32 num_components = 1  # Choose the number of components (e.g., 1 for visualization)
33
34 # Step 7: Project the data onto the selected principal components
35 projected_data = np.dot(data_std, eigenvectors[:, :num_components])
36 print('\nProjected data is\n',projected_data)
37
38 # Print explained variance ratio for the chosen components
39 explained_variance_ratio = eigenvalues[:num_components] / np.sum(eigenvalues)
40 print(f'\nExplained Variance Ratio: {explained_variance_ratio[0]*100:.4f}%')
41
42 # Step 8: Visualize the results (optional)
43 plt.scatter(projected_data, np.zeros_like(projected_data))
44 plt.xlabel('Principal Component 1')
45 plt.title('PCA Projection')
46 plt.show()
```

```
Mean = [8.  8.5]
Standard deviation = [3.24037035 4.15331193]
Standardized data is
 [[-1.2344268   0.60192927]
 [ 0.         -1.08347268]
 [ 1.5430335  -0.84270097]
 [-0.3086067   1.32424438]]

Covariance Matrix is
 [[ 1.33333333 -0.81734138]
 [-0.81734138  1.33333333]]

Eigenvalues = [2.15067471 0.51599195]

Corresponding Eigenvectors are
  [[ 0.70710678  0.70710678]
 [-0.70710678  0.70710678]]

Sorted Eigenvalues are [2.15067471 0.51599195]
Sorted Eigenvectors are
 [[ 0.70710678  0.70710678]
 [-0.70710678  0.70710678]]

Projected data is
 [[-1.29849983]
 [ 0.76613088]
 [ 1.68696902]
 [-1.15460007]]

Explained Variance Ratio: 80.6503%
```
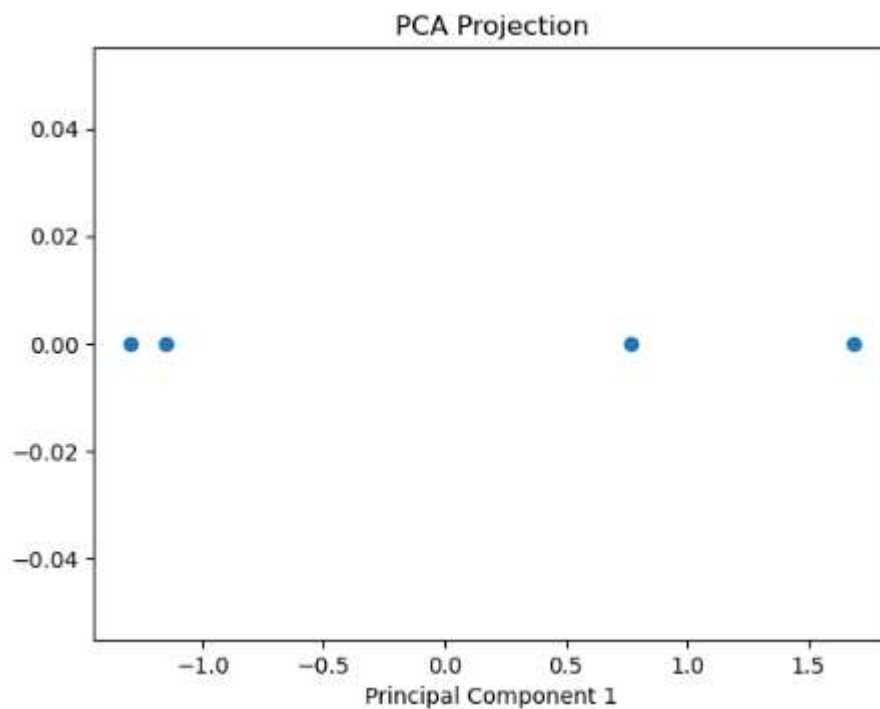


PCA Projection

## 2. Find the Principal Component Analysis (PCA) for the following dataset:

| X1 | 2.5 | 0.5 | 2.2 | 1.9 | 3.1 | 2.3 | 2.0 | 1.0 | 1.5 | 1.1 |
|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| X2 | 2.4 | 0.7 | 2.9 | 2.2 | 3.0 | 2.7 | 1.6 | 1.1 | 1.6 | 0.9 |

## Using Python Code

