

Multivariate statistics

Data

NULL SPACE

0.1 Organizing data

In multivariate analysis, we tend to analyze measurements made on several variables or features. We typically first select some p variables to record or measure and then the values of these variables are recorded for each distinct observation or item. We denote x_{jk} as the value taken by the k^{th} variable observed for the j^{th} item. So, if we have n observations or items and for each of these we have p measurements, we can easily package this set of observations into a $n \times p$ matrix as follows. Traditionally this is called a **data matrix**.

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \cdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \quad (1)$$

0.1.1 Descriptive statistics

In order to summarize information contained in vast amounts of data, we use certain measures known as **descriptive statistics**. For instance, **sample mean** gives us the central tendency in the data and **sample variance** gives us a measure of spread of the data.

- **Sample mean** for the k^{th} variable is given by:

$$\bar{x}_k = \frac{1}{n} \sum_{i=1}^n x_{ik} \quad (2)$$

- **Sample variance** for the k^{th} variable is given by:

$$s_k^2 = s_{kk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ik} - \bar{x}_k)^2 \quad (3)$$

- Suppose now that we have a pair of measurements concerning two variables - where vector indicates one measurement of the two variables:

$$\begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix}, \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix}, \cdots, \begin{bmatrix} x_{n1} \\ x_{n2} \end{bmatrix} \quad (4)$$

Generally we say that x_{j1} and x_{j2} are the variable 1 and variable 2 measurements for the j^{th} observation. Then a measure of linear association between the two variables is captured by the **sample covariance**:

$$s_{12} = \frac{1}{n-1} \sum_{j=1}^n (x_{j1} - \bar{x}_1)(x_{j2} - \bar{x}_2) \quad (5)$$

- The **sample Correlation coefficient** between two general variables x_i and x_k can be given by:

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}}\sqrt{s_{kk}}} = \frac{\sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)}{\sqrt{\sum_{j=1}^n (x_{ji} - \bar{x}_i)^2} \sqrt{\sum_{j=1}^n (x_{jk} - \bar{x}_k)^2}} \quad (6)$$

- When we have lots of variables and we want a convenient method to package all the **pairwise covariances and correlations** between all the p variables, then we use the **sample covariance and correlation matrices** as follows in the respective equations:

$$\mathbf{S}_n = \begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1p} \\ s_{21} & s_{22} & \cdots & s_{2p} \\ \vdots & \vdots & \cdots & \vdots \\ s_{p1} & s_{p2} & \cdots & s_{pp} \end{bmatrix} \quad (7)$$

$$\mathbf{R} = \begin{bmatrix} 1 & r_{12} & \cdots & r_{1p} \\ r_{21} & 1 & \cdots & r_{2p} \\ \vdots & \vdots & \cdots & \vdots \\ r_{p1} & r_{p2} & \cdots & 1 \end{bmatrix} \quad (8)$$

The above two matrix are **square** ($p \times p$) and **symmetric matrices**.

0.2 Revising essential matrix algebra

- The **norm** or the length of a vector $\mathbf{x} = (x_1, x_2)$ is calculated by taking the square root of the dot product of the vector with itself ($\sqrt{\mathbf{x} \cdot \mathbf{x}}$):

$$|\mathbf{x}| = \sqrt{x_1^2 + x_2^2} \quad (9)$$

We **normalize** a vector by dividing it by its length to obtain the corresponding **unit vector**. Usually when we compute **eigenvectors** - we usually normalize them and work with **unit eigenvectors**:

$$\mathbf{e} = \frac{\mathbf{x}}{\sqrt{\mathbf{x} \cdot \mathbf{x}}} \quad (10)$$

- The **spectral** decomposition of a matrix A is given below:

$$A = \lambda_1 e_1 e_1' + \lambda_2 e_2 e_2' + \cdots + \lambda_k e_k e_k' \quad (11)$$

In general we say that a $p \times p$ matrix A can be **diagonalized** or decomposed such that:

$$A = P \Lambda P' \quad (12)$$

Where P is the matrix of eigenvectors and Λ is the diagonal matrix of eigenvalues. Now for a $p \times p$ matrix there will be p eigenvalues and eigenvectors. But the spectral theorem states that we can take the **most significant eigenvalues and vectors and recompute the compressed form of the original matrix without losing much information**. This basically means that we can choose $k < p$ such that we can recompute the original matrix as given in equation 11.

- A **quadratic form** is nothing but a quadratic equation written in matrix form as $x'Ax$. The quadratic form is said to be **positive definite** if the following holds:

$$0 \leq x'Ax \quad (13)$$

- An **illustrative example for quadratic forms**. Consider the following quadratic equation:

$$3x_1^2 + 2x_2^2 - 2\sqrt{2}x_1x_2 \quad (14)$$

This can be packaged into a matrix form as follows:

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 3 & -\sqrt{2} \\ -\sqrt{2} & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x'Ax \quad (15)$$

Where A is nothing but the **coefficient matrix** of the quadratic equation. We know that the eigenvalues are solutions or roots to the characteristic equation given by:

$$|A - \lambda I| = (3 - \lambda)(2 - \lambda) - 2 = 0 \quad (16)$$

From the above example we find the eigenvalues to be $\lambda_1 = 4$ and $\lambda_2 = 1$. Using the **spectral decomposition** we can write:

$$A = \lambda_1 e_1 e_1' + \lambda_2 e_2 e_2' = 4e_1 e_1' + 1e_2 e_2' \quad (17)$$

