

Lecture 9

Principle Component Analysis

Reading: Zelterman Chapter 8.1-8.4; Izenman Chapter 7.1-7.2

DSA 8070 Multivariate Analysis
October 17-October 21, 2022

Whitney Huang
Clemson University

Principle Component Analysis



Background
Finding Principal Components
Principal Components Analysis in Practice

9.1

Notes

Agenda

- 1 Background
- 2 Finding Principal Components
- 3 Principal Components Analysis in Practice

Principle Component Analysis



Background
Finding Principal Components
Principal Components Analysis in Practice

9.2

Notes

History

- **Karl Pearson (1901):** a procedure for finding lines and planes which best fit a set of points in p -dimensional space
- **Harold Hotelling (1933):** to find a smaller "fundamental set of independent variables" that determines the values of the original set of p variables

[11] On Lines and Planes of Closest Fit to Systems of Points in Space. By KARL PEARSON, F.R.S., University College, London.

(1) In many physical, statistical, and biological investigations it is desirable to represent a system of points in three, four, or higher dimensional space by the "locating" straight line or plane. Analytically this consists in taking

$$px + qy + rz, \text{ or } x = a_0 + a_1y + a_2z,$$

or

$$x = a_0 + a_1y + a_2z + a_3w + \dots + a_{p-1}z_{p-1},$$

where $x, y, z, w, \dots, z_{p-1}$ are variables, and determining the "best" values for the constants $a_0, a_1, a_2, \dots, a_{p-1}$ in relation to the observed corresponding values of the variables. In nearly all the cases dealt with in the treatise of best squares, the variables on the right of our equations are independent, since on the left is the dependent variable. The results of this treatment is that we get one straight line or plane if we treat some one variable as independent, and a quite different one if we treat another variable as the independent variable. There is no general

ANALYSIS OF A COMPLEX OF STATISTICAL VARIABLES INTO PRINCIPAL COMPONENTS*

HAROLD HOTELLING

Columbia University

I. INTRODUCTION

Consider a variable attending to each individual of a population. These statistical variables x_1, x_2, \dots, x_p might for example be scores made by school children in tests of speed and skill in solving arithmetical problems or in reading or they might be various physical properties of telephone poles, or the rates of exchange among various currencies. The x 's will ordinarily be correlated. It is natural to ask whether there are more fundamental sets of independent variables which, perhaps there is another than the x 's, which determine the values the x 's will take. If so, x_1, x_2, \dots, x_p are not variables, we shall then have a set of relations of the form

$$x_i = f_i(y_1, y_2, \dots, y_k) \quad (i = 1, 2, \dots, p) \quad (1)$$

Quantities such as the y 's have been called latent factors in recent psychological literature. However in view of the prospect of application of these ideas outside of psychology, and the misleading usage attaching to the word "factor" in mathematics, it will be better simply to call the y 's components of the complex depicted by the tests.

Principle Component Analysis



Background
Finding Principal Components
Principal Components Analysis in Practice

9.3

Notes

Basic Idea

Reduce the **dimensionality** of a data set in which there is a large number (i.e., p is "large") of inter-related variables while retaining as much as possible the **variation** in the original set of variables

- The reduction is achieved by transforming the original variables to a new set of variables, "**principal components**", that are **uncorrelated**
- These principal components are **ordered** such that **the first few retains most of the variation present in the data**
- **Goals/Objectives**
 - Reduction and summary
 - Study the structure of **covariance/correlation matrix**



Notes

Some Applications

- Interpretation (by studying the structure of covariance/correlation matrix)
- Select a sub-set of the original variables, that are uncorrelated to each other, to be used in other multivariate procedures (e.g., multiple regression, classification)
- Detect outliers or clusters of multivariate observations



Notes

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:** $\bar{\mathbf{X}} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)^T$, where $\bar{X}_j = \frac{\sum_{i=1}^n X_{ji}}{n}$, $j = 1, \dots, p$
- **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 discuss how to find **principal components**