```python
1   import numpy as np
2   import matplotlib.pyplot as plt
3   # Define the dataset
4   data = np.array([[2.5, 2.4],[0.5, 0.7],[2.2, 2.9],[1.9, 2.2],[3.1, 3.0],[2.3, 2.7],
5                    [2.0, 1.6],[1.0, 1.1],[1.5, 1.6],[1.1, 0.9]])
6   # Step 2: Standardize the data
7   mean = np.mean(data, axis=0)
8   std_dev = np.std(data, axis=0)
9   data_std = (data - mean) / std_dev
10  print('Mean = ', mean)
11  print('Standard deviation =', std_dev)
12  print('Standardized data is\n', data_std)
13  # Step 3: Compute the covariance matrix
14  cov_matrix = np.cov(data_std, rowvar=False)
15  print('\nCovariance Matrix is \n',cov_matrix)
16  # Step 4: Find the eigenvalues and eigenvectors
17  eigenvalues, eigenvectors = np.linalg.eig(cov_matrix)
18  print('\nEigenvalues =',eigenvalues)
19  print('\nCorresponding Eigenvectors are\n ',eigenvectors)
20  # Step 5: Sort the eigenvalues and eigenvectors
21  sorted_indices = np.argsort(eigenvalues)[::-1]
22  eigenvalues = eigenvalues[sorted_indices]
23  eigenvectors = eigenvectors[:, sorted_indices]
24  print('\nSorted Eigenvalues are',eigenvalues)
25  print('Sorted Eigenvectors are\n',eigenvectors)
26  # Step 6: Choose the number of principal components to retain
27  num_components = 1  # Choose the number of components (e.g., 1 for visualization)
28  # Step 7: Project the data onto the selected principal components
29  projected_data = np.dot(data_std, eigenvectors[:, :num_components])
30  print('\nProjected data is\n',projected_data)
31  # Print explained variance ratio for the chosen components
32  explained_variance_ratio = eigenvalues[:num_components] / np.sum(eigenvalues)
33  print(f'\nExplained Variance Ratio: {explained_variance_ratio[0]*100:.4f}%')
34  # Step 8: Visualize the results (optional)
35  plt.scatter(projected_data, np.zeros_like(projected_data))
36  plt.xlabel('Principal Component 1')
37  plt.title('PCA Projection')
38  plt.show()
```

```
Mean =  [1.81 1.91]
Standard deviation = [0.7449161  0.80305666]
Standardized data is
 [[ 0.92627881  0.61016865]
 [-1.7585873  -1.506743  ]
 [ 0.52354889  1.23278973]
 [ 0.12081898  0.36112022]
 [ 1.73173864  1.35731394]
 [ 0.6577922   0.9837413 ]
 [ 0.25506228 -0.38602507]
 [-1.08737078 -1.00864614]
 [-0.41615425 -0.38602507]
 [-0.95312747 -1.25769457]]

Covariance Matrix is
 [[1.11111111 1.0288103 ]
 [1.0288103  1.11111111]]

Eigenvalues = [2.13992141 0.08230081]

Corresponding Eigenvectors are
 [[ 0.70710678 -0.70710678]
 [ 0.70710678  0.70710678]]

Sorted Eigenvalues are [2.13992141 0.08230081]
Sorted Eigenvectors are
 [[ 0.70710678 -0.70710678]
 [ 0.70710678  0.70710678]]

Projected data is
 [[ 1.08643242]
 [-2.3089372 ]
 [ 1.24191895]
 [ 0.34078247]
 [ 2.18429003]
 [ 1.16073946]
 [-0.09260467]
 [-1.48210777]
 [-0.56722643]
 [-1.56328726]]

