# Principal Component Analysis

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3  import pandas as pd
4  import random as rd
```

## 1. Data

- the data are given by the file data-pca.txt
- the data consist of a set of points $\{(x_i, y_i)\}_{i=1}^{n}$ where $z_i = (x_i, y_i)$ denotes a 2-dimensional point in the cartesian coordinate

load the data from the files

```
1  path = '/content/drive/My Drive/ML_Assignment/data/data-pca.txt'
2  data = np.loadtxt(path, delimiter=',')
3  x = data[:,0]
4  y = data[:,1]
```

Plot the original data points

```
1  fig_1 = plt.figure(figsize = (6,6))
2  plt.scatter(x, y, c='r', marker = '+')
3  plt.title('original data points')
4  plt.show()
5  fig_1.savefig('original data points.png')
```

## 2. Normalization

- the data is normalized to have the mean = 0 and the standard deviation = 1
- $x = \frac{x - \mu_x}{\sigma_x}$ and $y = \frac{y - \mu_y}{\sigma_y}$

  - $\mu_x$ denotes the mean of $x$
  - $\sigma_x$ denotes the standard deviation of $x$
  - $\mu_y$ denotes the mean of $y$
  - $\sigma_y$ denotes the standard deviation of $y$

define a function to normalize the input data points $x$ and $y$
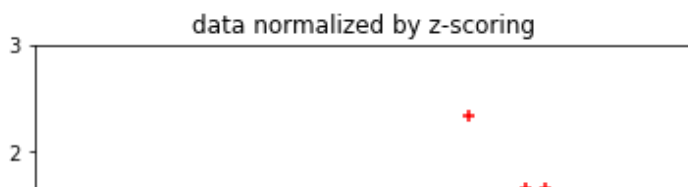
```
1   def normalize_data(x, y):
2
3       xn = (x - x.mean(axis=0)) / x.std(axis=0) # normalize x. the mean of xn is zero and the standard devia
4       yn = (y - y.mean(axis=0)) / y.std(axis=0) # normalize y. the mean of yn is zero and the standard devia
5
6       return xn, yn
```

Plot the normalized data points

- $z = \frac{z - \mu}{\sigma}$
- $\mu$ denotes the average and $\sigma$ denotes the standard deviation

```
1   xn, yn = normalize_data(x, y)
```

```
1   fig_2 = plt.figure(figsize = (6,6))
2   plt.scatter(xn, yn, c='r', marker = '+')
3   plt.title('data normalized by z-scoring')
4   plt.axis([-3, 3, -3, 3])
5   plt.show()
6   fig_1.savefig('normalized data points.png')
```

## 3. Covariance Matrix

- compute the co-variance matrix
- $\Sigma = \frac{1}{n} \sum_{i=1}^{n} z_i z_i^T = \frac{1}{n} Z^T Z$

  - $n$ denotes the number of data

  - $Z = \begin{bmatrix} z_1^T \\ \vdots \\ z_n^T \end{bmatrix}$

define a function to compute the co-variance matrix of the data

```
1   def compute_covariance(z):
2       # compute the covariance matrix #
3       covar = np.cov(z.T)
4       return covar
```

```
1   Z = np.c_[xn,yn]
2   covariance = compute_covariance(Z)  # return 2x2 metrix
```

## 4. Principal Components

- compute the eigen-values and the eigen-vectors of the co-variance matrix

define a function to compute the principal directions from the co-variance matrix

