

Lecture 24

Computer Experiments & Principal Component Analysis

STAT 8020 Statistical Methods II

November 19, 2020

Whitney Huang
Clemson University

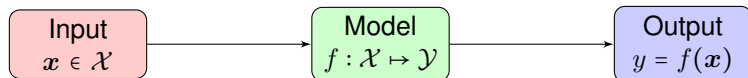
Agenda

- 1 **Computer Experiments**
- 2 **Multivariate Analysis**
- 3 **Principal component analysis (PCA)**

What is a Computer Experiment

In some situations it is economically, ethically, or simply not possible to run a **physical experiment**. Instead, the following scenario might be feasible:

- the physical process can be described by a mathematical model (e.g., a system of differential equations)
- computer code (simulator) can be written to compute the response from the mathematical model



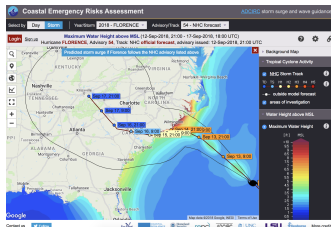
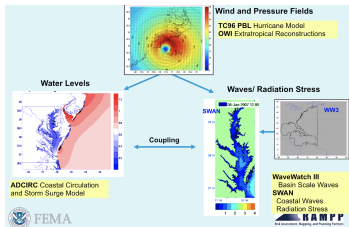
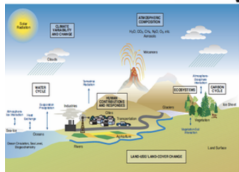
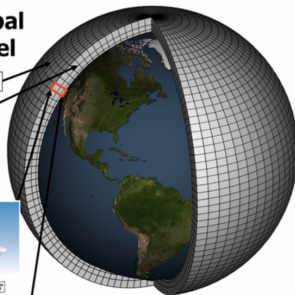
In this case, a researcher can conduct a **computer experiment** by running the computer code, which serves as a proxy for the physical process, to compute a “response” at any combination of values of the inputs

Examples of Computer Models

Schematic for Global Atmospheric Model

Horizontal Grid (Latitude-Longitude)

Vertical Grid (Height or Pressure)



Computer Experiments vs. Physical Experiments

- *“Experimental results are believed by everyone, except for the person who ran the experiment”*
- *“Computational results are believed by no one, except the person who wrote the code”*

Computer Experiments vs. Physical Experiments

- *“Experimental results are believed by everyone, except for the person who ran the experiment”*
- *“Computational results are believed by no one, except the person who wrote the code”*

Replication, randomization and blocking are irrelevant for a computer experiment because many **computer codes are deterministic** and **all the inputs to the code are known and can be controlled**

- **Design:**

where to make the runs, i.e., the selection of inputs $\{\mathbf{x}_i\}_{i=1}^n$
where $\mathbf{x}_i = (x_{1,i}, x_{2,i}, \dots, x_{d,i})$

- **Analysis:**

fit a statistical model using the model inputs-output
 $\{y_i, \mathbf{x}_i\}_{i=1}^n$ to “emulate” the simulator and to quantify the
prediction uncertainty for $y(\mathbf{x}_{\text{new}})$, usually via a **Gaussian
Process Model** $\text{GP}(m(\cdot), K(\cdot, \cdot))$, where

- $m(\mathbf{x}) = \text{E}[y(\mathbf{x})]$ is the **mean function**
- $K(\mathbf{x}, \mathbf{x}') = \text{Cov}(y(\mathbf{x}), y(\mathbf{x}'))$ is the **covariance function**

- 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)

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 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**

Example: Monthly Sea Surface Temperatures

Computer
Experiments &
Principal Component
Analysis

CLEMSON
UNIVERSITY

Computer Experiments

Multivariate Analysis

Principal component
analysis (PCA)

