

Multivariate data format and descriptive measures

Why multivariate

- As the name implies, multivariate statistical analysis is the study of several random variables simultaneously.
- Usually, there is certain dependence structure in these random variables. Incorporating the inter-dependence structure into data analysis is more appropriate than the one-at-a-time univariate analysis.
- For large, high dimensional data with many variables, it is often necessary to reduce the data to a lower dimensional space that still retains most of the desirable information in the data.

In fact, statistics has been about reducing observations to a few statistics holding essential information.

Common objectives in multivariate statistical analysis

- Understand and concisely summarize data structure, which is often obscured by noise, random perturbations, and measurements errors.
- Understand and concisely summarize the relationship of one part of data to another.
- Inference from a sample of data to a much larger population: via parameter estimations, hypothesis tests, and statistical models.
- Univariate and multivariate statistical inference care about similar issues in data analysis and statistical model, such as central locations and variations, difference in treatment effects, outlier detection, and check violations of distribution assumptions, thus the validity of the analysis.

Characterizing inter-dependence structure and dimension reduction are two key aspects of multivariate analysis.

1 Multivariate data form

Multivariate data analysis deals with measurements of random outcomes that consist of a set of variables for each observation.

Notations

Multivariate observations are typically denoted as

x_{jk} = Measurement (a.k.a. observation, outcome, response) of the k th variable on the j th item (or subject)

Commonly, multivariate data are displayed as an array, rows label individual observations, columns index variables.

		variable 1	variable 2	...	variable k	...	variable p
(observation 1)	item 1	x_{11}	x_{12}	...	x_{1k}	...	x_{1p}
(observation 2)	item 2	x_{21}	x_{22}	...	x_{2k}	...	x_{2p}
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
(observation j)	item j	x_{j1}	x_{j2}	...	x_{jk}	...	x_{jp}
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
(observation n)	item n	x_{n1}	x_{n2}	...	x_{nk}	...	x_{np}

(1)

The data array can be expressed in vector-matrix form.

$$\mathbf{X} = [x_{jk}]_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2k} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{j1} & x_{j2} & \cdots & x_{jk} & \cdots & x_{jp} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} & \cdots & x_{np} \end{bmatrix} = \begin{bmatrix} \mathbf{x}'_1 \\ \mathbf{x}'_2 \\ \vdots \\ \mathbf{x}'_j \\ \vdots \\ \mathbf{x}'_n \end{bmatrix}$$

\leftarrow 1st (multivariate) observation
 \leftarrow 2nd observation
 \vdots
 \leftarrow j th observation
 \vdots
 \leftarrow n th observation

The convention is to express a vector \mathbf{x} as a column. For a vector, the operation "transpose" means to transform a column vector to a row (or to transform a row to a column vector). The notation \mathbf{x}' (or \mathbf{x}^T) denotes the transpose of a vector \mathbf{x} . Hence the j th row in the data matrix

$$[x_{j1} \cdots x_{jp}] = \mathbf{x}'_j = \mathbf{x}_j^T$$

is expressed as the transpose of a (column) vector, consisting of the p components of the j th observation. The column vector corresponding to the j th observation is

$$\mathbf{x}_j = \begin{bmatrix} x_{j1} \\ \vdots \\ x_{jp} \end{bmatrix} \in \mathbb{R}^p$$

which is a p -vector in the p -dimensional space of real numbers \mathbb{R}^p , where the symbol \in means "belong to". Naturally, a vector is the transpose of transpose of itself,

$$\mathbf{x}_j = [x_{j1} \cdots x_{jp}]' = (\mathbf{x}'_j)' = (\mathbf{x}_j^T)^T$$

A p -vector corresponds to thus identified with a point in \mathbb{R}^p ,

$$\mathbf{x}_j \text{ (as well as } \mathbf{x}'_j) \overset{\text{corresponding to}}{=} (x_{j1}, \cdots, x_{jp}) \in \mathbb{R}^p$$

The last expression in the above is in terms of the endpoint coordinates of the vector placed at origin. The coordinate notation is often used interchangeably with the corresponding vector or its transpose when the context is clear.

