

Lecture 12

Principle Components Analysis

Reading: JWHT Chapter 10

DSA 8020 Statistical Methods II

March 29-April 2, 2021

Multivariate Analysis

Principal component analysis (PCA)

Principal Component Regression

Whitney Huang
Clemson University

Multivariate Analysis

Principal component
analysis (PCA)

Principal Component
Regression

1 Multivariate Analysis

2 Principal component analysis (PCA)

3 Principal Component Regression

An Overview of Multivariate Analysis

Principle
Components
Analysis



- In many studies, observations are collected on **several variables** on each experimental/observational unit
- Multivariate analysis is a collection of statistical methods for analyzing these multivariate data sets
- **Common Objectives**
 - Dimensionality reduction
 - Classification
 - Grouping (Clustering)

You will learn more on **Multivariate Analysis** in DSA 8070

Multivariate Analysis

Principal component analysis (PCA)

Principal Component Regression

Multivariate Data

We display a multivariate data that contains n units on p variables using a matrix

$$\mathbf{X} = \begin{pmatrix} X_{1,1} & X_{2,1} & \cdots & X_{p,1} \\ X_{1,2} & X_{2,2} & \cdots & X_{p,2} \\ \vdots & \cdots & \ddots & \vdots \\ X_{1,n} & X_{2,n} & \cdots & X_{p,n} \end{pmatrix}$$

Summary Statistics

- **Mean Vector:** $\bar{\mathbf{X}} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)^T$
- **Covariance Matrix:** $\Sigma = \{\sigma_{ij}\}_{i,j=1}^p$, where
 $\sigma_{ii} = \text{Var}(X_i)$, $i = 1, \dots, p$ and $\sigma_{ij} = \text{Cov}(X_i, X_j)$, $i \neq j$

Multivariate Analysis

Principal component analysis (PCA)

Principal Component Regression

Multivariate Data

We display a multivariate data that contains n units on p variables using a matrix

$$\mathbf{X} = \begin{pmatrix} X_{1,1} & X_{2,1} & \cdots & X_{p,1} \\ X_{1,2} & X_{2,2} & \cdots & X_{p,2} \\ \vdots & \cdots & \ddots & \vdots \\ X_{1,n} & X_{2,n} & \cdots & X_{p,n} \end{pmatrix}$$

Summary Statistics

- **Mean Vector:** $\bar{\mathbf{X}} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)^T$
- **Covariance Matrix:** $\Sigma = \{\sigma_{ij}\}_{i,j=1}^p$, where
 $\sigma_{ii} = \text{Var}(X_i)$, $i = 1, \dots, p$ and $\sigma_{ij} = \text{Cov}(X_i, X_j)$, $i \neq j$

Next, we are going to introduce **Principal Component Analysis (PCA)**, a useful tool for conducting dimension reduction

Multivariate Analysis

Principal component analysis (PCA)

Principal Component Regression

Example: Monthly Sea Surface Temperatures

Principle
Components
Analysis



Multivariate Analysis

Principal component
analysis (PCA)

Principal Component
Regression

Sea Surface Temperatures and Anomalies

Principle
Components
Analysis



- The “data” are gridded at a 2° by 2° resolution from $124^{\circ}E - 70^{\circ}W$ and $30^{\circ}S - 30^{\circ}N$. The dimension of this SST data set is 2303 (number of grid points in space) \times 552 (monthly time series from 1970 Jan. to 2015 Dec.)
- Sea-surface temperature anomalies are the temperature differences from the climatology (i.e. long-term monthly mean temperatures)
- We will demonstrate the use of Empirical Orthogonal Function (EOF) analysis to uncover the low-dimensional spatial structure of this spatio-temporal data set

Multivariate Analysis

Principal component analysis (PCA)

Principal Component Regression

The Empirical Orthogonal Function (EOF) Decomposition

Empirical orthogonal functions (EOFs) are the geophysicist's terminology for the eigenvectors in the eigen-decomposition of an empirical covariance matrix. In its discrete formulation, EOF analysis is simply Principal Component Analysis (PCA). EOFs are usually used

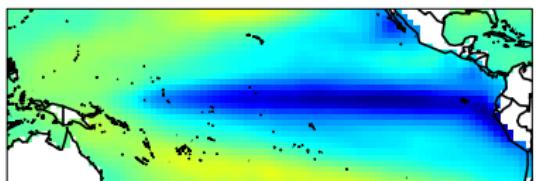
- To find principal spatial structures
- To reduce the dimension (spatially or temporally) in large spatio-temporal datasets

