

Multivariate Analysis

Principal Component Analysis

Shaobo Jin

Department of Mathematics

Intended Learning Outcome

Through this chapter, you should be able to

- derive PCA using matrix algebra
- conduct PCA

Motivation

Extract information from data, and achieve dimension reduction as an early step of an analytical process.

- A data set may contain a long list of variables.
- We want to reduce them to a smaller set of summary indices.
- Most of the information in the original set of variables are still preserved.

Task of Principal Component Analysis (PCA)

Rough Task

Let the random vector \mathbf{X} ($p \times 1$) have the covariance matrix $\mathbf{\Sigma} \geq 0$. Find linear combinations $Y_i = \mathbf{a}_i^T \mathbf{X}$ such that

\mathbf{a}_1 maximizes $\text{var}(\mathbf{a}_1^T \mathbf{X})$,

\mathbf{a}_2 maximizes $\text{var}(\mathbf{a}_2^T \mathbf{X})$, and $\text{cov}(\mathbf{a}_2^T \mathbf{X}, \mathbf{a}_1^T \mathbf{X}) = 0$,

\mathbf{a}_3 maximizes $\text{var}(\mathbf{a}_3^T \mathbf{X})$, and $\text{cov}(\mathbf{a}_3^T \mathbf{X}, \mathbf{a}_j^T \mathbf{X}) = 0$, $j < 3$,

\vdots

\mathbf{a}_p maximizes $\text{var}(\mathbf{a}_p^T \mathbf{X})$, and $\text{cov}(\mathbf{a}_p^T \mathbf{X}, \mathbf{a}_j^T \mathbf{X}) = 0$, $j < p$.

Restriction

Consider the linear combination $Y_i = \mathbf{a}_i^T \mathbf{X}$. We have

$$\begin{aligned}\text{var}(Y_i) &= \mathbf{a}_i^T \boldsymbol{\Sigma} \mathbf{a}_i, \\ \text{cov}(Y_i, Y_k) &= \mathbf{a}_i^T \boldsymbol{\Sigma} \mathbf{a}_k, \quad i \neq k.\end{aligned}$$

Consider the new linear combination $Z_i = cY_i = c\mathbf{a}_i^T \mathbf{X}$, for a constant c . We have

$$\begin{aligned}\text{var}(Z_i) &= c^2 \mathbf{a}_i^T \boldsymbol{\Sigma} \mathbf{a}_i, \\ \text{cov}(Z_i, Z_k) &= c^2 \mathbf{a}_i^T \boldsymbol{\Sigma} \mathbf{a}_k.\end{aligned}$$

Hence, we need to set the scale such as $\mathbf{a}_i^T \mathbf{a}_i = 1$.

Principal Components

Principal components are particular linear combinations of the p random variables X_1, X_2, \dots, X_p .

- 1 First principal component is the linear combination that maximizes $\text{var}(\mathbf{a}_1^T \mathbf{X})$ subject to $\mathbf{a}_1^T \mathbf{a}_1 = 1$.
- 2 Second principal component is the linear combination that maximizes $\text{var}(\mathbf{a}_2^T \mathbf{X})$ subject to $\mathbf{a}_2^T \mathbf{a}_2 = 1$ and $\text{cov}(\mathbf{a}_2^T \mathbf{X}, \mathbf{a}_1^T \mathbf{X}) = 0$
- 3 i th principal component is the linear combination that maximizes $\text{var}(\mathbf{a}_i^T \mathbf{X})$ subject to $\mathbf{a}_i^T \mathbf{a}_i = 1$ and $\text{cov}(\mathbf{a}_i^T \mathbf{X}, \mathbf{a}_k^T \mathbf{X}) = 0$ for all $k < i$.

It is not required to have a multivariate normal assumption for \mathbf{X} .