```
1   def compute_principal_direction(covariance):
2
3       e_value, e_vector = np.linalg.eig(covariance) # compute the principal directions from the co-variance 
4
5       return e_value, e_vector
```
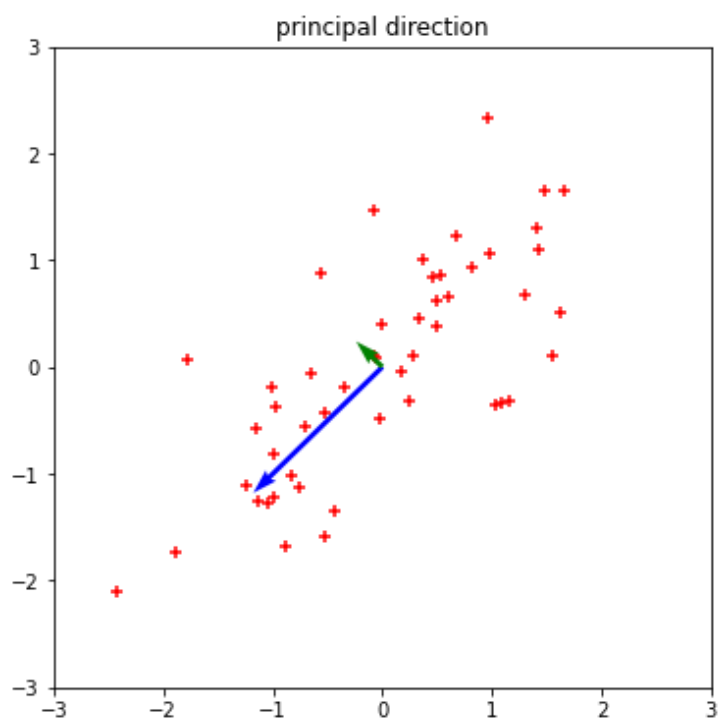
```
1   e_value, e_vector = compute_principal_direction(covariance)
```

```
1   # 내림차순 정렬
2   idx = np.flip(e_value.argsort())
3   e_value = e_value[idx]
4   e_vector = e_vector[:, idx]
```

## Plot the principal axes

- plot the normalized data points
- plot the first principal vector
- plot the second principal vector

```
1   fig_3 = plt.figure(figsize = (6,6))
2   plt.scatter(xn, yn, c='r', marker = '+')
3   plt.quiver([0], [0], e_vector[0, 0], e_vector[1, 0], color=['b'],angles='xy', scale_units='xy', scale=0.6)
4   plt.quiver([0], [0], e_vector[0, 1], e_vector[1, 1], color=['g'],angles='xy', scale_units='xy')
5   plt.title('principal direction')
6   plt.axis([-3, 3, -3, 3])
7   plt.show()
8   fig_3.savefig('principal direction.png')
```



## Make LInear example

```
1   def VectorToLinear(vector):
2     a = vector[1]/vector[0]
3     test_x = np.arange(-2.8, 2.8, 0.1)
4     test_y = a * test_x
5     return test_x, test_y
```

## define a function to compute the projection of the data point onto the principal axis

```
1   def compute_projection(point, axis):
2       # compute the projection of point on the axis #
3       pca = np.dot(point, axis)
4       trans_e_vector = axis[:,None].T
5       projection = np.dot(pca[:,None], trans_e_vector)
6       return projection
```
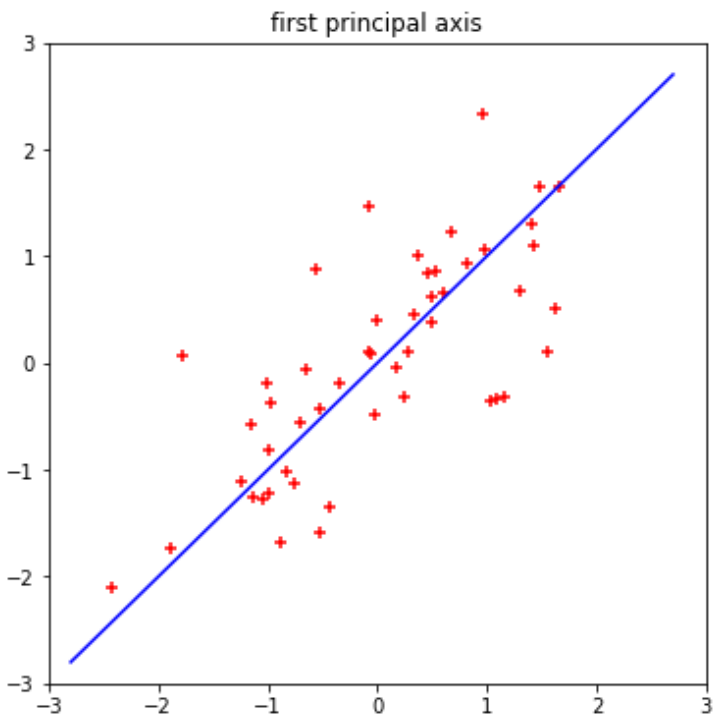
# Plot Output principal Axis

## First principal axis

Plot the first principal axis

- plot the normalized data points
- plot the first principal axis