Observed data and random variables

Conventionally, lowercase x often indicates an observed value or a realization of a random variable X , which is commonly denoted by a capital letter.

For example, x_{jk} often represents the j th observation of a random variable X_k .

However, lower case is often used as its corresponding random counterpart as well. Therefore, whether x_{jk} is a scalar or a random variable depending heavily on the context. It is important to check the definitions.

Remarks on multivariate data

- In the above n -by- p data array, each observation is a row \mathbf{x}'_j , the transpose of the observed p -variate vector.

An alternative layout of multivariate data exchanges the rows and columns: every column is an observation of p -components, every row consists of n observed values of one component variable of the p -variate vector. So the data matrix of n p -variate observations is displayed as

$$\begin{bmatrix} x_{11} & x_{21} & \cdots & x_{k1} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{k2} & \cdots & x_{n2} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{1p} & x_{2p} & \cdots & x_{kp} & \cdots & x_{np} \end{bmatrix}$$

The advantage of this alternative layout is that the observed p -multivariate vector remains to be a column vector. Some statistical models, such as factor models, may be more convenient to use this data layout.

- Multivariate data analysis are useful when the component variables are correlated.
- Usually the measurements x_{ik} 's are real numbers, unless specified otherwise.
- Classical multivariate application focuses on the case $n > p$, with fixed p and fixed n .
Classical multivariate asymptotic results considers the case when $n \rightarrow \infty$, for fixed p .
- Commonly the notations \mathbf{X} , X , and x are used for multivariate matrix, random variable or vector, and observed outcome, respectively.

2 Descriptive statistics

For $k = 1, \dots, p$, the sample mean for the k th variable based on n observations is

$$\bar{x}_k = \frac{1}{n} \sum_{j=1}^n x_{jk}$$

(Question: Is it meaningful to also consider sample mean for each observation across vector components $k = 1, \dots, p$?)

The overall sample mean vector $\bar{\mathbf{x}} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix}$.

For $k = 1, \dots, p$, the sample variance for the k th variable based on n observations is

$$s_k^2 = \frac{1}{n-1} \sum_{j=1}^n (x_{jk} - \bar{x}_k)^2, \quad \left(\text{Alternatively } s_k^2 = \frac{1}{n} \sum_{j=1}^n (x_{jk} - \bar{x}_k)^2 \right)$$

The denominator $n-1$ instead of n is to ensure unbiasedness of s_k^2 as an estimator of the population variance σ_k^2 .

The sample covariance between the i th and k th variables based on n observations is

$$s_{ik} = \frac{1}{n-1} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)$$

Correspondingly, the sample correlation between the i th and k th variables based on n observations is

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}}\sqrt{s_{kk}}} = \frac{s_{ik}}{s_i \times s_k}$$

The overall sample variance-covariance matrix and sample correlation matrix based on n observations are

$$\mathbf{S} = [s_{ik}]_{p \times p} \quad \text{with} \quad s_{kk} = s_k^2, \quad \mathbf{R} = [r_{ik}]_{p \times p} \quad \text{with} \quad r_{kk} \equiv 1.$$

From matrix algebra point of view, the sample covariance and correlation matrix can be viewed as linear mappings of white noise to the observed data with given dependence structure in terms of variance-covariance.

For anyone yearning for a single number summary of data variation,

$$|\mathbf{S}| = \det(\mathbf{S})$$

is the generalized sample variance, which is the determinant of the sample variance-covariance matrix.

Vector-matrix representation of sample statistics

- The overall sample mean vector can be written as a sum of n vectors:

$$\bar{\mathbf{x}} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^n x_{j1}/n \\ \vdots \\ \sum_{j=1}^n x_{jp}/n \end{bmatrix} = \frac{1}{n} \sum_{j=1}^n \begin{bmatrix} x_{j1} \\ \vdots \\ x_{jp} \end{bmatrix} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

where \mathbf{x}_j is the j th observation, a p -vector $\in \mathbb{R}^p$.