Multivariate Analysis

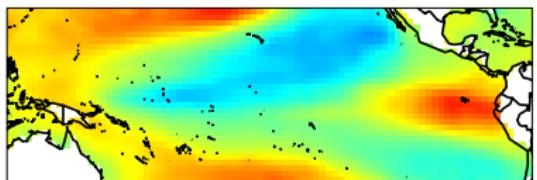
Principal component analysis (PCA)

Principal Component Regression

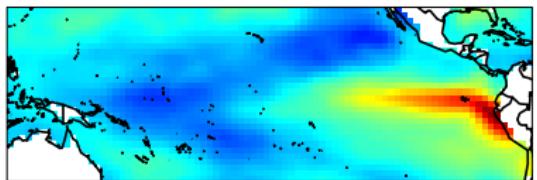
Perform EOF Decomposition and Plot the First Three Modes



EOF1: The classic
ENSO pattern



EOF2: A modulation
of the center



EOF3: Messing with
the coast of SA and
the Northern Pacific.

Multivariate Analysis

Principal component analysis (PCA)

Principal Component Regression

1998 Jan El Niño Event

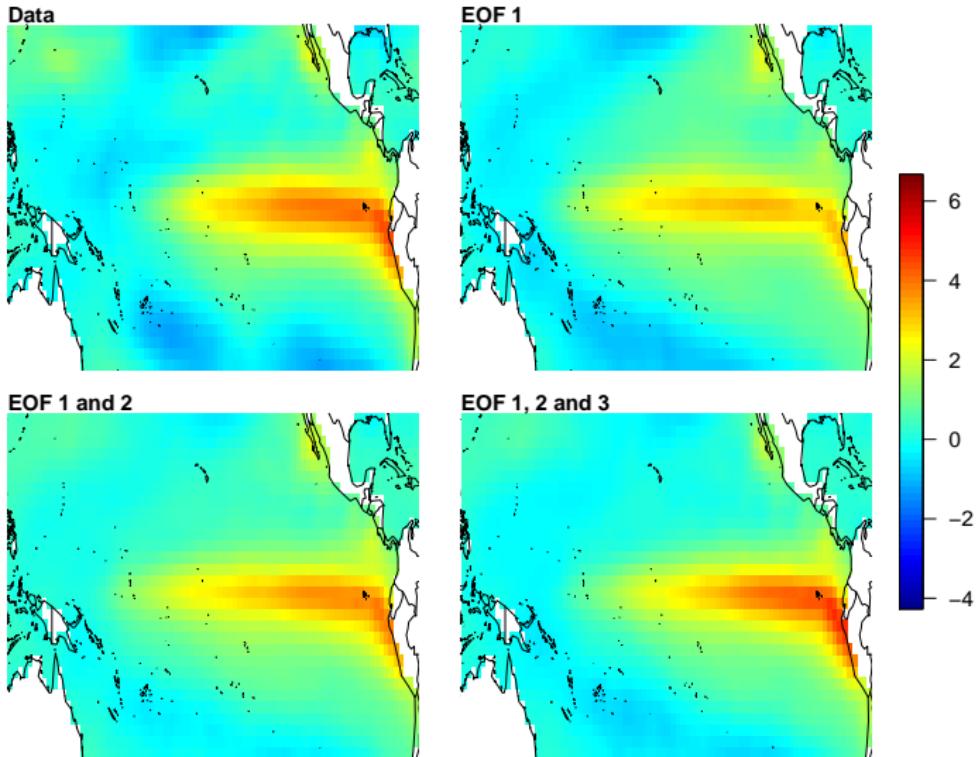
Principle
Components
Analysis



Multivariate Analysis

Principal component analysis (PCA)

Principal Component Regression



Principal Component Analysis

Given a random sample from a p -dimensional random vector

$$\mathbf{X}_i = \{X_{1,i}, X_{2,i}, \dots, X_{p,i}\}, \quad i = 1, \dots, n$$

- Dimension reduction technique
 - Large number of variables (p)
 - Number of variables (p) may be greater than number of observations (n)
- Create new, uncorrelated variables (principal components) for the follow up analysis
 - Principal Component Regression
 - Interpretation of principal components can be difficult in some situations

Multivariate Analysis

Principal component
analysis (PCA)Principal Component
Regression

Finding Principal Components

Principal Components (PCs) are uncorrelated **linear combinations** $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_p$ determined sequentially, as follows:

- ① The first PC is the linear combination

$$\tilde{X}_1 = \mathbf{c}_1^T \mathbf{X} = \sum_{i=1}^p c_{1i} X_i \text{ that maximize } \text{Var}(\tilde{X}_1) \text{ subject to}$$
$$\mathbf{c}_1^T \mathbf{c}_1 = 1$$

