

Principle Component Analysis - Example

NIRDOSH BHATNAGAR

This example is on principal component analysis (PCA). The set of data points is $\{x_1, x_2, x_3\}$. These are

$$x_1 = \begin{bmatrix} 1 & 0 \end{bmatrix}^T, \quad x_2 = \begin{bmatrix} 0 & -1 \end{bmatrix}^T, \quad x_3 = \begin{bmatrix} -1 & 1 \end{bmatrix}^T$$

Determine the principal component of these data points. Also determine the corresponding error in the approximation. The value of k (number of largest eigenvalues used in the approximation of a data point) is 1.

Hint: The student is expected to complete the intermediate steps / calculations.

Step 1: The number of data points is $n = 3$. Let the average value of the data points be

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \begin{bmatrix} 0 & 0 \end{bmatrix}^T$$

Define a column vector

$$\theta_i = (x_i - \bar{x}), \quad 1 \leq i \leq 3$$

and a 2×3 matrix Θ as

$$\Theta = \begin{bmatrix} \theta_1 & \theta_2 & \theta_3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & -1 \\ 0 & -1 & 1 \end{bmatrix}$$

Step 2: The covariance matrix $\tilde{\Sigma}$ of the data points is

$$\tilde{\Sigma} = \frac{1}{(n-1)} \Theta \Theta^T = \begin{bmatrix} 1 & -1/2 \\ -1/2 & 1 \end{bmatrix}$$

The 2×2 covariance matrix $\tilde{\Sigma}$ is symmetric and positive-definite (in this problem). Why?

Step 3: Therefore $\tilde{\Sigma} = \Psi \Lambda \Psi^T$ where Λ is a diagonal matrix with eigenvalues of the matrix $\tilde{\Sigma}$ on its main diagonal. As the matrix $\tilde{\Sigma}$ is symmetric and positive-definite, its eigenvalues are positive. Let the eigenvalues of the matrix $\tilde{\Sigma}$ be $\lambda_j > 0, 1 \leq j \leq 2$. The columns of the matrix Ψ are the mutually orthogonal eigenvectors of the covariance matrix $\tilde{\Sigma}$. Assume that these eigenvectors are orthonormal. Therefore $\Psi \Psi^T = I$, where I is an identity matrix of size 2. We determine the eigenvalues λ_1 and λ_2 and the corresponding eigenvectors ψ_1 and ψ_2 in this step.

(a) The eigenvalues of the matrix $\tilde{\Sigma}$ are $\lambda_1 = 1/2$, and $\lambda_2 = 3/2$. Note that these eigenvalues are positive.

The corresponding eigenvectors of length one are

$$\psi_1 = \pm \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \end{bmatrix}^T, \quad \text{and} \quad \psi_2 = \pm \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \end{bmatrix}^T$$

respectively. Observe that both the eigenvectors have a length of unity, and $\psi_1^T \psi_2 = 0$. That is, these vectors are orthonormal.

(b) Thus

$$\Psi = [\psi_1 \quad \psi_2] = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

It can be verified that $\Psi\Psi^T = I$, where I is an identity matrix of size 2.

(c) Verify by direct calculation, that $\Psi\Lambda\Psi^T = \tilde{\Sigma}$.

Step 4: As $\tilde{\Sigma} = \Psi\Lambda\Psi^T$, we have

$$\begin{aligned} \Lambda &= \Psi^T \tilde{\Sigma} \Psi = \frac{1}{(n-1)} \Psi^T \Theta \Theta^T \Psi \\ &= \frac{1}{(n-1)} W W^T \end{aligned}$$

where

$$\begin{aligned} W &= \Psi^T \Theta = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 0 & -1 & 1 \end{bmatrix} \\ &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 & 0 \\ 1 & 1 & -2 \end{bmatrix} \end{aligned}$$

is a 2×3 matrix. Also

$$\Theta = \Psi W$$

Let

$$W = [w_1 \quad w_2 \quad w_3]$$

Thus

$$\theta_i = \Psi w_i, \quad 1 \leq i \leq 3$$

Step 5: We next approximate Θ by $\hat{\Theta}$, where

$$\hat{\Theta} = [\hat{\theta}_1 \quad \hat{\theta}_2 \quad \hat{\theta}_3]$$

Recall that eigenvalues of the covariance matrix $\tilde{\Sigma}$ are $\lambda_1 = 1/2$, and $\lambda_2 = 3/2$. As $\lambda_2 > \lambda_1$, we approximate $\Theta = \Psi W$ as

$$\hat{\Theta} = \hat{\Psi} W$$

where

$$\hat{\Psi} = \frac{1}{\sqrt{2}} \begin{bmatrix} 0 & 1 \\ 0 & -1 \end{bmatrix}$$

Thus

$$\begin{aligned} \hat{\Theta} &= \hat{\Psi} W = \frac{1}{2} \begin{bmatrix} 0 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & -1 & 0 \\ 1 & 1 & -2 \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} 1 & 1 & -2 \\ -1 & -1 & 2 \end{bmatrix} \end{aligned}$$

Therefore

$$\hat{\Theta} = \frac{1}{2} \begin{bmatrix} 1 & 1 & -2 \\ -1 & -1 & 2 \end{bmatrix} = \begin{bmatrix} \hat{\theta}_1 & \hat{\theta}_2 & \hat{\theta}_3 \end{bmatrix}$$

We summarize θ_i and $\hat{\theta}_i$, for $1 \leq i \leq 3$ as

$$\begin{aligned} \theta_1 &= \begin{bmatrix} 1 & 0 \end{bmatrix}^T, \quad \text{and} \quad \hat{\theta}_1 = \begin{bmatrix} 1/2 & -1/2 \end{bmatrix}^T \\ \theta_2 &= \begin{bmatrix} 0 & -1 \end{bmatrix}^T, \quad \text{and} \quad \hat{\theta}_2 = \begin{bmatrix} 1/2 & -1/2 \end{bmatrix}^T \\ \theta_3 &= \begin{bmatrix} -1 & 1 \end{bmatrix}^T, \quad \text{and} \quad \hat{\theta}_3 = \begin{bmatrix} -1 & 1 \end{bmatrix}^T \end{aligned}$$

Therefore, the mean-squared error E_{err} in the approximation of a data point is

$$E_{err} = \frac{1}{(3-1)} \sum_{i=1}^3 \left\| \theta_i - \hat{\theta}_i \right\|^2$$

where $\|\cdot\|$ is the Euclidean norm. It is

$$E_{err} = \frac{1}{2} = \lambda_1$$

Step 6: Therefore the principal component is $\lambda_2 = 3/2$, and its direction is given by the eigenvector ψ_2 . The error in the approximation is $E_{err} = 0.5$. \square