

## Principal Component Analysis (PCA)

1. Consider the set of points in the plane:

$$x_1 = (4, 3)^T, x_2 = (-4, 6)^T, x_3 = (7, -2)^T, x_4 = (1, 1)^T, x_5 = (0, -2)^T$$

- Set up a corresponding data matrix  $X \in \mathbb{R}^{2 \times 5}$ .
- Find the principal components of  $X$  (remember to center the data points first).
- Compute the projections  $y_1, \dots, y_5 \in \mathbb{R}^1$  of  $x_1, \dots, x_5 \in \mathbb{R}^2$  on the first principal component (remember to correct for the mean).
- Compute the reconstructions of  $x_1, \dots, x_5$  using the first principal components, denoted  $\hat{x}_1, \dots, \hat{x}_5 \in \mathbb{R}^2$ .
- In a 2D axis system, plot the following:
  - The original points  $x_1, \dots, x_5$ .
  - The reconstructed points  $\hat{x}_1, \dots, \hat{x}_5$ .
  - The principal components (directions).

2. Consider the following points in  $\mathbb{R}^5$ :

$$x_1 = (2, 3, 1, 0, -1)^T, x_2 = (-4, 2, 4, 4, 1)^T, x_3 = (-1, -1, 3, 9, 0)^T.$$

Find the projections  $y_1, y_2, y_3 \in \mathbb{R}^2$  of these points, that makes them uncorrelated and with maximal variance (remember to center the points first).

## Solution

1. (a) The data matrix is:

$$X = \begin{pmatrix} 4 & -4 & 7 & 1 & 0 \\ 3 & 6 & -2 & 1 & -2 \end{pmatrix}$$

- (b) The mean is:

$$\bar{x}^T = \frac{1}{5} ((4, 3) + (-4, 6) + (7, -2) + (1, 1) + (0, -2)) = (8/5, 6/5).$$

The centered data matrix is:

$$X' = X - (8/5, 6/5)^T = \frac{1}{5} \begin{pmatrix} 12 & -28 & 27 & -3 & -8 \\ 9 & 24 & -16 & -1 & -16 \end{pmatrix}.$$

The SVD decomposition of  $X'$  is  $X' = U\Sigma V^T$  where

$$U = \begin{pmatrix} -0.8087 & -0.5882 \\ 0.5882 & -0.8087 \end{pmatrix},$$
$$\Sigma = \begin{pmatrix} 9.7143 & 0 & 0 & 0 & 0 \\ 0 & 4.6511 & 0 & 0 & 0 \end{pmatrix},$$
$$V = \begin{pmatrix} -0.0908 & -0.6165 & -0.3752 & 0.1174 & 0.6761 \\ 0.7568 & -0.1263 & 0.5664 & -0.0002 & 0.3008 \\ -0.6433 & -0.1266 & 0.7282 & 0.0455 & 0.1944 \\ 0.0378 & 0.1107 & 0.0073 & 0.9912 & -0.0620 \\ -0.0606 & 0.7587 & -0.0902 & -0.0417 & 0.6409 \end{pmatrix}$$

(we don't really need  $\Sigma$  and  $V$  here). The principal components are therefore

$$\hat{u}_1 = (-0.8087, 0.5882)^T, \hat{u}_2 = (-0.5882, -0.8087)^T.$$

- (c) The projections are given by

$$Y = \hat{u}_1^T X' = (-0.8820, 7.3522, -6.2493, 0.3676, -0.5884).$$

So that

$$y_1 = -0.8820, y_2 = 7.3522, y_3 = -6.2493, y_4 = 0.3676, y_5 = -0.5884.$$

- (d) The (centered) reconstructions are given by

$$\hat{u}_1 Y = \hat{u}_1 \hat{u}_1^T X' = \begin{pmatrix} 0.7133 & -5.9457 & 5.0537 & -0.2973 & 0.4759 \\ -0.5188 & 4.3248 & -3.6760 & 0.2162 & -0.3461 \end{pmatrix}.$$