- ② The second PC is the linear combination

$$\tilde{X}_2 = \mathbf{c}_2^T \mathbf{X} = \sum_{i=1}^p c_{2i} X_i \text{ that maximize } \text{Var}(\tilde{X}_2) \text{ subject to}$$
$$\mathbf{c}_2^T \mathbf{c}_2 = 1 \text{ and } \mathbf{c}_2^T \mathbf{c}_1 = 0$$

⋮

- ③ The j_{th} PC is the linear combination

$$\tilde{X}_j = \mathbf{c}_j^T \mathbf{X} = \sum_{i=1}^p c_{ji} X_i \text{ that maximize } \text{Var}(\tilde{X}_j) \text{ subject to}$$
$$\mathbf{c}_j^T \mathbf{c}_j = 1 \text{ and } \mathbf{c}_j^T \mathbf{c}_k = 0, \forall k < j$$

Multivariate Analysis

Principal component
analysis (PCA)Principal Component
Regression

Finding Principal Components by Decomposing Covariance Matrix

- Let Σ , the covariance matrix of \mathbf{X} , have eigenvalue-eigenvector pairs $(\lambda_i, \mathbf{e}_i)_{i=1}^p$ with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Then, the k^{th} principal component is given by

$$\tilde{\mathbf{X}}_k = \mathbf{e}_k^T \mathbf{X} = e_{k1} X_1 + e_{k2} X_2 + \dots + e_{kp} X_p$$

- Then,

$$\text{Var}(\tilde{X}_i) = \lambda_i, \quad i = 1, \dots, p$$

$$\text{Cov}(\tilde{X}_j, \tilde{X}_k) = 0, \quad \forall j \neq k$$

PCA and Proportion of Variance Explained

Principle
Components
Analysis



- It can be shown that

$$\sum_{i=1}^p \text{Var}(\tilde{X}_i) = \lambda_1 + \lambda_2 + \cdots + \lambda_p = \sum_{i=1}^p \text{Var}(X_i)$$

- The proportion of the total variance associated with the k_{th} principal component is given by

$$\frac{\lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_p}$$

- If a large proportion of the total population variance (say 80% or 90%) is explained by the first k PCs, then we can restrict attention to the first k PCs without much loss of information

Multivariate Analysis

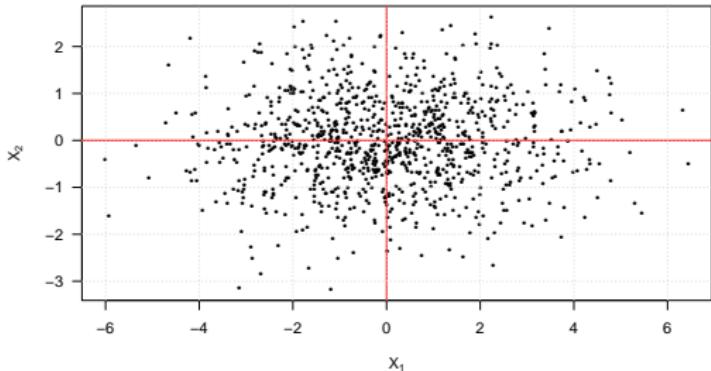
Principal component analysis (PCA)

Principal Component Regression

Toy Example 1

Suppose we have $\mathbf{X} = (X_1, X_2)^T$ where $X_1 \sim N(0, 4)$, $X_2 \sim N(0, 1)$ are independent

- Total variation = $\text{Var}(X_1) + \text{Var}(X_2) = 5$
- X_1 axis explains 80% of total variation
- X_2 axis explains the remaining 20% of total variation



Multivariate Analysis

Principal component
analysis (PCA)

Principal Component
Regression

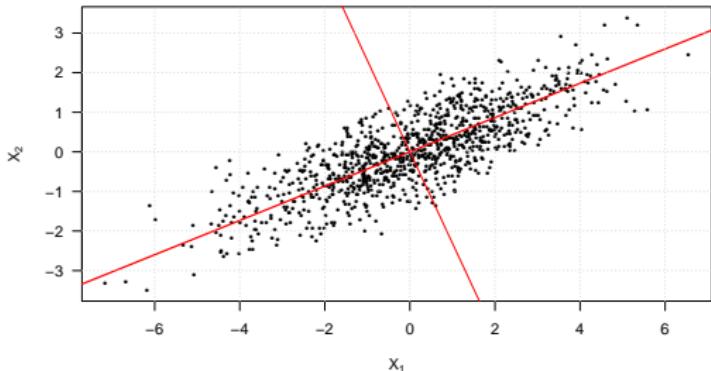
Toy Example 2

Suppose we have $\mathbf{X} = (X_1, X_2)^T$ where $X_1 \sim N(0, 4)$, $X_2 \sim N(0, 1)$ and $\text{Cor}(X_1, X_2) = 0.8$