```
1    first_x , first_y = VectorToLinear(e_vector[:,0].T)
```

```
1    fig_4 = plt.figure(figsize = (6,6))
2    plt.scatter(xn, yn, c='r', marker = '+')
3    plt.title('first principal axis')
4    plt.plot(first_x, first_y, c='b')
5    plt.axis([-3, 3, -3, 3])
6    plt.show()
7    fig_4.savefig('first principal axis.png')
```



```
1    first_pca = compute_projection(Z, e_vector[:,0])
```

Plot the project of the normalized data points onto the first principal axis
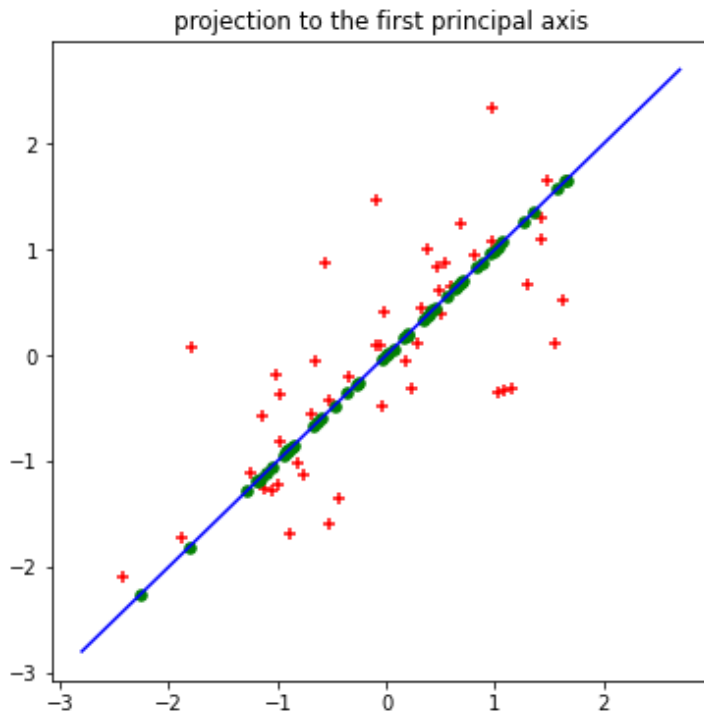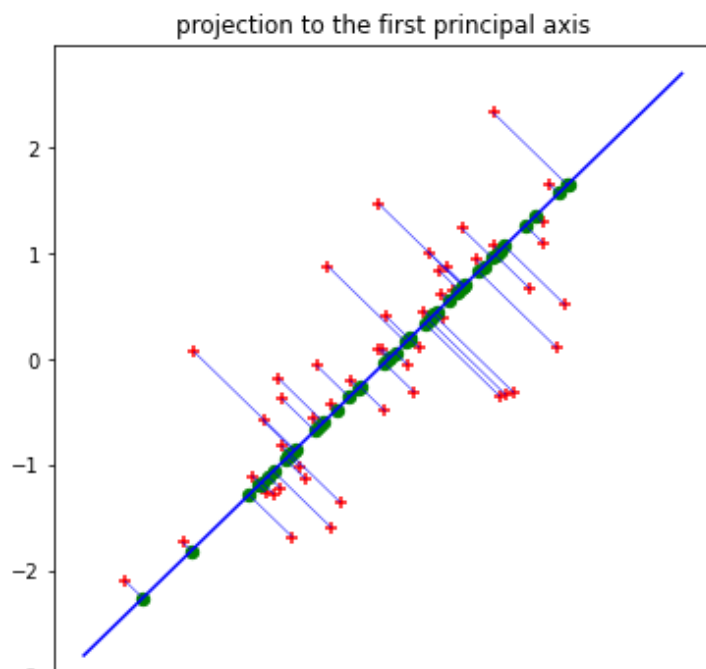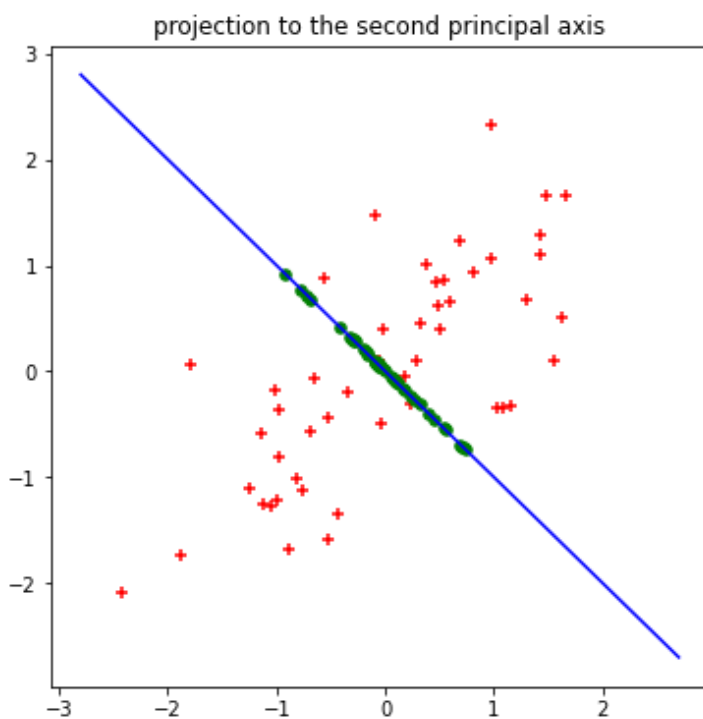
- plot the normalized data points
- plot the first principal axis
- plot the projected points from the normalized data points onto the first principal axis