Notes

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 $\tilde{X}_1 = \mathbf{c}_1^T \mathbf{X} = \sum_{i=1}^p c_{1i} X_i$ that maximize $\text{Var}(\tilde{X}_1)$ subject to $\mathbf{c}_1^T \mathbf{c}_1 = 1$
- 2 The second PC is the linear combination $\tilde{X}_2 = \mathbf{c}_2^T \mathbf{X} = \sum_{i=1}^p c_{2i} X_i$ that maximize $\text{Var}(\tilde{X}_2)$ subject to $\mathbf{c}_2^T \mathbf{c}_2 = 1$ and $\mathbf{c}_2^T \mathbf{c}_1 = 0$
- ⋮
- p The p th PC is the linear combination $\tilde{X}_p = \mathbf{c}_p^T \mathbf{X} = \sum_{i=1}^p c_{pi} X_i$ that maximize $\text{Var}(\tilde{X}_p)$ subject to $\mathbf{c}_p^T \mathbf{c}_p = 1$ and $\mathbf{c}_p^T \mathbf{c}_k = 0, \forall k < p$

Principle Component Analysis
 CLEMSON UNIVERSITY
 Background
 Finding Principal Components
 Principal Components Analysis in Practice
 9.7

Notes

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

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

⇒ we can perform a single matrix operation to get the coefficients to form all the PCs!

- Then,

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

Moreover $\text{Var}(\tilde{X}_1) \geq \text{Var}(\tilde{X}_2) \geq \dots \geq \text{Var}(\tilde{X}_p) \geq 0$

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

⇒ different PCs are **uncorrelated** with each other

Principle Component Analysis
 CLEMSON UNIVERSITY
 Background
 Finding Principal Components
 Principal Components Analysis in Practice
 9.8

Notes

PCA and Proportion of Variance Explained

- It can be shown that

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

- The proportion of the total variance associated with the k th principal component is given by

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

- 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 ⇒ **we achieve dimension reduction by considering $k < p$ uncorrelated components rather than the original p correlated variables**

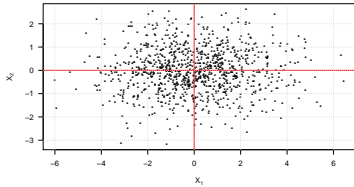
Principle Component Analysis
 CLEMSON UNIVERSITY
 Background
 Finding Principal Components
 Principal Components Analysis in Practice
 9.9

Notes

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



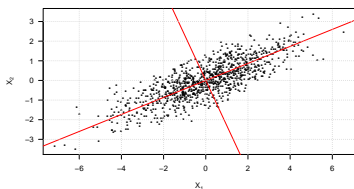
Principle Component Analysis
 CLEMSON UNIVERSITY
 Background
 Finding Principal Components
 Principal Components Analysis in Practice
 9.10

Notes

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



Principle Component Analysis
 CLEMSON UNIVERSITY
 Background
 Finding Principal Components
 Principal Components Analysis in Practice
 9.11

Notes

PCs of Standardized versus Original Variables

If we use standardized variables, i.e.,

$Z_j = \frac{X_j - \mu_j}{\sqrt{\sigma_{jj}}}$ $j = 1, \dots, p$ ("z-scores"). Then we are going to work with the **correlation matrix** instead of the **covariance matrix** of $(X_1, \dots, X_p)^T$

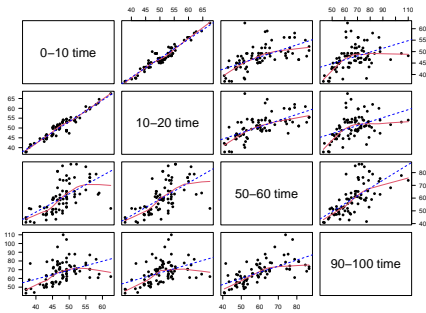
- We can obtain PCs of standardized variables by applying spectral decomposition of the correlation matrix
- However, the PCs (and the proportion of variance explained) are, in general, different than those from original variables
- If units of p variables are comparable, covariance PCA may be more informative, if units of p variables are incomparable, correlation PCA may be more appropriate

Principle Component Analysis
 CLEMSON UNIVERSITY
 Background
 Finding Principal Components
 Principal Components Analysis in Practice
 9.12

Notes

Example: Men's 100k Road Race

The data consists of the times (in minutes) to complete successive 10k segments ($p = 10$) of the race. There are 80 racers in total ($n = 80$)