- Total variation

$$= \text{Var}(X_1) + \text{Var}(X_2) = \text{Var}(\tilde{X}_1) + \text{Var}(\tilde{X}_2) = 5$$

- $\tilde{X}_1 = .9175X_1 + .3975X_2$ explains 93.9% of total variation
- $\tilde{X}_2 = .3975X_1 - .9176X_2$ explains the remaining 6.1% of total variation



Principal Component Regression

Principle
Components
Analysis



We are going to use Longley's data set, which provides a well-known example of multicollinearity, to illustrate Principal Component Regression

Correlation matrix

```
round(cor(longley[, -7]), 3)
```

```
##           GNP.deflator   GNP Unemployed Armed.Forces Population Year
## GNP.deflator      1.000 0.992     0.621     0.465    0.979 0.991
## GNP              0.992 1.000     0.604     0.446    0.991 0.995
## Unemployed       0.621 0.604     1.000    -0.177    0.687 0.668
## Armed.Forces     0.465 0.446    -0.177     1.000    0.364 0.417
## Population        0.979 0.991     0.687     0.364    1.000 0.994
## Year              0.991 0.995     0.668     0.417    0.994 1.000
```

Variance inflation factor

```
vif(longley[, -7])
```

```
## GNP.deflator          GNP Unemployed Armed.Forces Population Year
## 135.53244   1788.51348    33.61889    3.58893   399.15102  758.98060
```

Multivariate Analysis

Principal component analysis (PCA)

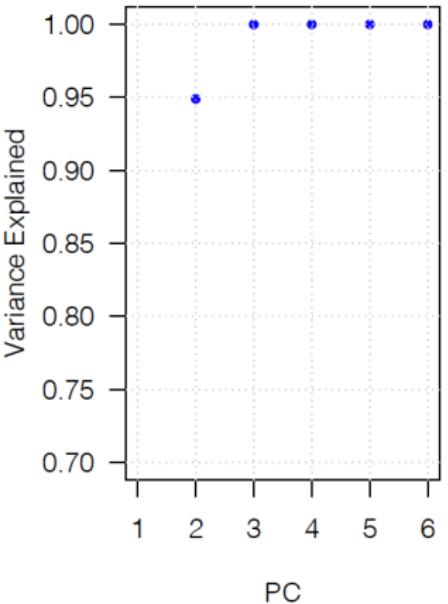
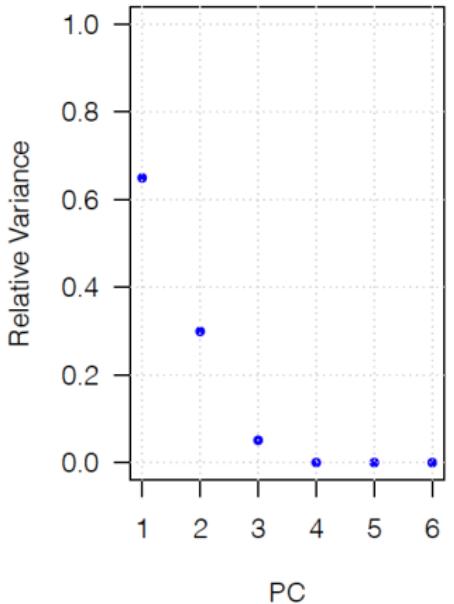
Principal Component Regression

How Many PCs to Use?

Multivariate Analysis

Principal component
analysis (PCA)

Principal Component
Regression



Principal Component Regression (PCR)

In PCR, instead of regressing the y on $x = (x_1, \dots, x_p)^T$ directly, \tilde{x} are used as regressors.

```
library(pls)
pcrFit <- pcr(Employed ~ ., data = longley, validation = "cv")
summary(pcrFit)
```

```
## Data: X dimension: 16 6
## Y dimension: 16 1
## Fit method: svdpc
## Number of components considered: 6
## TRAINING: % variance explained
##          1 comps 2 comps 3 comps 4 comps 5 comps 6 comps
## X         64.96   94.90   99.99  100.00  100.00  100.00
## Employed 78.42   89.73   98.51   98.56   98.83   99.55
```

Multivariate Analysis

Principal component analysis (PCA)

Principal Component Regression