- Writing the data array as $\mathbf{X} = [x_{jk}]_{n \times p}$, another vector-matrix expression for the mean vector is

$$\bar{\mathbf{x}} = \frac{1}{n} \mathbf{X}' \mathbf{1}_n$$

where $\mathbf{1}_n$ is the n -variate vector with elements $\equiv 1$.

- The sample covariance matrix \mathbf{S} can be expressed as a sum of n matrices.

$$\begin{aligned} \mathbf{S} = [s_{ik}]_{p \times p} &= \left[\frac{1}{n-1} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) \right]_{p \times p} \\ &= \frac{1}{n-1} \sum_{j=1}^n [(x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)]_{p \times p} = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \end{aligned}$$

where we use the vector product

$$(\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' = \begin{bmatrix} x_{j1} - \bar{x}_1 \\ x_{j2} - \bar{x}_2 \\ \vdots \\ x_{jp} - \bar{x}_p \end{bmatrix} [x_{j1} - \bar{x}_1, x_{j2} - \bar{x}_2, \dots, x_{jp} - \bar{x}_p] = [(x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)]_{i,k=1, \dots, p}$$

Note that

$$(\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'$$

is a $p \times p$ matrix for each fixed j , \mathbf{x}_j is the j th observation vector, $\bar{\mathbf{x}}$ is the sample mean vector. Therefore

$$\mathbf{S} = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'$$

is a sum of n $p \times p$ matrices, In the case of $p = 1$, we get back to the familiar univariate sample variance

$$S^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2$$

- Another vector-matrix expression is

$$\mathbf{S} = \frac{1}{n-1} (\mathbf{X} - \mathbf{1}_n \bar{\mathbf{x}}')^T (\mathbf{X} - \mathbf{1}_n \bar{\mathbf{x}}')$$

where $\mathbf{1}_n$ is the n -vector with elements 1 as previously defined. Recall $\bar{\mathbf{x}}' = [\bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_p]$, thus $\mathbf{1}_n \bar{\mathbf{x}}$ is

$$\mathbf{1}_n \bar{\mathbf{x}}' = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} [\bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_p] = \begin{bmatrix} \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_p \\ \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_p \\ \vdots & \vdots & \ddots & \vdots \\ \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_p \end{bmatrix}$$

Vector-matrix expressions are indispensably useful. For example, they are used in the proof of positive-definiteness of the sample covariance matrix \mathbf{S} , and in the derivation of properties of \mathbf{S} and the correlation matrix \mathbf{R} .

Properties of sample correlation of two component variables

- r_{jk} is the sample correlation of component variables j and k , which is not affected by whether n or $n - 1$ is used in the calculation of component variable sample variance s_k^2 .
- $r_{ik} \in [-1, 1]$.
- r_{ik} is a scale-invariant measure.
- r_{ik} is the Pearson correlation coefficient, which measures linear and only linear correlation.
- There exist other more general measures of dependence. For example, Kendall's τ and Spearman's ρ , as described in the section below.

Properties of sample covariance and correlation matrix

- Sample covariance matrix \mathbf{S} and sample correlation matrix \mathbf{R} are symmetric ($\mathbf{S}' = \mathbf{S}, \mathbf{R}' = \mathbf{R}$).
- Positive semi-definiteness of \mathbf{S} and \mathbf{R}

\mathbf{S} and \mathbf{R} are positive semi-definite, which means $\mathbf{v}'\mathbf{S}\mathbf{v} \geq 0$ and $\mathbf{v}'\mathbf{R}\mathbf{v} \geq 0$, for any p -vector \mathbf{v} .

Proof. Use the useful expression of $\mathbf{S} = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'$,

$$\mathbf{v}'\mathbf{S}\mathbf{v} = \frac{1}{n-1} \sum_{j=1}^n \mathbf{v}'(\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'\mathbf{v} = \frac{1}{n-1} \sum_{j=1}^n |\mathbf{v}'(\mathbf{x}_j - \bar{\mathbf{x}})|^2 \geq 0$$

for any p -vector \mathbf{v} . Therefore covariance matrix \mathbf{S} is positive semi-definite.