Two Useful Lemma

Lemma

Consider the function $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$, where \mathbf{x} is a vector. Then,

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}.$$

If \mathbf{A} is also symmetric, then,

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{A}\mathbf{x}.$$

Lemma

Let \mathbf{A} and $\mathbf{B} > 0$ be two symmetric matrices. The maximum value of $\mathbf{x}^T \mathbf{A} \mathbf{x}$ subject to $\mathbf{x}^T \mathbf{B} \mathbf{x} = 1$ is attained when \mathbf{x} is the eigenvector of $\mathbf{B}^{-1} \mathbf{A}$ corresponding to the largest eigenvalue of $\mathbf{B}^{-1} \mathbf{A}$. Its maximum value is the largest eigenvalue of $\mathbf{B}^{-1} \mathbf{A}$.

Find Principal Components

In order to find the first principal component, we consider

$$\max \text{var}(\mathbf{a}_1^T \mathbf{X}) \quad \text{s.t.} \quad \mathbf{a}_1^T \mathbf{a}_1 = 1.$$

In other words, we need to optimize

$$f(\mathbf{a}_1) = \mathbf{a}_1^T \Sigma \mathbf{a}_1 - \lambda (\mathbf{a}_1^T \mathbf{a}_1 - 1).$$

In order to find the first principal component, we consider

$$\max \text{var}(\mathbf{a}_2^T \mathbf{X}) \quad \text{s.t.} \quad \mathbf{a}_2^T \mathbf{a}_2 = 1, \text{cov}(\mathbf{a}_2^T \mathbf{X}, \mathbf{a}_1^T \mathbf{X}) = 0.$$

In other words, we need to optimize

$$f(\mathbf{a}_2) = \mathbf{a}_2^T \Sigma \mathbf{a}_2 - \lambda_1 (\mathbf{a}_2^T \mathbf{a}_2 - 1), \quad \mathbf{a}_2 \notin \text{span}\{\mathbf{a}_1\}.$$

Principal Components

Result 8.1: Simply An Eigen Decomposition

Let Σ be the covariance matrix associated with the $p \times 1$ random vector \mathbf{X} . Let Σ have the eigenvalue-eigenvector pairs $(\lambda_1, \mathbf{e}_1), \dots, (\lambda_p, \mathbf{e}_p)$, where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Then, the i th principal component is $Y_i = \mathbf{e}_i^T \mathbf{X}$. With these choices

$$\begin{aligned}\text{var}(Y_i) &= \mathbf{e}_i^T \Sigma \mathbf{e}_i = \lambda_i, \\ \text{cov}(Y_i, Y_k) &= 0.\end{aligned}$$

Total Variation Explained by Principal Components

Result 8.2

Let \mathbf{X} have covariance matrix $\mathbf{\Sigma}$ with eigenvalue-eigenvector pairs $(\lambda_1, \mathbf{e}_1), \dots, (\lambda_p, \mathbf{e}_p)$, where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Let $Y_i = \mathbf{e}_i^T \mathbf{X}$, $i = 1, \dots, p$, be the unique principal components. Then,

$$\sum_{i=1}^p \sigma_{ii} = \sum_{i=1}^p \lambda_i = \sum_{i=1}^p \text{var}(Y_i).$$

This result says that the **total population variance**, $\sum_{i=1}^p \text{var}(X_i)$, is the same as the **total principal component variance**. Hence, the proportion of total variance explained by the k th principal component is

$$\frac{\lambda_k}{\sum_{i=1}^p \lambda_i}.$$

Importance

Result 8.3

If $Y_i = \mathbf{e}_i^T \mathbf{X}$, $i = 1, 2, \dots, p$, are the principal components obtained from the covariance matrix $\mathbf{\Sigma}$, then

$$\rho_{Y_i, X_k} = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}, \quad i, k = 1, 2, \dots, p,$$

are the correlation coefficients between the components Y_i and the variables X_k .

The magnitude of e_{ik} measures the importance of X_k to the i th principal component Y_i .

Principal Components From Correlation Matrix

Suppose that we standardize all X_i by $Z_i = (X_i - \mu_i) / \sqrt{\sigma_{ii}}$. All our previous results apply to $\text{cov}(\mathbf{Z}) = \text{corr}(\mathbf{X})$.