```
1    fig_5 = plt.figure(figsize = (6,6))
2    plt.scatter(xn, yn, c='r', marker = '+',s=30)
```

```
3   zero = np.zeros(50)
4   plt.plot(first_x, first_y, c='b')
5   plt.scatter(first_pca[:,0], first_pca[:,1], c = 'g', s=30)
6   plt.title('projection to the first principal axis')
7   plt.show()
8   fig_5.savefig('projection to the first principal.png')
```



projection to the first principal axis

Plot the lines between the normalized data points and their projection points on the first principal axis

- plot the normalized data points
- plot the first principal axis
- plot the projected points from the normalized data points onto the first principal axis
- plot the lines that connect between the normalized data points and their projection points on the first principal axis

```
1   # 유클리드 좌표계에서 점사이의 거리
2   def compute_distance(x1, y1, x2, y2):
3       distance = np.array([[x1, y1],[x2, y2]])
4       return distance
```

```
1    fig_6 = plt.figure(figsize = (6,6))
2    plt.scatter(xn, yn, c='r', marker = '+',s=30)
3    zero = np.zeros(50)
4    plt.plot(first_x, first_y, c='b')
5    plt.scatter(first_pca[:,0], first_pca[:,1], c = 'g', s=40)
6    for i in range (50):
7        dist = np.array([[xn[i], yn[i]],[first_pca[i,0], first_pca[i,1]]])
8        plt.plot(dist[:,0], dist[:,1], c = 'b', linewidth = 0.5)
9    plt.title('projection to the first principal axis')
10   plt.show()
11   fig_6.savefig('distance to the first projection axis.png')
```

projection to the first principal axis

## Second principal axis

Plot the second principal axis

- plot the normalized data points
- plot the second principal axis

```
1   second_x , second_y = VectorToLinear(e_vector[:,1].T)
```

```
1   fig_7 = plt.figure(figsize = (6,6))
2   plt.scatter(xn, yn, c='r', marker = '+')
3   plt.title('second principal axis')
4   plt.plot(second_x , second_y, c='b')
5   plt.axis([-3, 3, -3, 3])
6   plt.show()
7   fig_7.savefig('second principal axis.png')
```

Plot the project of the normalized data points onto the second principal axis

- plot the normalized data points
- plot the second principal axis
- plot the projected points from the normalized data points onto the second principal axis

