

Lecture Note on Principal Component Analysis

This version: September 19, 2022

1 Covariance

Definition 1. Let $\mathbf{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_p \end{pmatrix}$ be a random vector. The covariance matrix of \mathbf{X} , denoted as Σ , is defined as

$$\Sigma_{p \times p} = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])'] = \left(\sigma_{ij} \right)_{p \times p},$$

where

$$\sigma_{i,j}^2 = \begin{cases} \text{Var}(X_i), & \text{if } i = j, \\ \text{Cov}(X_i, X_j), & \text{if } i \neq j. \end{cases}$$

Example 1. Let Z_1 and Z_2 be two independent $N(0, 1)$. Suppose that $X_1 = Z_1 + 2Z_2$ and $X_2 = Z_2$. Then

$$\text{Var}(X_1) = \text{Var}(Z_1 + 2Z_2) = \text{Var}(Z_1) + 4\text{Var}(Z_2) = 1 + 4 = 5$$

$$\text{Var}(X_2) = \text{Var}(Z_2) = 1$$

$$\text{Cov}(X_1, X_2) = \mathbb{E}[(X_1 - \mathbb{E}[X_1])(X_2 - \mathbb{E}[X_2])]$$

$$= \mathbb{E}[X_1 X_2]$$

$$= \mathbb{E}[(Z_1 + 2Z_2)(Z_2)]$$

$$= \mathbb{E}[2Z_2^2]$$

(Z_1, Z_2 independent)

$$= 2(\text{Var}(Z_2) + (\mathbb{E}[Z_2])^2) = 2$$

Hence, the covariance matrix of (X_1, X_2) is

$$\Sigma = \begin{pmatrix} 5 & 2 \\ 2 & 1 \end{pmatrix}.$$

Remark 1. The off-diagonal terms of Σ are $\text{Cov}(X_i, X_j)$ for any $1 \leq i \neq j \leq p$, so if each pair of components of the random vector \mathbf{X} is mutually independent, then Σ is a diagonal matrix.

Remark 2. In matrix notation,

$$\Sigma = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])'].$$

Remark 3. Because

$$\Sigma' = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])'] = \Sigma,$$

Σ is symmetric. Or, in scalar notation,

$$\sigma_{i,j} = \text{Cov}(X_i, X_j) = \text{Cov}(X_j, X_i) = \sigma_{j,i}.$$

2 Principal Component Analysis (PCA)

PCA attempts to summarize variation in the random vector \mathbf{X} with few principal components (PC). In below, we will define what are *principal components* and explain how to derive them.

The first PC, PC_1 , can be obtained from the following maximization problem:

$$\max_{\mathbf{a}_1 \in \mathbb{R}^p} \text{Var}(\mathbf{a}_1' \mathbf{X}) \quad \text{s.t.} \quad \mathbf{a}_1' \mathbf{a}_1 = 1.$$

The first principal component $PC_1 = \mathbf{a}_1' \mathbf{X} = a_{11}X_1 + a_{12}X_2 + \cdots + a_{1p}X_p$, where \mathbf{a}_1 is the solution to the above problem. We call \mathbf{a}_1 the **coefficient** of the first principal component.

The constraint requires \mathbf{a} has to be length 1:

$$\mathbf{a}'\mathbf{a} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1p} \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \\ \vdots \\ a_{1p} \end{pmatrix} = a_{11}^2 + a_{12}^2 + \dots + a_{1p}^2 = 1.$$

Remark 4. The restriction $\mathbf{a}'_i\mathbf{a}_i = 1$ is necessary, otherwise $\text{Var}(\mathbf{a}'_iX)$ can be made arbitrary large by multiplying the coefficient with a constant.

The rest of the principal components can be defined iteratively.¹ The second PC, $PC_2 = \mathbf{a}'_2\mathbf{x}$, has coefficients \mathbf{a}_2 given by the optimization problem below:

$$\begin{aligned} \max_{\mathbf{a}_2 \in \mathbb{R}^p} \quad & \text{Var}(\mathbf{a}'_2\mathbf{X}) \\ \text{s.t.} \quad & \mathbf{a}'_2\mathbf{a}_2 = 1, \\ & \text{Cov}(\mathbf{a}'_1\mathbf{X}, \mathbf{a}'_2\mathbf{X}) = 0. \end{aligned}$$

Notice that we require the second PC has to be uncorrelated with the first PC. We can think of it as we are trying to have the second PC explain the variation that is not explained by the first PC.