- A general notation of the spectral decomposition is given as:

$$A = \sum_{i=1}^k \lambda_i e_i e_i' \quad (18)$$

Where the e_i vectors are **normalized eigenvectors**. Now if we put all the normalized eigenvectors into a matrix as column vectors of the matrix called P such that:

$$P = \begin{bmatrix} e_1 & \cdots & e_k \end{bmatrix} \quad (19)$$

Then we can write the spectral decomposition in matrix form as:

$$\mathbf{A} = \sum_{i=1}^k \lambda_i \mathbf{e}_i \mathbf{e}_i' = \mathbf{P} \mathbf{\Lambda} \mathbf{P}' \quad (20)$$

Where the matrix of eigenvectors \mathbf{P} is an **orthonormal matrix** that has the property that $\mathbf{P}\mathbf{P}' = \mathbf{P}'\mathbf{P} = \mathbf{I}$ and $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues:

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & \lambda_k \end{bmatrix} \quad (21)$$

- We can denote $\mathbf{\Lambda}^{1/2}$ as the diagonal matrix with $\sqrt{\lambda_i}$ as its diagonal entries. Hence we can then write the **square root matrix of \mathbf{A}** as follows:

$$\mathbf{A}^{1/2} = \sum_{i=1}^k \sqrt{\lambda_i} \mathbf{e}_i \mathbf{e}_i' = \mathbf{P} \mathbf{\Lambda}^{1/2} \mathbf{P}' \quad (22)$$

0.3 The geometry of Data

First of all we note that since multivariate data is presented in the form of matrices and vector and these are nothing but geometric objects plotted in some coordinate space, we can get important insights from the data just by looking at the geometric properties of these data vectors and matrices. Let us have a good look at our data matrix that we presented in the start:

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \cdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \quad (23)$$

Each row of \mathbf{X} represents a multivariate observation - we say that the entire dataset is a sample of n size from a p -variate population - which in turn means that our data has n measurements, each of which has p components. If we think of all **observations as p dimensional vectors** then we can represent the data matrix as composed of n , p -dimensional row vectors:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}'_1 \\ \mathbf{x}'_2 \\ \vdots \\ \mathbf{x}'_n \end{bmatrix} \quad (24)$$

Alternatively we can think of another geometrical interpretation of this data matrix by considering each column as a **feature vector** with n components. Here each

column vector corresponds to **all the measurements for each variable**. So we will have p column vector each of n components. The data matrix then becomes:

$$\mathbf{X} = \begin{bmatrix} \mathbf{y}_1 & \mathbf{y}_2 & \cdots & \mathbf{y}_p \end{bmatrix} \quad (25)$$

Just to clarify, in the above matrix, $\mathbf{y}_1 = [x_{11}, x_{21}, \dots, x_{n1}]$ are the n measurements on the first variable. Now we will look at a general geometric method of finding the sample mean.

- Let a vector of all 1s be represented as $\mathbf{1}$, where $\mathbf{1}_n$ is the n component vector that contains only 1s - $[1, 1, \dots, 1]$.
- Now to get the sample mean of the i^{th} feature or variable we carry out the following computation:

$$\bar{x}_i = \frac{1}{n} \mathbf{y}_i \cdot \mathbf{1} = \frac{(x_{1i} + x_{2i} + \dots + x_{ni})}{n} \quad (26)$$

If we further multiply this vector with the $\mathbf{1}$ vector again, we would get a vector with all the values as \bar{x}_i :

$$\bar{x}_i \mathbf{1} = \begin{bmatrix} \bar{x}_i \\ \vdots \\ \bar{x}_i \end{bmatrix} \quad (27)$$

- To find the **deviation vector** we carry out the following manipulations:

$$\mathbf{d}_i = \mathbf{y}_i - \bar{x}_i \mathbf{1} = \begin{bmatrix} x_{1i} - \bar{x}_i \\ x_{2i} - \bar{x}_i \\ \vdots \\ x_{ni} - \bar{x}_i \end{bmatrix} \quad (28)$$

- Some general geometric matrix transformations are given below:

$$\text{mean vector: } \bar{\mathbf{x}} = \frac{1}{n} \mathbf{X}' \mathbf{1} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} \quad (29)$$

$$\text{mean matrix: } \mathbf{1} \bar{\mathbf{x}}' = \frac{1}{n} \mathbf{1} \mathbf{1}' \mathbf{X} = \begin{bmatrix} \bar{x}_1 & \cdots & \bar{x}_p \\ \vdots & \cdots & \vdots \\ \bar{x}_1 & \cdots & \bar{x}_p \end{bmatrix} \quad (30)$$

$$\text{deviation matrix: } \mathbf{X} - \frac{1}{n} \mathbf{1} \mathbf{1}' \mathbf{X} = \begin{bmatrix} x_{11} - \bar{x}_1 & \cdots & x_{1p} - \bar{x}_p \\ \vdots & \cdots & \vdots \\ x_{n1} - \bar{x}_1 & \cdots & x_{np} - \bar{x}_p \end{bmatrix} \quad (31)$$

$$\text{deviation covariance matrix: } \mathbf{S} = \frac{1}{n-1} \mathbf{X}' \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}' \right) \mathbf{X} \quad (32)$$

References

- [1] Richard Arnold Johnson - Applied multivariate Statistical analysis