To show that correlation matrix \mathbf{R} is positive semi-definite, replace each x_{jk} by $x_{jk}^* = x_{jk}/s_k$ for $k = 1, \dots, p$, and $j = 1, \dots, n$. Then the sample covariance matrix of the n p -variate observations of \mathbf{x}_j^* is $\mathbf{S}^* = \mathbf{R}$, the sample correlation matrix of the original \mathbf{x}_j 's. Thus

$$\mathbf{v}'\mathbf{R}\mathbf{v} = \mathbf{v}'\mathbf{S}^*\mathbf{v} \geq 0$$

by the positive semi-definite property of covariance matrix. Therefore correlation matrix \mathbf{R} is also positive semi-definite. □

- An alternative proof of the positive semi-definite property

Notice that $y_j = \mathbf{v}'\mathbf{x}_j = \mathbf{x}_j'\mathbf{v}$ can be viewed as the j th observation of a univariate variable Y for $j = 1, \dots, n$. Then for any p -vector \mathbf{v} ,

$$\begin{aligned} \mathbf{v}'\mathbf{S}\mathbf{v} &= \frac{1}{n-1} \sum_{j=1}^n \mathbf{v}'(\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'\mathbf{v} \\ &= \frac{1}{n-1} \sum_{j=1}^n (\mathbf{v}'\mathbf{x}_j - \mathbf{v}'\bar{\mathbf{x}})(\mathbf{x}_j'\mathbf{v} - \bar{\mathbf{x}}'\mathbf{v}) \\ &= \frac{1}{n-1} \sum_{j=1}^n (y_j - \bar{y})(y_j - \bar{y}) = \frac{1}{n-1} \sum_{j=1}^n (y_j - \bar{y})^2 \end{aligned}$$

This is the sample covariance of the univariate variable y , which is always non-negative. Consequently the above equations provide another proof that \mathbf{S} is positive semi-definite.

- Definition of (strict) positive definiteness

If $\mathbf{v}'\mathbf{S}\mathbf{v} > 0$ for any p -vector $\mathbf{v} \neq 0$, then \mathbf{S} is a positive definite matrix.

Analogously, if $\mathbf{v}'\mathbf{R}\mathbf{v} > 0$ for any p -vector $\mathbf{v} \neq 0$, then \mathbf{R} is positive definite.

- Positive definite vs positive semi-definite

Positive definiteness is a desirable or even required property for covariance matrix.

\mathbf{S} and \mathbf{R} are either both (strictly) positive definite, or both positive semi-definite,

When \mathbf{S} is only positive semi-definite but not positive definite, then by definition, there is a non-zero vector $\mathbf{v} \in \mathbb{R}^p \setminus \{0_p\}$, such that $\mathbf{v}'\mathbf{S}\mathbf{v} = 0$.

Note that $\mathbf{y} = \mathbf{X}\mathbf{v} \in \mathbb{R}^n$ is a linear combination of the p column vectors of \mathbf{X} . We can derive (exercise) the sample variance of \mathbf{y} ,

$$\text{Var}(\mathbf{y}) = \text{Var}(\mathbf{X}\mathbf{v}) = \mathbf{v}'\text{Cov}(\mathbf{X})\mathbf{v} = \mathbf{v}'\mathbf{S}\mathbf{v} = 0$$

Thus $\mathbf{y} = \mathbf{X}\mathbf{v} = c\mathbf{1}$ has component $\equiv c$, for some constant $c \in \mathbb{R}$ (with probability 1).

In other words, the p columns of \mathbf{X} are linearly dependent, \mathbf{X} actually lies on a lower dimension hyperplane, perpendicular to a non-zero vector in \mathbb{R}^p .

This means the p variables are linearly correlated. At least one variable can be written as a linear combination of the other $p - 1$ variables. Therefore, at least one of the linear dependent variables should be taken out to remove the redundancy.

\mathbf{S} and \mathbf{R} are (strictly) positive definite when all $p < n$ variables are linearly independent.