Explained Variance Ratio: 96.2965%
```
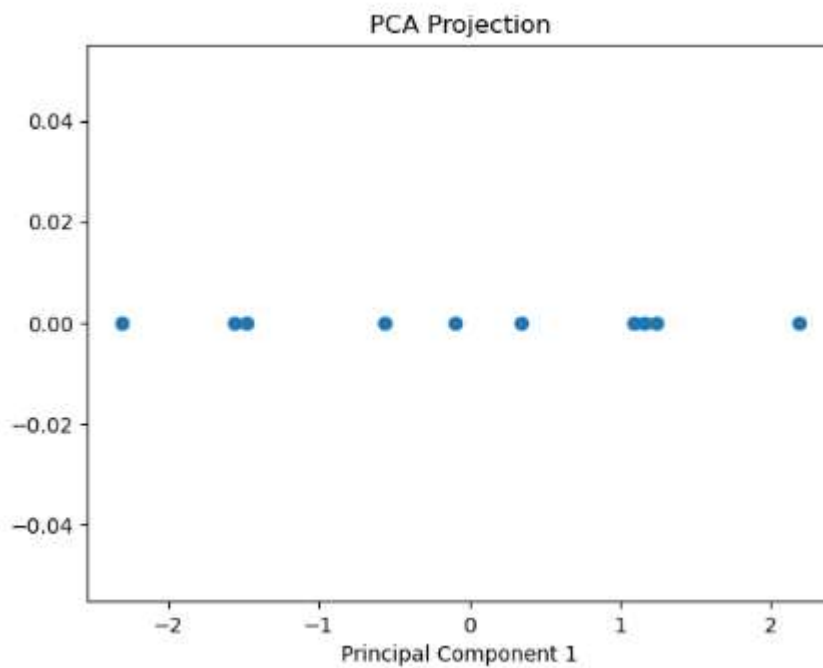
PCA Projection

**To perform Principal Component Analysis (PCA) in Python using the sklearn library:**

In this code:

- We first import the necessary libraries.
- We assume `X` is your data.
- We standardize the data using `StandardScaler` from `sklearn.preprocessing`. This subtracts the mean and divides by the standard deviation for each variable.
- We perform PCA using `PCA` from `sklearn.decomposition`. The `n_components` parameter specifies the number of components to keep. In this case, we are keeping all components.
- The transformed data is stored in `principalComponents`.

# Examples

```python
1  from sklearn.decomposition import PCA
2  from sklearn.preprocessing import StandardScaler
3  import numpy as np
4
5  # Assume X is your data
6  X = np.array([[1, 2], [3, 4], [5, 6]])
7
8  # Standardize the data
9  scaler = StandardScaler()
10 X_std = scaler.fit_transform(X)
11
12 # Perform PCA
13 pca = PCA(n_components=1)
14 principalComponents = pca.fit_transform(X_std)
15
16 print(principalComponents)
17
```