Result 8.4

The i th principal component of the standardized variables \mathbf{Z} with $\text{cov}(\mathbf{Z}) = \boldsymbol{\rho}$ is given by

$$Y_i = \mathbf{e}_i^T \mathbf{Z} = \mathbf{e}_i^T \mathbf{V}^{-1/2} (\mathbf{X} - \boldsymbol{\mu}).$$

Moreover, $\sum_{i=1}^p \text{var}(Y_i) = \sum_{i=1}^p \text{var}(Z_i) = p$, and $\rho_{Y_i, Z_k} = e_{ik} \sqrt{\lambda_i}$. In this case, $(\lambda_1, \mathbf{e}_1), \dots, (\lambda_p, \mathbf{e}_p)$ are the eigenvalue-eigenvector pairs for $\boldsymbol{\rho}$, with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.

Standardization Matters

The PC derived from Σ are different from those derived from ρ . One set of PC is not a simple function of the other.

```
A <- matrix(c(2.0, 0.5, 0.4,
              0.5, 1.5, 0.3,
              0.4, 0.3, 1.0), 3, 3, byrow = TRUE)
eigen(A)
```

```
## eigen() decomposition
## $values
## [1] 2.477083 1.195800 0.827117
##
## $vectors
##      [,1]      [,2]      [,3]
## [1,] 0.8000667 0.5626808 -0.2080470
## [2,] 0.5075924 -0.8197795 -0.2651631
## [3,] 0.3197549 -0.1065451 0.9414908
```

```
D <- diag(1.0 / sqrt(c(2.0, 1.5, 1)))
D %*% A %*% D
```

```
##      [,1]      [,2]      [,3]
## [1,] 1.0000000 0.2886751 0.2828427
## [2,] 0.2886751 1.0000000 0.2449490
## [3,] 0.2828427 0.2449490 1.0000000
```

```
eigen(D %*% A %*% D)
```

```
## eigen() decomposition
## $values
## [1] 1.5447573 0.7552427 0.7000000
##
## $vectors
##      [,1]      [,2]      [,3]
## [1,] -0.5958111 -0.0463897 0.8017837
## [2,] -0.5700908 -0.6787568 -0.4629100
## [3,] -0.5656904 0.7328965 -0.3779645
```

Sample Principal Components

Suppose the data $\mathbf{x}_1, \dots, \mathbf{x}_n$ represent n independent drawings from some p -dimensional population with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, not necessarily a multivariate population.

- ① First **sample principal component** is the linear combination that maximizes the sample variance of $\mathbf{a}_1^T \mathbf{x}$ subject to $\mathbf{a}_1^T \mathbf{a}_1 = 1$.
- ② Second **sample principal component** is the linear combination that maximizes the sample variance of $\mathbf{a}_2^T \mathbf{x}$ subject to $\mathbf{a}_2^T \mathbf{a}_2 = 1$ and zero sample covariance between $\mathbf{a}_2^T \mathbf{x}$ and $\mathbf{a}_1^T \mathbf{x}$.
- ③ i th **sample principal component** is the linear combination that maximizes the sample variance of $\mathbf{a}_i^T \mathbf{x}$ subject to $\mathbf{a}_i^T \mathbf{a}_i = 1$ and zero sample covariance between $\mathbf{a}_i^T \mathbf{x}$ and $\mathbf{a}_k^T \mathbf{x}$ for all $k < i$.

Apply Sample Covariance

By [Result 2.5](#),

- 1 the linear combination $\mathbf{a}_i^T \mathbf{x}$ has sample mean $\mathbf{a}_i^T \bar{\mathbf{x}}$, and sample variance $\mathbf{a}_i^T \mathbf{S} \mathbf{a}_i$,
- 2 the sample covariance between $\mathbf{a}_i^T \mathbf{x}$ and $\mathbf{a}_k^T \mathbf{x}$ is $\mathbf{a}_i^T \mathbf{S} \mathbf{a}_k$.