Similarly, the i th PC, $\mathbf{a}'_i\mathbf{x}$ is the solution to

$$\begin{aligned} \max_{\mathbf{a}_i \in \mathbb{R}^p} \quad & \text{Var}(\mathbf{a}'_i\mathbf{X}) \\ \text{s.t.} \quad & \mathbf{a}'_i\mathbf{a}_i = 1, \\ & \text{Cov}(\mathbf{a}'_j\mathbf{X}, \mathbf{a}'_i\mathbf{X}) = 0 \text{ for } j = 1, 2, \dots, i-1. \end{aligned}$$

Remark 5. Recall that

$$\mathbf{a}'_i\mathbf{x} = \langle \mathbf{a}, \mathbf{x} \rangle = \begin{pmatrix} a_{i1} & a_{i2} & \dots & a_{ip} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} = a_{i1}x_1 + a_{i2}x_2 + \dots + a_{ip}x_p,$$

¹If the random vector \mathbf{X} contains p variables, then there are at most p principal components.

so $\mathbf{a}'\mathbf{x}$ is a scalar.

2.1 Solve PCA when Σ Diagonal

Remark 6. We can verify that

$$\text{Var}(\mathbf{a}'\mathbf{X}) = \mathbf{a}'\Sigma\mathbf{a}.$$

See the following example.

Example 2. Let $\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$ and $\mathbf{a} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$.

$$\begin{aligned} \text{Var}(\mathbf{a}'\mathbf{X}) &= \text{Var}\left(\begin{pmatrix} 1 & 2 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}\right) \\ &= \text{Var}(X_1 + 2X_2) \\ &= \text{Var}(X_1) + 4\text{Var}(X_2) \\ &= 1 + 4 \cdot 2 = 9 \end{aligned}$$

Alternatively, we can calculate $\text{Var}(\mathbf{a}'\mathbf{X})$ by

$$\mathbf{a}'\Sigma\mathbf{a} = \begin{pmatrix} 1 & 2 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} = 9.$$

Indeed, we end up with the same result, as the previous remarks implied.

By remark 6, the optimization problem can be written as

$$\max_{\mathbf{a}_1 \in \mathbb{R}^p} \mathbf{a}_1'\Sigma\mathbf{a}_1 \quad \text{s.t.} \quad \mathbf{a}_1'\mathbf{a}_1 = 1.$$

In his book “How to Solve It”, mathematician and Probabilist George Pólya famously said

“If you can’t solve a problem, then there is an easier problem you can solve: find it.”

Let’s follow his suggestion by assuming the covariance matrix

$$\Sigma = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix}_{p \times p},$$

where $\lambda_1 > \lambda_2 > \dots > \lambda_p$. That is, we assume Σ is now a diagonal matrix with positive, decreasing diagonal elements.

Now we are ready to solve the optimization problem for PC_1 :

$$\max_{\mathbf{a}_1 \in \mathbb{R}^p} \mathbf{a}_1' \Sigma \mathbf{a}_1 \quad \text{s.t.} \quad \mathbf{a}_1' \mathbf{a}_1 = 1.$$

Notice the objective function is now

$$\begin{aligned} \mathbf{a}_1' \Sigma \mathbf{a}_1 &= \begin{pmatrix} a_{11} & \dots & a_{1p} \end{pmatrix} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} a_{11} \\ \vdots \\ a_{1p} \end{pmatrix} \\ &= \lambda_1 a_{11}^2 + \lambda_2 a_{12}^2 + \dots + \lambda_p a_{1p}^2. \end{aligned}$$

By redefining $b_{1i} = a_{1i}^2$, the maximization problem becomes

$$\max_{\mathbf{b}} \lambda_1 b_{11} + \dots + \lambda_p b_{1p} \quad \text{s.t.} \quad b_{11} + \dots + b_{1p} = 1.$$

It is plain to see that $b_{11} = 1, b_{12} = \dots = b_{1p} = 0$ attains the maximum, while subjecting to the constraint. Hence, $a_{11} = \pm 1, a_{12} = \dots = a_{1p} = 0$ is the solution to the PC_1 problem.²

²PCs are only uniquely defined up to reflection over the origin. It is not hard to see that if \mathbf{a}^* is a solution to the PCA problem, then $-\mathbf{a}^*$ is also a solution.

The first PC is hence

$$PC_1 = \mathbf{a}'\mathbf{X} = \begin{pmatrix} 1 & 0 & \dots & 0 \end{pmatrix} \begin{pmatrix} X_1 \\ \vdots \\ X_p \end{pmatrix} = X_1,$$

which points to the direction with largest variance.