Principle Component Analysis

CLEMSON UNIVERSITY

Background

Finding Principal Components

Principal Components Analysis in Practice

9.13

Notes

Eigenvalues of Σ

	Eigenvalue	Proportion	Cumulative
PC1	735.77	0.75	0.75
PC2	98.47	0.10	0.85
PC3	53.27	0.05	0.90
PC4	37.30	0.04	0.94
PC5	26.04	0.03	0.97
PC6	17.25	0.02	0.98
PC7	8.03	0.01	0.99
PC8	4.25	0.00	1.00
PC9	2.40	0.00	1.00
PC10	1.29	0.00	1.00

Much of the total variance can be explained by the first three PCs

Principle Component Analysis

CLEMSON UNIVERSITY

Background

Finding Principal Components

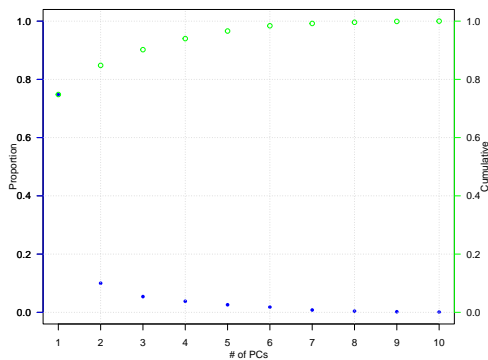
Principal Components Analysis in Practice

9.14

Notes

How Many Components to Retain?

A **scree plot** displays the variance explained by each component



Principle Component Analysis

CLEMSON UNIVERSITY

Background

Finding Principal Components

Principal Components Analysis in Practice

9.15

Notes

Men's 100k Road Race Component Weights

	Comp.1	Comp.2	Comp.3
0-10 time	0.13	0.21	0.36
10-20 time	0.15	0.25	0.42
20-30 time	0.20	0.31	0.34
30-40 time	0.24	0.33	0.20
40-50 time	0.31	0.30	-0.13
50-60 time	0.42	0.21	-0.22
60-70 time	0.34	-0.05	-0.19
70-80 time	0.41	-0.01	-0.54
80-90 time	0.40	-0.27	0.15
90-100 time	0.39	-0.69	0.35

What these numbers mean?

Principle Component Analysis

CLEMSON UNIVERSITY

Background

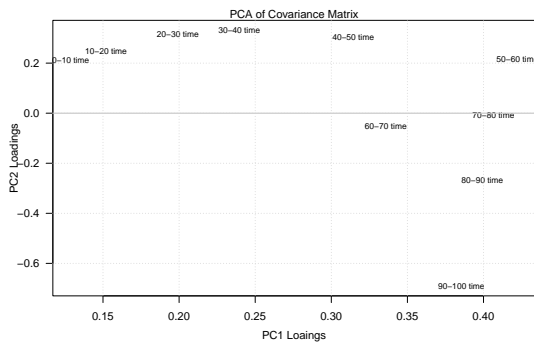
Finding Principal Components

Principal Components Analysis in Practice

9.16

Notes

Visualizing the Weights to Gain Insight



First component: overall speed
 Second component: contrast long and short races

Principle Component Analysis

CLEMSON UNIVERSITY

Background

Finding Principal Components

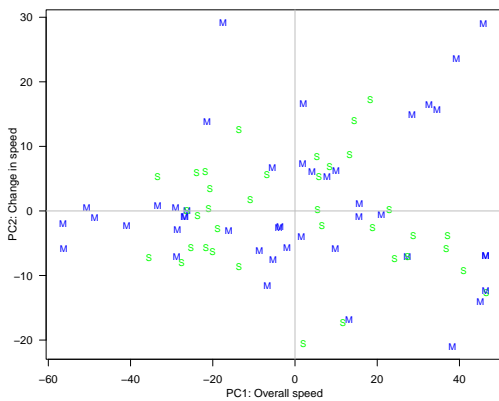
Principal Components Analysis in Practice

9.17

Notes

Looking for Patterns

Mature runners: Age < 40 (M); Senior runners: Age >= 40 (S)