```
[[ 1.73205081]
 [-0.        ]
 [-1.73205081]]
```

```python
1  import numpy as np
2  import matplotlib.pyplot as plt
3  from sklearn.decomposition import PCA
4  from sklearn.preprocessing import StandardScaler
5  # Create a sample dataset (you can replace this with your dataset)
6  data = np.array([[2.5, 2.4],[0.5, 0.7],[2.2, 2.9],[1.9, 2.2],[3.1, 3.0],
7           [2.3, 2.7],[2.0, 1.6],[1.0, 1.1],[1.5, 1.6],[1.1, 0.9]])
8  # Step 3: Standardize the data
9  scaler = StandardScaler()
10 data_std = scaler.fit_transform(data)
11 print('Standardized data is:\n',data_std)
12 # Step 4: Initialize PCA
13 pca = PCA(n_components=2)   # Choose the number of components you want
14 print(pca)
15 # Step 5: Fit the PCA model to the standardized data
16 pca.fit(data_std)
17 # Step 6: Transform the data into the PCA space
18 data_pca = pca.transform(data_std)
19 print("Principal component Analysis is \n",data_pca)
20 # Step 7: Visualize the results (optional)
21 plt.figure(figsize=(8, 4))
22 # Original data
23 plt.subplot(1, 2, 1)
24 plt.scatter(data_std[:, 0], data_std[:, 1], c='blue')
25 plt.title('Original Data')
26 # PCA-transformed data
27 plt.subplot(1, 2, 2)
28 plt.scatter(data_pca[:, 0], data_pca[:, 1], c='red')
29 plt.title('PCA-Transformed Data')
30 plt.tight_layout()
31 plt.show()
32 # Print explained variance ratios
33 explained_variance_ratio = pca.explained_variance_ratio_
34 print(f'\nExplained Variance Ratios: {explained_variance_ratio[0]*100:.3f}% and {explained_variance_ratio[1]*100:.3f}%')
```
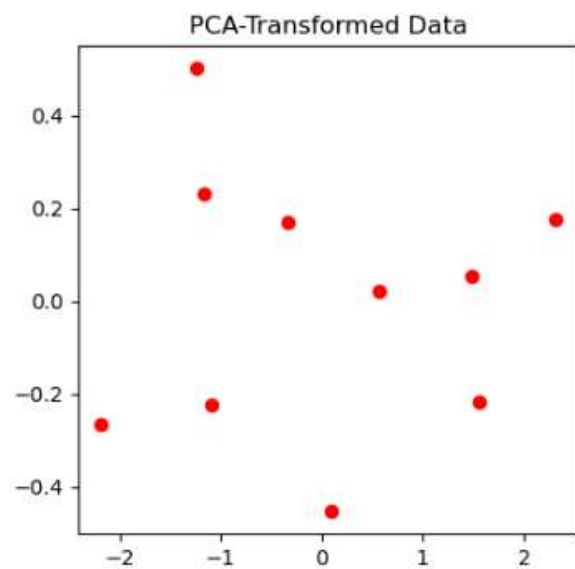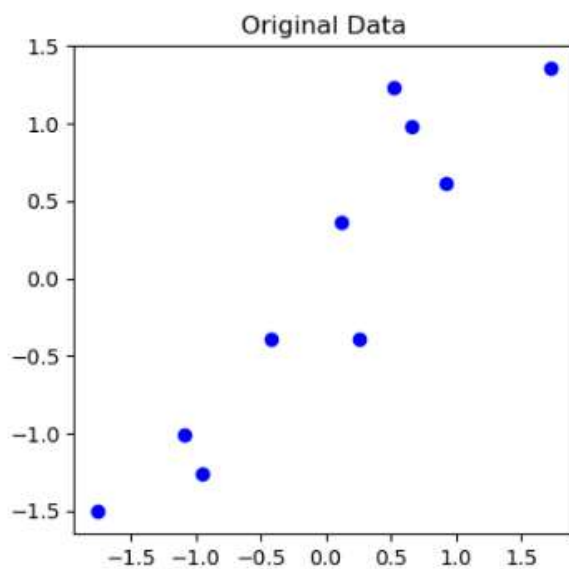
```
Standardized data is:
 [[ 0.92627881  0.61016865]
 [-1.7585873  -1.506743  ]
 [ 0.52354889  1.23278973]
 [ 0.12081898  0.36112022]
 [ 1.73173864  1.35731394]
 [ 0.6577922   0.9837413 ]
 [ 0.25506228 -0.38602507]
 [-1.08737078 -1.00864614]
 [-0.41615425 -0.38602507]
 [-0.95312747 -1.25769457]]
PCA(n_components=2)
Principal component Analysis i
 [[-1.08643242 -0.22352364]
 [ 2.3089372   0.17808082]
 [-1.24191895  0.501509  ]
 [-0.34078247  0.16991864]
 [-2.18429003 -0.26475825]
 [-1.16073946  0.23048082]
 [ 0.09260467 -0.45331721]
 [ 1.48210777  0.05566672]
 [ 0.56722643  0.02130455]
 [ 1.56328726 -0.21536146]]
```



Explained Variance Ratios: 96.296% and 3.704%

```python
1  import numpy as np
2  from sklearn.decomposition import PCA
3  import matplotlib.pyplot as plt
4  # Given data
5  data = np.array([[4, 11], [8, 4], [13, 5], [7, 14]])
6  # Initialize PCA
7  pca = PCA(n_components=2)
8  # Fit the PCA model to your data
9  pca.fit(data)
10 # The principal components
11 print("Principal components:")
12 print(pca.components_)
13 # The explained variance
14 print("Explained variance:")
15 print(pca.explained_variance_)
16 # Transform the data to the first principal components
17 transformed_data = pca.transform(data)
18 print("Transformed data:")
19 print(transformed_data)
20 # Create a scatter plot of the transformed data
21 plt.figure(figsize=(6,6))
22 plt.scatter(transformed_data[:, 0], transformed_data[:, 1])
23 plt.title('PCA result')
24 plt.xlabel('First Principal Component')
25 plt.ylabel('Second Principal Component')
26 plt.grid(True)
27 plt.show()
```