- Relation of \mathbf{S} and \mathbf{R}

Let $\text{diag}(\mathbf{S}) = \text{diag}(s_1^2, \dots, s_p^2)$ be the diagonal matrix with s_1^2, \dots, s_p^2 as the diagonal elements and zero elsewhere. Define $D = (\text{diag}(\mathbf{S}))^{1/2} = \text{diag}(s_1, \dots, s_p)$. Then (exercise)

$$\mathbf{S} = \mathbf{D}\mathbf{R}\mathbf{D}, \quad \mathbf{R} = \mathbf{D}^{-1}\mathbf{S}\mathbf{D}^{-1}$$

Remarks

- If s_{ik} uses denominator n instead of $n - 1$, the notation for the corresponding covariance matrix $[s_{ik}]_{p \times p}$ is \mathbf{S}_n instead of \mathbf{S} , another matrix that is useful in multivariate statistics.
- Often it is more convenient to derive theoretical results using centered data with variable means $= 0$, that is, using

$$\mathbf{X}_c = [x_{jk} - \bar{x}_j]_{n \times p} = \mathbf{X} - \mathbf{1}_n \bar{\mathbf{x}}'$$

Then the covariance matrix of \mathbf{X} has a simplified and useful expression

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}_c' \mathbf{X}_c$$

- A geometric interpretation of the generalized variance is $|\mathbf{S}| = V_p^2 / (n - 1)^p$, where V_p is the volume generated by the p deviation (or centered) vectors.

3 Alternative measures of correlatedness and dependence

The sample correlation coefficient formula is the Pearson correlation coefficient which measures linear correlation. There are other metrics aim at measuring different dependence structures between two variables.

3.1 Kendall's rank correlation coefficient τ

Kendall's rank correlation coefficient (*a.k.a.* Kendall's τ) measures similarity of two variables by comparing the relative ordering of the two sets of ranks.

Let $(x_i, y_i), i = 1, \dots, n$ be observations of a bivariate random vector (X, Y) .

Consider $\{(x_i, y_i), (x_j, y_j)\}$ with $i < j$. There are $\binom{n}{2} = \frac{1}{2}n(n-1)$ many such pairs.

The pair is called concordant if $(x_j - x_i)(y_j - y_i) > 0$ and discordant if $(x_j - x_i)(y_j - y_i) < 0$.

Let

n_c = the number of concordant pairs,
 n_d = the number of discordant pairs.

Kendall's tau for the data $\{(x_i, y_i), i = 1, \dots, n\}$ is defined as

$$\tau = \frac{n_c - n_d}{\frac{1}{2}n(n-1)}$$

Remarks on Kendall's τ

- $\tau \in [-1, 1]$.
 $\tau = 1$ is the case of perfect agreement in the rankings of x and y .
 $\tau = -1$ is the case of complete ranking reversal.

- There are various modifications of τ in the case of ties.

- $n_c - n_d$ can be expressed as

$$n_c - n_d = \sum_{i < j} \text{sign}(x_j - x_i) \text{sign}(y_j - y_i)$$

where the *sign* function is defined as

$$\text{sign}(x) = \begin{cases} 1, & x > 0, \\ 0, & x = 0, \\ -1, & x < 0. \end{cases}$$

Concordant pairs have $\text{sign}(x_j - x_i) = \text{sign}(y_j - y_i)$, discordant pairs $\text{sign}(x_j - x_i) = -\text{sign}(y_j - y_i)$. Thus Kendall's τ for n pairs of observed data can also be written as

$$\tau = \frac{\sum_{i < j} \text{sign}(x_j - x_i) \cdot \text{sign}(y_j - y_i)}{\frac{1}{2}n(n-1)}$$

If we view (x_i, y_i) as observed values from independent random vectors (X_i, Y_i) , then kendal's τ for an i.i.d. (independently identically distributed) random sample can be written as

$$\tau = \frac{\sum_{i < j} \text{sign}(X_j - X_i) \cdot \text{sign}(Y_j - Y_i)}{\frac{1}{2}n(n-1)}$$