Sea Surface Temperatures and Anomalies

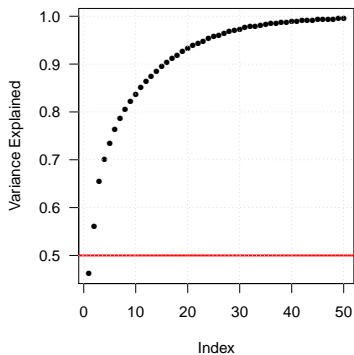
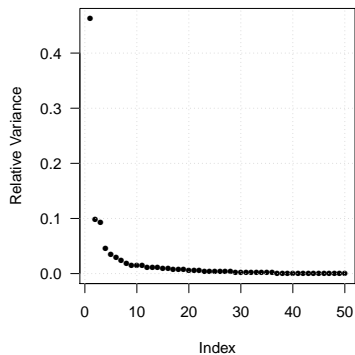
- The “data” are gridded at a 2° by 2° resolution from $124^\circ E - 70^\circ W$ and $30^\circ S - 30^\circ N$. The dimension of this SST data set is 2303 (number of grid points in space) \times 552 (monthly time series from 1970 Jan. to 2015 Dec.)
- Sea-surface temperature anomalies are the temperature 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

The Empirical Orthogonal Function (EOF) Decomposition

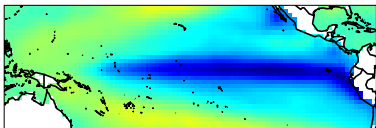
Empirical orthogonal functions (EOFs) are the geophysicist's terminology for the eigenvectors in the eigen-decomposition of an empirical covariance matrix. In its discrete formulation, EOF analysis is simply **Principal Component Analysis (PCA)**. EOFs are usually used

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

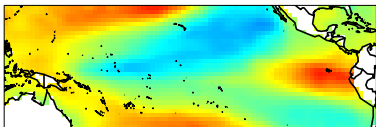
Screen Plot for EOFs



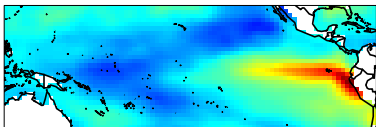
Perform EOF Decomposition and Plot the First Three Modes



EOF1: The classic ENSO pattern



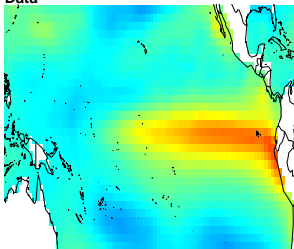
EOF2: A modulation of the center



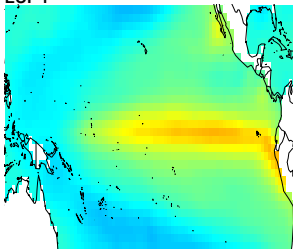
EOF3: Messing with the coast of SA and the Northern Pacific.

1998 Jan El Niño Event

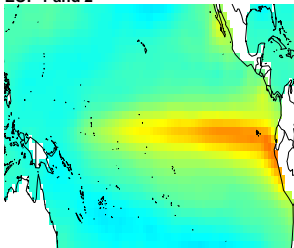
Data



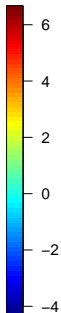
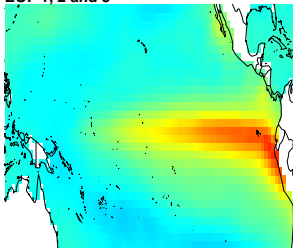
EOF 1



EOF 1 and 2



EOF 1, 2 and 3



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

Finding Principal Components

Principal Components (PC) 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 $\mathbf{c}_2^T \mathbf{c}_1 = 0$

⋮

3 The j_{th} PC is the linear combination $\tilde{X}_j = \mathbf{c}_j^T \mathbf{X} = \sum_{i=1}^p c_{ji} X_i$ that maximize $\text{Var}(\tilde{X}_j)$ subject to $\mathbf{c}_j^T \mathbf{c}_j = 1$ and $\mathbf{c}_j^T \mathbf{c}_k = 0 \forall k < j$

Principal Components

- 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

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

- 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

- 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

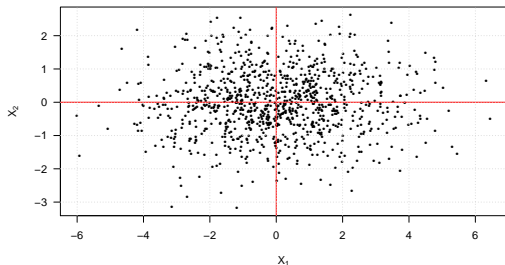
$$\frac{\lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_p}$$

- 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

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



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