```
Principal components:
[[ 0.55738997 -0.83025082]
 [-0.83025082 -0.55738997]]
Explained variance:
[30.38486432  6.61513568]
Transformed data:
[[-4.30518692  1.92752836]
 [ 3.73612869  2.50825486]
 [ 5.69282771 -2.20038921]
 [-5.12376947 -2.23539401]]
```
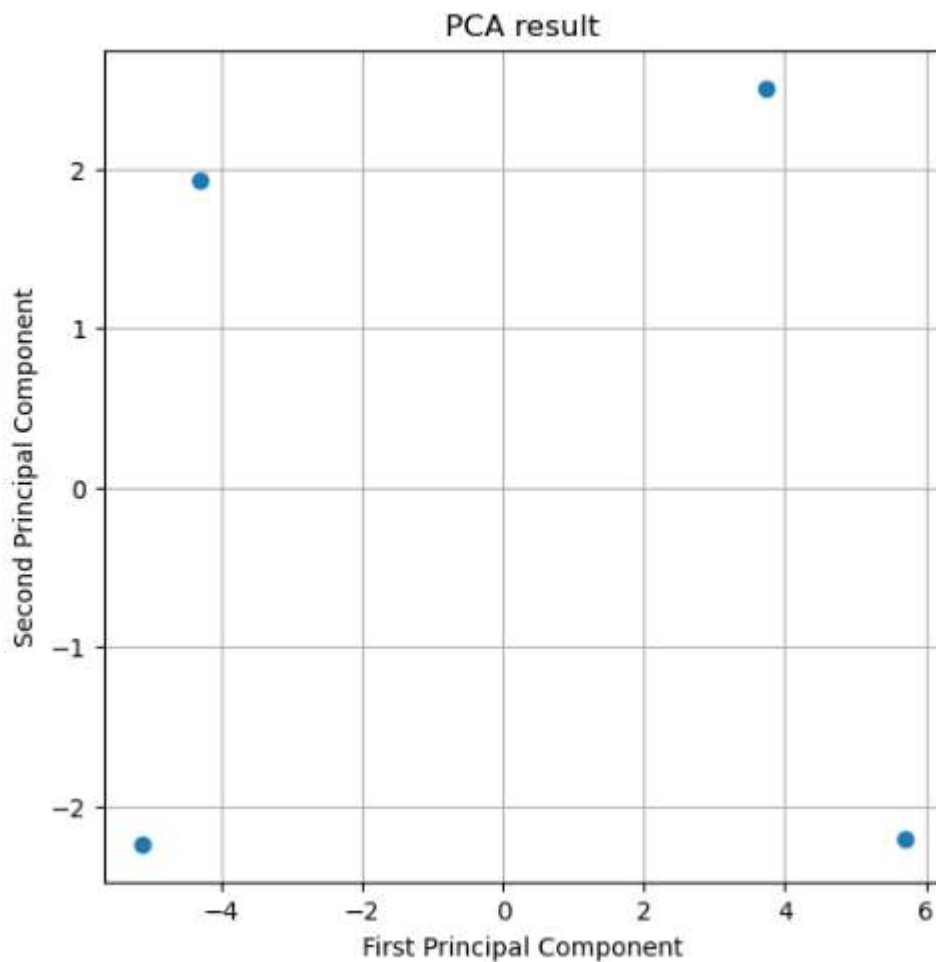
PCA result

## PCA has several applications:

Dimensionality Reduction: It reduces the number of variables in a dataset, making it easier to analyse and visualize while retaining most of the important information.

Feature Engineering: In machine learning, PCA can be used for feature selection or feature extraction to improve model performance.

Noise Reduction: PCA can help remove noise from data by focusing on the principal components with the highest variance.

Data Visualization: PCA can be used to visualize high-dimensional data in a lower-dimensional space, allowing for easier data exploration and pattern recognition.

**NOTE**

Keep in mind that while PCA is a powerful technique, it's not always suitable for every dataset. The choice of the number of principal components to retain should be made carefully, and the interpretation of the reduced dimensions might not always be straightforward. PCA is also sensitive to outliers, so data preprocessing and outlier handling are essential.