- For independent bivariate random vectors (X_1, Y_1) and (X_2, Y_2) with the same continuous cumulative joint distribution function, Kendall's τ can be defined by the joint probability function,

$$\tau = \mathbb{P}[(X_1 - X_2)(Y_1 - Y_2) > 0] - \mathbb{P}[(X_1 - X_2)(Y_1 - Y_2) < 0]$$

- For two (continuous) independent bivariate random variables (X_1, Y_1) and (X_2, Y_2) of the same distribution, an alternative expression is (exercise)

$$\tau = 2\mathbb{P}[(X_1 - X_2)(Y_1 - Y_2) > 0] - 1$$

- When (X_1, Y_1) and (X_2, Y_2) are independent and have the same continuous cumulative distribution function (CDF)

$$F(x, y) = \mathbb{P}(X_i < x, Y_i < y), \quad i = 1, 2,$$

then Kendall's τ can also be defined by the joint CDF as (exercise)

$$\tau = 4 \iint_{\mathbb{R}^2} F(x, y) dF(x, y) - 1$$

- Relation of τ and ρ for normal sample

If (X, Y) are of bivariate normal distribution with correlation coefficient r , then Kendall's τ gives slightly weaker correlation than the Pearson correlation r (other than the extreme of $\tau = 0$ and ± 1), with the relation

$$\tau = \frac{2}{\pi} \arcsin r$$

(The proof involves a few steps of variable transformations and inverse trigonometric function relations and is omitted here.)

The formula holds for more general distributions such as elliptical distributions (e.g., see Lindskog et al. 2001, Kendall's tau for Elliptical distributions), and provides a useful alternative for the estimation of r .

- Invariance property

By definition, Kendall's τ depends on relative orders of data within each variable only, thus it is invariant under strictly monotone transformation of the variables. This property is useful in the analysis of large data, such as in estimation of covariance and precision matrices for samples beyond Gaussian distributions.

3.2 Spearman's rank correlation coefficient ρ

Let $(x_i, y_i), i = 1, \dots, n$ be the component-wise ranks of n observations $(X_i, Y_i), i = 1, \dots, n$ of a bivariate random vector (X, Y) , each individual variable is ordered and indexed by i , so that x_i is the i th largest among all x 's, and y_i is the i th largest among all y 's, (or vice versa, meaning letting x_i be the i th smallest among all x 's, and y_i be the i th smallest among all y 's).

The Spearman's rank correlation coefficient is the Pearson correlation coefficient on the rank values of the observations.

$$\rho_s = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{j=1}^n (y_j - \bar{y})^2}}$$

In the absence of ties, Spearman's ρ can be expressed as (exercise)

$$\rho_s = 1 - \frac{6 \sum_{i=1}^n (x_i - y_i)^2}{n(n^2 - 1)}$$

For bivariate random vector (X_1, X_2) of continuous cumulative distribution function $F(x_1, x_2) = P(X_1 \leq x_1, X_2 \leq x_2)$ and of marginal cumulative distribution functions $F_1(x_1) = P(X_1 \leq x_1)$, $F_2(x_2) = P(X_2 \leq x_2)$, Spearman's ρ is the correlation of $F_1(X_1)$ and $F_2(X_2)$,

$$\rho_s = \text{Corr}(F_1(X_1), F_2(X_2))$$

Then

$$\rho_s = 12 \iint F_1(x_1) F_2(x_2) dF(x_1, x_2) - 3$$

Proof. For $i = 1, 2$, for any $u \in [0, 1]$, since F_i is necessarily continuous,

$$P(F_i(X_i) \leq u) = P(X_i \leq F_i^{-1}(u)) = F_i(F_i^{-1}(u)) = u, \quad \text{where } F_i^{-1}(u) = \inf\{t, F_i(t) \geq u\}.$$

Therefore $F_i(X_i) \sim U(0, 1)$, the uniform distribution on $(0, 1)$, which is of mean $\frac{1}{2}$ and variance $\frac{1}{12}$.