Principle Component Analysis

CLEMSON UNIVERSITY

Background

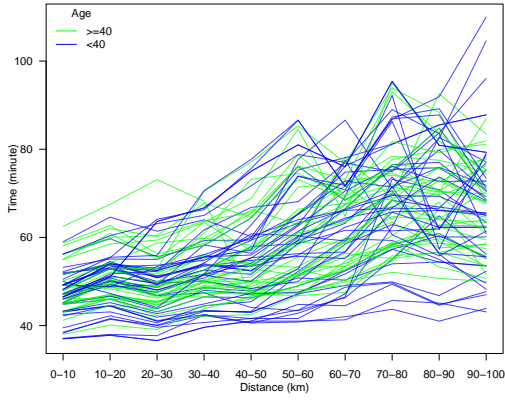
Finding Principal Components

Principal Components Analysis in Practice

9.18

Notes

Relating to Original Data: Profile Plot



Principle Component Analysis

CLEMSON UNIVERSITY

Background

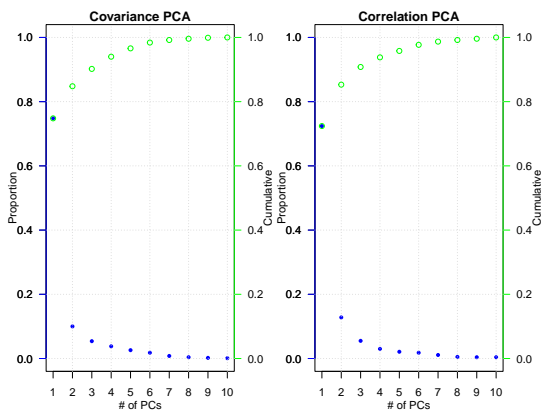
Finding Principal Components

Principal Components Analysis in Practice

9.19

Notes

Correlation PCA versus Covariance PCA



Principle Component Analysis

CLEMSON UNIVERSITY

Background

Finding Principal Components

Principal Components Analysis in Practice

9.20

Notes

Example: Monthly Sea Surface Temperatures

Principle Component Analysis

CLEMSON UNIVERSITY

Background

Finding Principal Components

Principal Components Analysis in Practice

9.21

Notes

Sea Surface Temperatures and Anomalies

- 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 structure of this spatio-temporal data set

Principle Component Analysis

CLEMSON UNIVERSITY

Background

Finding Principal Components

Principal Components Analysis in Practice

9.22

Notes

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

Principle Component Analysis

CLEMSON UNIVERSITY

Background

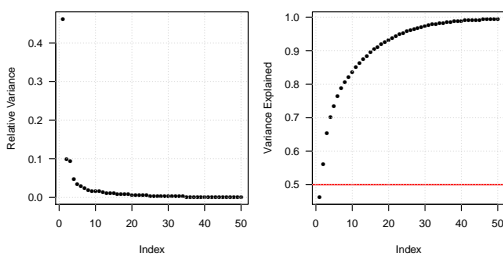
Finding Principal Components

Principal Components Analysis in Practice

9.23

Notes

Screen Plot for EOFs



Principle Component Analysis

CLEMSON UNIVERSITY

Background

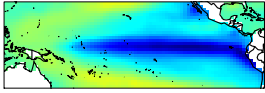
Finding Principal Components

Principal Components Analysis in Practice

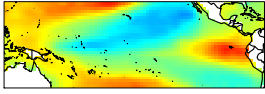
9.24

Notes

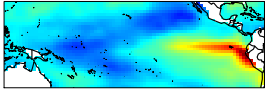
Perform EOF Decomposition and Plot the First Three Modes



EOF1: The classic ENSO pattern



EOF2: A modulation of the center



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

Principle Component Analysis

CLEMSON UNIVERSITY

Background

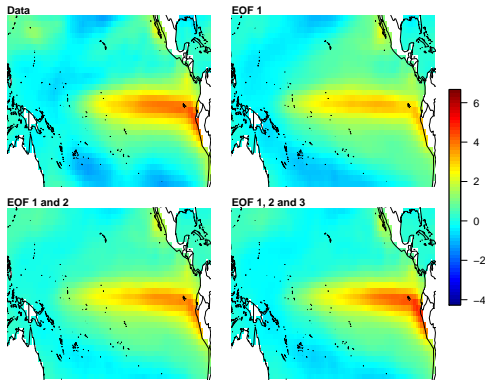
Finding Principal Components

Principal Components Analysis in Practice

9.25

Notes

1998 Jan El Niño Event



Principle Component Analysis

CLEMSON UNIVERSITY

Background

Finding Principal Components

Principal Components Analysis in Practice

9.26

Notes

Notes