Hence,

- 1 First [sample principal component](#) is the linear combination that maximizes $\mathbf{a}_1^T \mathbf{S} \mathbf{a}_1$ subject to $\mathbf{a}_1^T \mathbf{a}_1 = 1$.
- 2 Second [sample principal component](#) is the linear combination that maximizes $\mathbf{a}_2^T \mathbf{S} \mathbf{a}_2$ subject to $\mathbf{a}_2^T \mathbf{a}_2 = 1$ and zero sample covariance between $\mathbf{a}_2^T \mathbf{x}$ and $\mathbf{a}_1^T \mathbf{x}$.
- 3 i th [sample principal component](#) is the linear combination that maximizes $\mathbf{a}_i^T \mathbf{S} \mathbf{a}_i$ subject to $\mathbf{a}_i^T \mathbf{a}_i = 1$ and zero sample covariance between $\mathbf{a}_i^T \mathbf{x}$ and $\mathbf{a}_k^T \mathbf{x}$ for all $k < i$.

PCA

Result: Simply An Eigen Decomposition

If \mathbf{S} is a $p \times p$ sample covarinace matrix with eigenvalue-eigenvector pairs $(\hat{\lambda}_i, \hat{\mathbf{e}}_i)$, $i = 1, \dots, p$, the i th sample principal component is

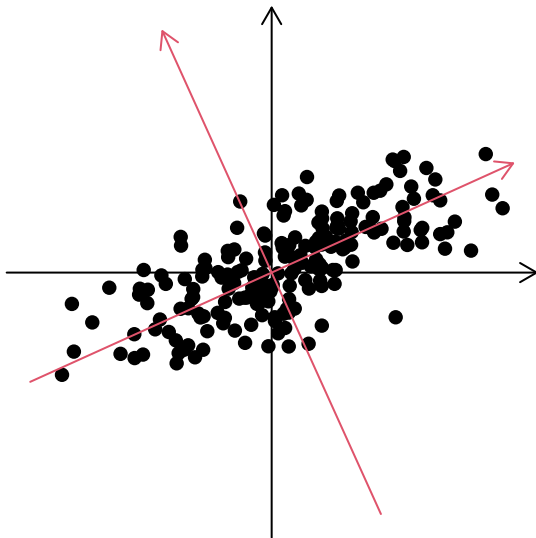
$$\hat{y}_i = \hat{\mathbf{e}}_i^T \mathbf{x},$$

where $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$. The values $(\hat{y}_1 \dots \hat{y}_p)$ are the **principal component scores**. The sample variance of \hat{y}_i is $\hat{\lambda}_i$, and the sample covariance between \hat{y}_i and \hat{y}_k is 0. In addition, the total sample variance satisfies $\sum_{i=1}^p s_{ii} = \sum_{i=1}^p \hat{\lambda}_i$, and

$$r_{\hat{y}_i, x_k} = \frac{\hat{e}_{ik} \sqrt{\hat{\lambda}_i}}{\sqrt{s_{kk}}}.$$

Further, the PCs from \mathbf{S} and the sample correlation matrix are not the same.

Rotation of Axes



Eigen Decomposition and Singular Value Decomposition

PCA is simply eigen decomposition of \mathbf{S} as $\mathbf{S} = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^T$.

- The eigen decomposition is only applicable to a symmetric matrix. Every matrix has a **singular value decomposition (SVD)**. For any $m \times n$ matrix \mathbf{A} of rank r , there exist an $m \times m$ orthogonal matrix \mathbf{U} and an $n \times n$ orthogonal matrix \mathbf{V} such that

$$\mathbf{A}_{m \times n} = \mathbf{U}_{m \times m} \left(\begin{array}{cc} \mathbf{D}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{array} \right)_{m \times n} \mathbf{V}_{n \times n}^T,$$

where \mathbf{D}_1 is an $r \times r$ diagonal matrix with diagonal elements that are positive.

- Diagonal elements in $\left(\begin{array}{cc} \mathbf{D}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{array} \right)$ are called singular values.