```
1   second_pca = compute_projection(Z, e_vector[:,1])
```

```
1   fig_8 = plt.figure(figsize = (6,6))
2   plt.scatter(xn, yn, c='r', marker = '+',s=30)
3   zero = np.zeros(50)
4   plt.plot(second_x , second_y, c='b')
5   plt.scatter(second_pca[:,0], second_pca[:,1], c = 'g', s=30)
6   plt.title('projection to the second principal axis')
7   plt.show()
8   fig_8.savefig('projection to the second principal.png')
```



Plot the lines between the normalized data points and their projection points on the second principal axis

- plot the normalized data points
- plot the second principal axis
- plot the projected points from the normalized data points onto the second principal axis
- plot the lines that connect between the normalized data points and their projection points on the second principal axis

```
1   fig_9 = plt.figure(figsize = (6,6))
2   plt.scatter(xn, yn, c='r', marker = '+',s=30)
3   zero = np.zeros(50)
4   plt.plot(second_x, second_y, c='b')
5   plt.scatter(second_pca[:,0], second_pca[:,1], c = 'g', s=40)
6   for i in range (50):
```
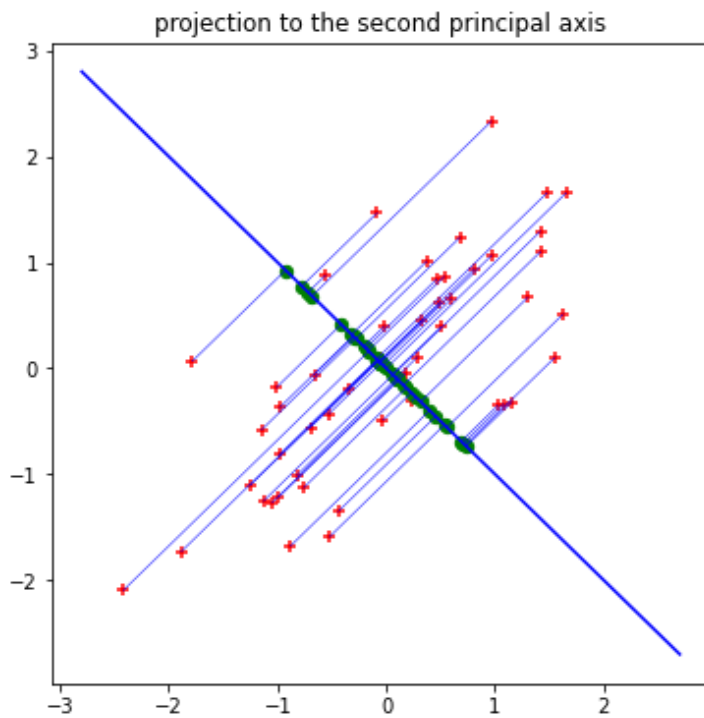
```
7          dist = np.array([[xn[i], yn[i]],[second_pca[i,0], second_pca[i,1]]])
8          plt.plot(dist[:,0], dist[:,1], c = 'b', linewidth = 0.5)
9     plt.title('projection to the second principal axis')
10    plt.show()
11    fig_9.savefig('distance to the second projection axis.png')
```
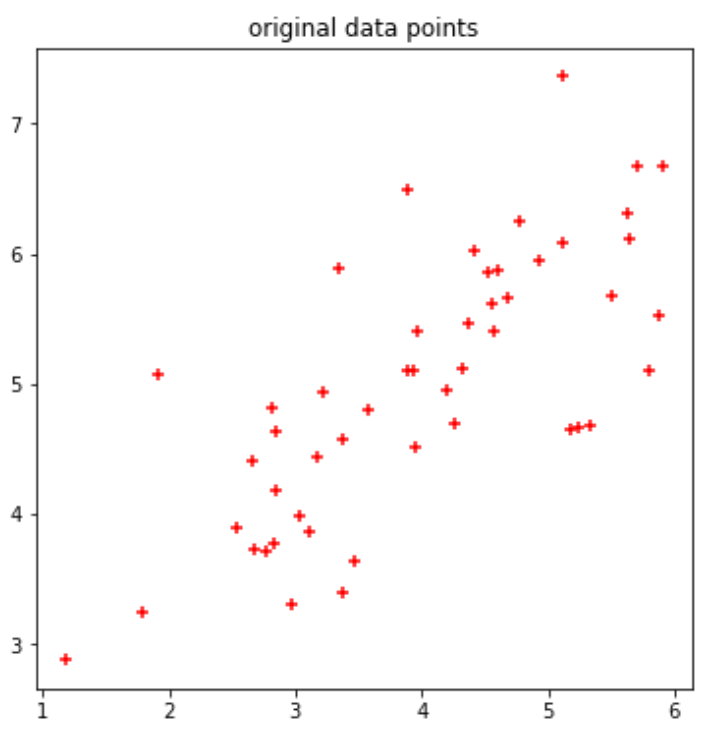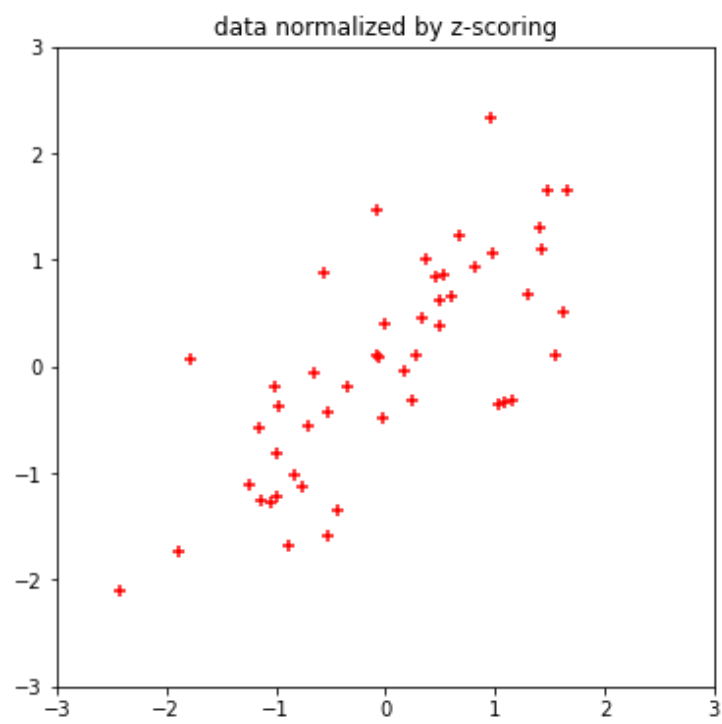


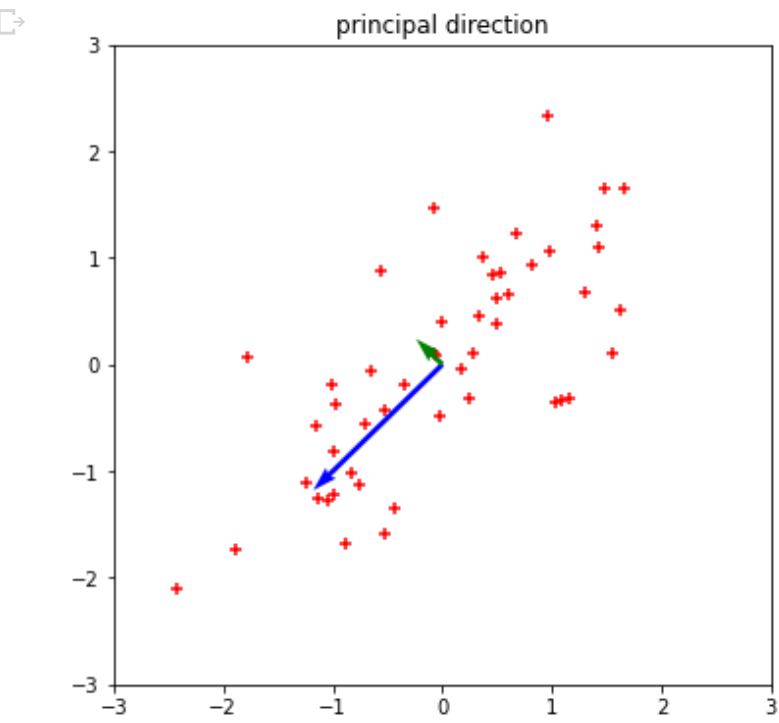## Output

## 1. Plot the original data points [1pt]