By the definition of ρ_s ,

$$\rho_s = \frac{\text{Cov}[F_1(X_1), F_2(X_2)]}{\sqrt{V(F_1(X_1))V(F_2(X_2))}} = 12 \times \text{Cov}[F_1(X_1), F_2(X_2)]$$

Using the notation $F(dx) = dF(x) = f(x)dx$ (when corresponding density f exists),

$$\begin{aligned} \text{Cov}[F_1(X_1), F_2(X_2)] &= E\left(F_1(X_1) - \frac{1}{2}\right)\left(F_2(X_2) - \frac{1}{2}\right) \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} \left(F_1(x_1) - \frac{1}{2}\right)\left(F_2(x_2) - \frac{1}{2}\right) F(dx_1, dx_2) \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} F_1(x_1) F_2(x_2) F(dx_1, dx_2) - \frac{1}{2} \int_{\mathbb{R}} \int_{\mathbb{R}} F_1(x_1) F(dx_1, dx_2) \\ &\quad - \frac{1}{2} \int_{\mathbb{R}} \int_{\mathbb{R}} F_2(x_2) F(dx_1, dx_2) + \frac{1}{4} \int_{\mathbb{R}} \int_{\mathbb{R}} F(dx_1, dx_2) \end{aligned}$$

In the last expression, the integral in the second term

$$\int_{\mathbb{R}} \int_{\mathbb{R}} F_1(x_1) F(dx_1, dx_2) = \int_{x_1 \in \mathbb{R}} F_1(x_1) \int_{x_2 \in \mathbb{R}} F(dx_1, dx_2) = \int_{\mathbb{R}} F_1(x_1) F_1(dx_1) = \int_0^1 u du = \frac{1}{2}$$

Similarly, the integral in the third term

$$\int_{\mathbb{R}} \int_{\mathbb{R}} F_2(x_2) F(dx_1, dx_2) = \int_{x_2 \in \mathbb{R}} F_2(x_2) \int_{x_1 \in \mathbb{R}} F(dx_1, dx_2) = \int_0^1 F_2(x_2) F_2(dx_2) = \int_0^1 u du = \frac{1}{2}$$

The last term $= \frac{1}{4}$, and the rest follows. \square

An example — Comparisons of correlation measures

Consider the heights (h_i 's in cm) and weights (w_i 's in kg) of four subjects. (*Height, Weight*) are dependent bivariate random variables. The observations (measurements) are in the table below, along with the correlations calculated by Pearson, Kendall, and Spearman methods. The values are quite different. Why, exactly?

Subject i	(h_i, w_i)	h_i rank	w_i rank	Concordant pairs	Discordant pairs	Correlation
A	(160, 45)	1	1	AB, AC, AD	BC, BD, CD	Pearson's $r = 0.7$
B	(170, 61)	2	4			Kendall's $\tau = 0.0$
C	(180, 60)	3	3			Spearman's $\rho = 0.2$
D	(190, 59)	4	2			

Discussions

If we change the value 45 in Subject A, in which situations that only one of the three correlations would change?

Example R commands for various correlations:

```
> cor(c(45,61,60,59),c(160,170,180,190),method="pearson") # default
> cor(c(45,61,60,59),c(160,170,180,190),method="spearman")
> cor(c(45,61,60,59),c(160,170,180,190),method="kendall")
```

Another example

With slight modifications to the example above, the following example produces different signs using different metrics of correlation.

Subject i	(h_i, w_i)	h_i rank	w_i rank	Concordant pairs	Discordant pairs	Correlation
A	(150, 59.1)	1	2	AB, AC	AD, BC, BD, CD	Pearson's $r = 0.005$
B	(170, 61.0)	2	4			Kendall's $\tau = -1/3$
C	(180, 60.0)	3	3			Spearman's $\rho = -0.4$
D	(190, 59.0)	4	1			

Remarks

Two variables can be dependent or “related” in various manners. It is important to know the measure used to quantify the dependence structure.

Note Related contents in the book by Johnson & Wichern: chapter 1, parts of chapters 2 and 3.