Next we solve the PC_2 problem

$$\begin{aligned} \max_{\mathbf{a}_2 \in \mathbb{R}^p} \quad & Var(\mathbf{a}_2'\mathbf{X}) \\ \text{s.t.} \quad & \mathbf{a}_2'\mathbf{a}_2 = 1 \\ & Cov(\mathbf{a}_1'\mathbf{X}, \mathbf{a}_2'\mathbf{X}) = 0 \end{aligned}$$

The constraint

$$\begin{aligned} Cov(\mathbf{a}_1'\mathbf{X}, \mathbf{a}_2'\mathbf{X}) &= Cov(X_1, \mathbf{a}_2'\mathbf{X}) \\ &= \mathbb{E}[X_1(a_{21}X_1 + \dots + a_{2p}X_p)] \\ &= a_{21}\mathbb{E}[X_1^2] = 0, \end{aligned}$$

which implies $a_{21} = 0$. The last equality exploits $\mathbb{E}[X_iX_j] = 0$ for all $i \neq j$, and $\mathbb{E}[X_i] = 0$ for all i .

So $\mathbf{a}_2 = \begin{pmatrix} 0 \\ a_{22} \\ \vdots \\ a_{2p} \end{pmatrix}$, and the PC_2 problem reduces to

$$\begin{aligned} \mathbf{a}_2' \Sigma \mathbf{a}_2 &= \begin{pmatrix} 0 & a_{22} & \dots & a_{2p} \end{pmatrix} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} 0 \\ a_{22} \\ \vdots \\ a_{2p} \end{pmatrix} \\ &= \begin{pmatrix} a_{22} & \dots & a_{2p} \end{pmatrix} \begin{pmatrix} \lambda_2 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} a_{22} \\ \vdots \\ a_{2p} \end{pmatrix}. \end{aligned}$$

The problem is now the same as finding PC_1 , with one less variable. Hence the solution to the second component is $a_{22} = \pm 1$, $a_{21} = a_{23} = \dots = a_{2p} = 0$, and $PC_2 = \mathbf{a}_2' \mathbf{X} = X_2$.

Similarly, we have $PC_3 = X_3, PC_4 = X_4, \dots, PC_p = X_p$.

Conclusion: When Σ is diagonal, PC_i is X_i , for $i = 1, \dots, p$. So when X_i 's are uncorrelated, PCA is basically finding X_i 's with largest variances.

2.2 Solve PCA when Σ is not Diagonal

What if Σ is not diagonal? Luckily, we have the following theorem:

Theorem 1 (Real Spectral Theorem). *If Σ is symmetric, then there exists a $p \times p$ matrix \mathbf{P} such that*

$$\Sigma = \mathbf{P} \mathbf{D} \mathbf{P}^{-1},$$

where \mathbf{D} is a diagonal matrix, and

$$\mathbf{P}' \mathbf{P} = \mathbf{I}_{p \times p},$$

i.e., $\mathbf{P}^{-1} = \mathbf{P}'$.

Example 3.

$$\begin{pmatrix} 34 & 12 \\ 12 & 41 \end{pmatrix} = \begin{pmatrix} \frac{3}{5} & \frac{-4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{pmatrix} \begin{pmatrix} 50 & 0 \\ 0 & 25 \end{pmatrix} \begin{pmatrix} \frac{3}{5} & \frac{4}{5} \\ \frac{-4}{5} & \frac{3}{5} \end{pmatrix}.$$

Since Σ is symmetric, it can now be written as $\Sigma = \mathbf{P}\mathbf{D}\mathbf{P}'$, where the notation is followed by the Real Spectral theorem. Let $\mathbf{b}_i = \mathbf{P}'\mathbf{a}_i$,

$$\mathbf{a}_i'\Sigma\mathbf{a}_i = \mathbf{a}_i'\mathbf{P}\mathbf{D}\mathbf{P}'\mathbf{a}_i = (\mathbf{P}'\mathbf{a}_i)'\mathbf{D}(\mathbf{P}'\mathbf{a}_i) = \mathbf{b}_i'\mathbf{D}\mathbf{b}_i.$$

$$\mathbf{a}_i'\mathbf{a}_i = \mathbf{a}_i'(\mathbf{P}\mathbf{P}^{-1})\mathbf{a}_i = \mathbf{a}_i'(\mathbf{P}\mathbf{P}')\mathbf{a}_i = (\mathbf{P}'\mathbf{a}_i)'(\mathbf{P}'\mathbf{a}_i) = \mathbf{b}_i'\mathbf{b}_i$$

Restate the maximization problem in terms of \mathbf{b}_i gives

$$\max_{\mathbf{b}_i \in \mathbb{R}^p} \mathbf{b}_i'\mathbf{D}\mathbf{b}_i \quad \text{s.t.} \quad \mathbf{b}_i'\mathbf{b}_i = 1,$$

and we know how to solve it since \mathbf{D} is diagonal.