```
1    fig_1
```

## 2. Plot the normalized data points [1pt]

```
1    fig_2
```



data normalized by z-scoring

## 3. Plot the principal axes [2pt]

```
1    fig_3
```



principal direction

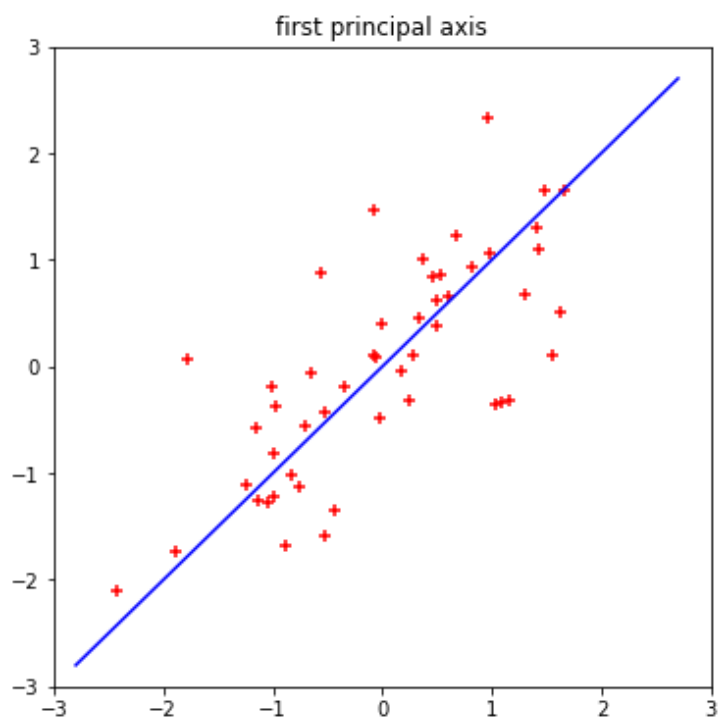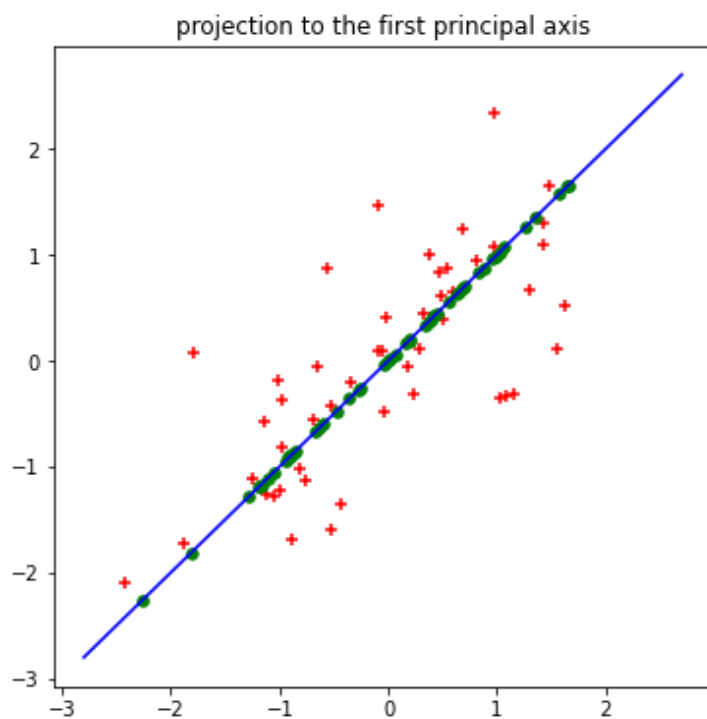## 4. Plot the first principal axis [3pt]

1    fig_4



first principal axis

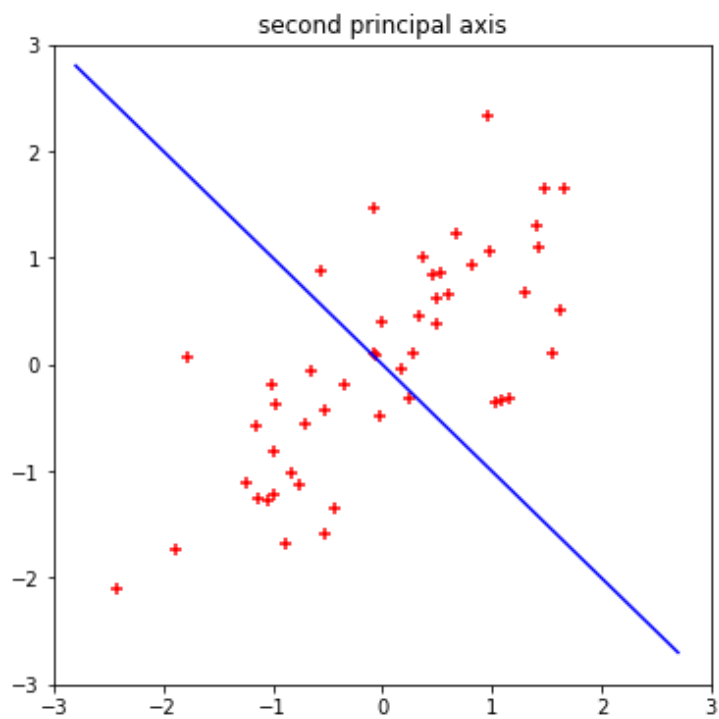5. Plot the project of the normalized data points onto the first principal axis [4pt]

1    fig_5



projection to the first principal axis