Find Singular Value Decomposition

- 1 The diagonal elements in \mathbf{D}_1 are the positive square root of the nonzero eigenvalues (not necessarily distinct) of $\mathbf{A}^T \mathbf{A}$ or $\mathbf{A} \mathbf{A}^T$.
- 2 Let the columns of \mathbf{V}_1 be the eigenvectors corresponding to the nonzero eigenvalues of $\mathbf{A}^T \mathbf{A}$, and the columns of \mathbf{V}_2 be the eigenvectors corresponding to the 0 eigenvalues. Then $\mathbf{V} = [\mathbf{V}_1 \quad \mathbf{V}_2]$.
- 3 Let the columns of \mathbf{U}_1 be the eigenvectors corresponding to the nonzero eigenvalues of $\mathbf{A} \mathbf{A}^T$, and the columns of \mathbf{U}_2 be the eigenvectors corresponding to the 0 eigenvalues. Then $\mathbf{U} = [\mathbf{U}_1 \quad \mathbf{U}_2]$.

The SVD is

$$\mathbf{A} = \mathbf{U} \begin{pmatrix} \mathbf{D}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{V}^T = \mathbf{U}_1 \mathbf{D}_1 \mathbf{V}_1^T.$$

And \mathbf{U}_1 actually must satisfy $\mathbf{U}_1 = \mathbf{A} \mathbf{V}_1 \mathbf{D}_1^{-1}$.

SVD in R

```

A <- matrix(c(2, 0, 1, 0, 1, 2), 3, 2)
edV <- eigen(t(A) %*% A)
D1 <- diag(sqrt(edV$values))
D <- rbind(D1, matrix(0, 1, 2))
V <- edV$vectors
edU <- eigen(A %*% t(A))
U <- edU$vectors
round(cbind(U %*% D %*% t(V), NA, U %*% D %*% t(-1.0 * V), NA,
            (A %*% V %*% solve(D1)) %*% D1 %*% t(V)), 6)

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
## [1,]   -2    0  NA    2    0  NA    2    0
## [2,]    0   -1  NA    0    1  NA    0    1
## [3,]   -1   -2  NA    1    2  NA    1    2

```

SVD in R

```
A <- matrix(c(2, 0, 1, 0, 1, 2), 3, 2)
svd(A)

## $d
## [1] 2.645751 1.732051
##
## $u
##           [,1]      [,2]
## [1,] -0.5345225  0.8164966
## [2,] -0.2672612 -0.4082483
## [3,] -0.8017837 -0.4082483
##
## $v
##           [,1]      [,2]
## [1,] -0.7071068  0.7071068
## [2,] -0.7071068 -0.7071068
```

PCA and SVD

Suppose that we have a data matrix \mathbf{X} that has been demeaned. \mathbf{X} can be SVD decomposed as

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T \Rightarrow \mathbf{X}\mathbf{V} = \mathbf{U}\mathbf{D},$$

where the m columns of \mathbf{U} are orthonormal eigenvectors of $\mathbf{X}\mathbf{X}^T$. Then, we should have

$$\mathbf{X}^T\mathbf{X}/(n-1) = \mathbf{V}\mathbf{D}\mathbf{U}^T\mathbf{U}\mathbf{D}\mathbf{V}^T/(n-1) = \mathbf{V}[\mathbf{D}^2/(n-1)]\mathbf{V}^T.$$

- the PCA loading matrix: \mathbf{V} .
- the principal component scores: $\mathbf{U}\mathbf{D}$.
- the variances of principal components: $\mathbf{D}^2/(n-1)$.

SVD and PCA

```
dx <- x - matrix(colMeans(x), nrow = N, ncol = 2, byrow = TRUE)
sqrt( (svd(dx)$d^2) / (N - 1) )
```

```
## [1] 0.7600400 0.3254688
```

```
svd(dx)$v
```

```
##           [,1]      [,2]
## [1,] -0.8377819 -0.5460050
## [2,] -0.5460050  0.8377819
```

```
prcomp(x)
```

```
## Standard deviations (1, .., p=2):
```

```
## [1] 0.7600400 0.3254688
```

```
##
```

```
## Rotation (n x k) = (2 x 2):
```

```
##           PC1      PC2
```

```
## [1,] -0.8377819 -0.5460050
```

```
## [2,] -0.5460050  0.8377819
```

How Many Components?

There is no definite answer on how many PCs we should choose. Some popular methods are

- Scree plot: the point (elbow) before where the curve flattens.
- Choose the number of PCs such that a specified percentage of variance been explained.