2

4

Multivariate Statistical Analysis

Shyh-Kang Jeng Department of Electrical Engineering/ Graduate Institute of Communication/ Graduate Institute of Networking and Multimedia



- Introduction
- Data Displays and Pictorial Representations
- Distances
- Reading Assignments

Outline

3

- Introduction
- → Organization of Data
- Data Displays and Pictorial Representations
- Distances
- Reading Assignments

Questions

- What is a model?
- →How to model Nature?







Why to Learn Multivariate Analysis?

- Explanation of a social or physical phenomenon must be tested by gathering and analyzing data
- Complexities of most phenomena require an investigator to collect observations on many different variables

10

12

Major Uses of Multivariate Analysis

- Data reduction or structural simplification
- Sorting and grouping
- Investigation of the dependence among variables
- Prediction
- Hypothesis construction and testing

Course Outline

- Introduction
- Matrix Algebra and Random Vectors
- Sample Geometry and Random Samples
- Multivariate Normal Distribution
- Inference about a Mean Vector
- Comparison of Several Multivariate Means

11

 Multivariate Linear Regression Models

Application Examples

- +Is one product better than the other?
- Which factor is the most important to determine the performance of a system?
- How to classify the results into clusters?
- What are the relationships between variables?

Course Outline

- Principal Components
- Factor Analysis and Inference for Structured Covariance Matrices
- Canonical Correlation Analysis*
- Discrimination and Classification*
- Clustering, Distance Methods, and Ordination*

14

16

Major Multivariate Techniques Not Included

- Structural Equation Models
- Multidimensional Scaling

Text Book and Website

13

15

- ★R. A. Johnson and D. W. Wichern, Applied Multivariate Statistical Analysis, 6th ed., Pearson Education, 2007. (雙葉)
- http://cc.ee.ntu.edu.tw/~skjeng/ MultivariateAnalysis2010.htm

Feature of This Course

 Uses matrix algebra to introduce theories and practices of multivariate statistical analysis

References

- ★林震岩,多變量分析-SPSS的操作與應用, 智勝, 2007
- J. F. Hair, Jr., B. Black, B. Babin, R. E. Anderson, and R. L. Tatham, Multivariate Data Analysis, 6th ed., Prentice Hall, 2006. (華泰)
- ✤D. C. Montgomery, Design and Analysis of Experiments, 6th ed., John Wiley, 2005. (歐亞)



Some Important Laws

- ✤First things first
- +80 20 Law
- ✤Fast prototyping and evolution
- ◆物有本末,事有始终,知所先後,則近道 矣。





Questions

- How to represent the measurement data for multivariate analysis?
- How to summarize the measurement data?
- How to determine if two variables are related?

Descriptive Statistics

21

- Summary numbers to assess the information contained in data
- Basic descriptive statistics
 - -Sample mean
 - Sample variance
 - Sample standard deviation
 - Sample covariance
 - -Sample correlation coefficient

Array of Data $\mathbf{x} = \begin{vmatrix} x_{21} & x_{22} & \cdots & x_{2k} & \cdots & x_{2p} \\ \vdots & \vdots & & \vdots & & \vdots \\ x_{j1} & x_{j2} & \cdots & x_{jk} & \cdots & x_{jp} \\ \vdots & \vdots & & \vdots & & \vdots \end{vmatrix}$ x_{n1} x_{n2} \cdots x_{nk} \cdots x_{np} 22



$\begin{aligned} & \text{Sample Covariance and} \\ & \text{sample Correlation Coefficient} \\ & s_{ik} = \frac{1}{n} \sum_{j=1}^{n} \left(x_{ji} - \overline{x}_i \right) \left(x_{jk} - \overline{x}_k \right) \\ & r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}} = \frac{\sum_{j=1}^{n} \left(x_{ji} - \overline{x}_i \right) \left(x_{jk} - \overline{x}_k \right)}{\sqrt{\sum_{j=1}^{n} \left(x_{ji} - \overline{x}_i \right)^2} \sqrt{\sum_{j=1}^{n} \left(x_{jk} - \overline{x}_k \right)^2} \\ & i = 1, 2, \cdots, p; \quad k = 1, 2, \cdots, p \\ & s_{ik} = s_{ki}, \quad r_{ik} = r_{ki} \end{aligned}$

Properties of Sample Correlation Coefficient

- Value is between -1 and 1
- Magnitude measure the strength of the linear association
- Sign indicates the direction of the association
- ✓ Value remains unchanged if all x_{ji}'s and x_{jk}'s are changed to y_{ji} = a x_{ji} + b and y_{jk} = c x_{jk} + d, respectively, provided that the constants a and c have the same sign

Standardized Values
(or Standardized Scores)
• Centered at zero
• Unit standard deviation
• Sample correlation coefficient can be
regarded as a sample covariance of
two standardized variables

$$\frac{x_{jk} - \overline{x}_k}{\sqrt{s_{kk}}}$$





















Lizard	Mass	SVL	HLS	Lizard	Mass	SVL	HLS
1	5.526	59.0	113.5	14	10.067	73.0	136.5
2	10.401	75.0	142.0	15	10.091	73.0	135.5
3	9.213	69.0	124.0	16	10.888	77.0	139.0
4	8.953	67.5	125.0	17	7.610	61.5	118.0
5	7.063	62.0	129.5	18	7.733	66.5	133.5
6	6.610	62.0	123.0	19	12.015	79.5	150.0
7	11.273	74.0	140.0	20	10.049	74.0	137.0
8	2.447	47.0	97.0	21	5.149	59.5	116.0
9	15.493	86.5	162.0	22	9.158	68.0	123.0
10	9.004	69.0	126.5	23	12.132	75.0	141.0
11	8.199	70.5	136.0	24	6.978	66.5	117.0
12	6.601	64.5	116.0	25	6.890	63.0	117.0
13	7.622	67.5	135.0				













Euclidean Distance

Each coordinate contributes equally to the distance

$$P(x_1, x_2, \dots, x_p), \quad Q(y_1, y_2, \dots, y_p)$$
$$d(P, Q) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_p - y_p)^2}$$

45

Statistical Distance for Uncorrelated Data $P(x_1, x_2), \quad O(0,0)$ $x_1^* = x_1 / \sqrt{s_{11}}, \quad x_2^* = x_2 / \sqrt{s_{22}}$ $d(O, P) = \sqrt{(x_1^*)^2 + (x_2^*)^2} = \sqrt{\frac{x_1^2}{s_{11}} + \frac{x_2^2}{s_{22}}}$













Necessary Conditions for Statistical Distance Definitions

d(P,Q) = d(Q,P) $d(P,Q) > 0 \text{ if } P \neq Q$ d(P,Q) = 0 if P = Q $d(P,Q) \leq d(P,R) + d(R,Q)$ (Triangle inequality)



Reading Assignments

→ Text book

-pp. 49-59 (Sections 2.1~2.2)

-pp. 82-96 (Supplement 2A)

55