6. Plot the lines between the normalized data points and their projection points on the first principal axis [3pt]

projection to the first principal axis
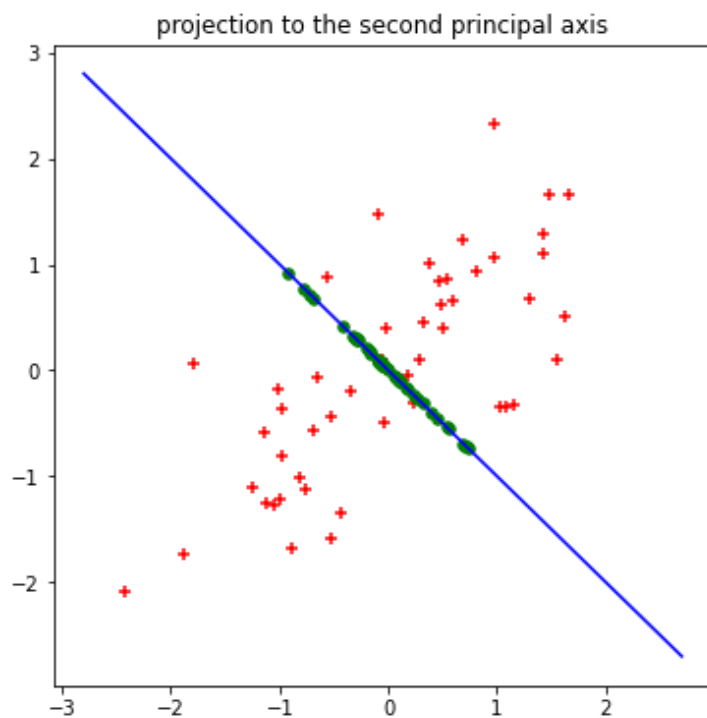
## 7. Plot the second principal axis [3pt]

second principal axis

## 8. Plot the project of the normalized data points onto the second principal axis [4pt]

projection to the second principal axis

9. Plot the lines between the normalized data points and their projection points on the second principal axis [3pt]

1    fig_9



projection to the second principal axis