In other words, the reconstructed points (after correcting for the mean):

$$\hat{X} = \begin{pmatrix} 2.3133 & -4.3457 & 6.6537 & 1.3027 & 2.0759 \\ 0.6812 & 5.5248 & -2.4760 & 1.4162 & 0.8539 \end{pmatrix}.$$

In other words,  $\hat{x}_i$  is the  $i$ -th column of  $\hat{X}$ ,

$$\hat{x}_1 = (2.3133, 0.6812)^T$$

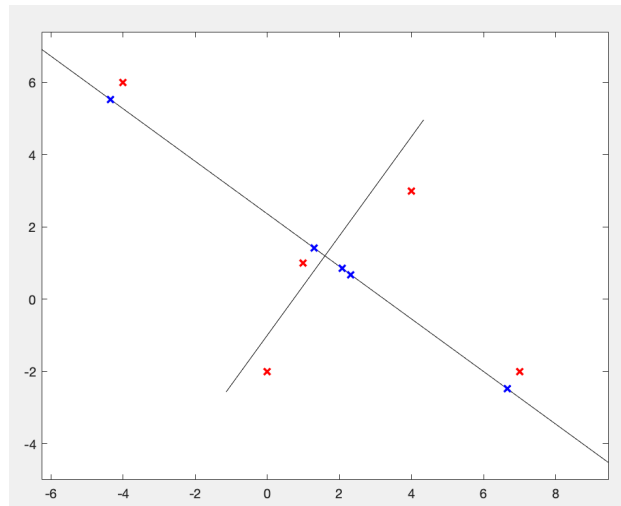
$$\hat{x}_2 = (-4.3457, 5.5248)^T$$

$$\hat{x}_3 = (6.6537, -2.4760)^T$$

$$\hat{x}_4 = (1.3027, 1.4162)^T$$

$$\hat{x}_5 = (2.0759, 0.8539)^T$$

(e) The red points are  $x_1, \dots, x_5$ , the blue are  $\hat{x}_1, \dots, \hat{x}_5$ , and the black lines are the principal components.



2. We first write the data matrix:

$$X = \begin{pmatrix} 2 & -4 & -1 \\ 3 & 2 & -1 \\ 1 & 4 & 3 \\ 0 & 4 & 9 \\ -1 & 1 & 0 \end{pmatrix}$$

The mean is  $\bar{x}^T = \frac{1}{3}(-3, 4, 8, 13, 0)^T$ . Therefore,

$$X' = \begin{pmatrix} 3.0000 & -3.0000 & 0 \\ 1.6667 & 0.6667 & -2.3333 \\ -1.6667 & 1.3333 & 0.3333 \\ -4.3333 & -0.3333 & 4.6667 \\ -1.0000 & 1.0000 & 0 \end{pmatrix}$$

The matrix  $U$  in the SVD of  $X'$  is:

$$U = \begin{pmatrix} -0.3559 & -0.7867 & -0.2949 & -0.2701 & 0.3076 \\ -0.3546 & 0.2991 & -0.7930 & 0.3898 & -0.0634 \\ 0.2201 & 0.3319 & -0.3832 & -0.8334 & 0.0042 \\ 0.8277 & -0.3359 & -0.3691 & 0.2544 & -0.0320 \\ 0.1186 & 0.2622 & 0.0318 & 0.1259 & 0.9488 \end{pmatrix}$$

The projection on the first 2 PCs is given by:

$$Y = \hat{U}_2 X' = \begin{pmatrix} -5.7308 & 0.9674 & 4.7634 \\ -1.2212 & 3.3760 & -2.1548 \end{pmatrix},$$

where

$$\hat{U}_2 = \begin{pmatrix} -0.3559 & -0.7867 \\ -0.3546 & 0.2991 \\ 0.2201 & 0.3319 \\ 0.8277 & -0.3359 \\ 0.1186 & 0.2622 \end{pmatrix}.$$

In other words,

$$y_1 = (-5.7308, -1.2212)^T, y_2 = (0.9674, 3.3760)^T, y_3 = (4.7634, -2.1548)^T.$$