

第十章 典型相关分析

- §10.1 引言
- §10.2 总体典型相关
- §10.3 样本典型相关
- §10.4 典型相关系数的显著性检验

§10.1 引言

- 典型相关分析 (canonical correlation analysis) 是研究两组变量之间相关关系的一种统计分析方法, 它能够有效地揭示两组变量之间的相互线性依赖关系。
- 典型相关分析是由霍特林 (Hotelling, 1935, 1936) 首先提出的。

§10.2 总体典型相关

- 一、典型相关的定义及导出
- 二、典型相关变量的性质
- 三、从相关矩阵出发计算典型相关

一、典型相关的定义及导出

- 设 $\mathbf{x}=(x_1, x_2, \dots, x_p)'$ 和 $\mathbf{y}=(y_1, y_2, \dots, y_q)'$ 是两组随机变量, 且 $V(\mathbf{x})=\Sigma_{11}(>0)$, $V(\mathbf{y})=\Sigma_{22}(>0)$, $\text{Cov}(\mathbf{x}, \mathbf{y})=\Sigma_{12}$, 即有

$$V \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

其中 $\Sigma_{21}=\Sigma_{12}'$ 。

- 我们研究 $u=\mathbf{a}'\mathbf{x}$ 与 $v=\mathbf{b}'\mathbf{y}$ 之间的相关关系, 其中

$$\mathbf{a}=(a_1, a_2, \dots, a_p)', \quad \mathbf{b}=(b_1, b_2, \dots, b_q)'$$

现来计算一下 u 与 v 的相关系数。

$$\text{Cov}(u, v) = \text{Cov}(\mathbf{a}'\mathbf{x}, \mathbf{b}'\mathbf{y}) = \mathbf{a}'\text{Cov}(\mathbf{x}, \mathbf{y})\mathbf{b} = \mathbf{a}'\Sigma_{12}\mathbf{b}$$

$$V(u) = V(\mathbf{a}'\mathbf{x}) = \mathbf{a}'V(\mathbf{x})\mathbf{a} = \mathbf{a}'\Sigma_{11}\mathbf{a}$$

$$V(v) = V(\mathbf{b}'\mathbf{y}) = \mathbf{b}'V(\mathbf{y})\mathbf{b} = \mathbf{b}'\Sigma_{22}\mathbf{b}$$

所以, u 与 v 的相关系数

$$\rho(u, v) = \frac{\mathbf{a}'\Sigma_{12}\mathbf{b}}{\sqrt{\mathbf{a}'\Sigma_{11}\mathbf{a}}\sqrt{\mathbf{b}'\Sigma_{22}\mathbf{b}}}$$

由于对任意非零常数 k_1 和 k_2 , 有

$$\rho(k_1u, k_2v) = \rho(u, v)$$

因此, 为避免不必要的结果重复, 我们常常限定 u 与 v 均为标准化的变量, 即附加约束条件

$$V(u)=1, \quad V(v)=1$$

即

$$\mathbf{a}'\Sigma_{11}\mathbf{a}=1, \quad \mathbf{b}'\Sigma_{22}\mathbf{b}=1$$

在此约束条件下, 求 $\mathbf{a} \in R^p$ 和 $\mathbf{b} \in R^q$, 使得

$$\rho(u, v) = \mathbf{a}'\Sigma_{12}\mathbf{b}$$

达到最大。

容易证明, $\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ 和 $\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$ 有着相同的非零特征值, 且皆为正, 其个数为 $m = \text{rank}(\Sigma_{12})$ 。将这些正特征值分别记为 $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_m^2 > 0$ 。设 a_1, a_2, \dots, a_m 为 $\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ 的相应于 $\rho_1^2, \rho_2^2, \dots, \rho_m^2$ 的特征向量, 且满足标准化条件

$$a_i' \Sigma_{11} a_i = 1, \quad i=1, 2, \dots, m$$

令 $b_i = \frac{1}{\rho_i} \Sigma_{22}^{-1} \Sigma_{21} a_i$, 则有

$$\begin{aligned} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} b_i &= \frac{1}{\rho_i} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} (\Sigma_{22}^{-1} \Sigma_{21} a_i) \\ &= \frac{1}{\rho_i} \Sigma_{22}^{-1} \Sigma_{21} (\rho_i^2 a_i) = \rho_i^2 b_i \end{aligned}$$

从而 b_1, b_2, \dots, b_m 为 $\Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$ 的相应于 $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_m^2 > 0$ 的特征向量, 并且满足

$$\begin{aligned}
 \mathbf{b}'_i \boldsymbol{\Sigma}_{22} \mathbf{b}_i &= \frac{1}{\rho_i^2} \mathbf{a}'_i \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{22} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \mathbf{a}_i \\
 &= \frac{1}{\rho_i^2} \mathbf{a}'_i \boldsymbol{\Sigma}_{11} (\boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \mathbf{a}_i) \\
 &= \frac{1}{\rho_i^2} \mathbf{a}'_i \boldsymbol{\Sigma}_{11} (\rho_i^2 \mathbf{a}_i) = 1, \quad i = 1, 2, \dots, m
 \end{aligned}$$

- 可以证明, 当取 $\mathbf{a}=\mathbf{a}_1, \mathbf{b}=\mathbf{b}_1$ 时, $\rho(u,v)=\mathbf{a}'\boldsymbol{\Sigma}_{12}\mathbf{b}$ 达到最大值 ρ_1 (显然 $\rho_1 \leq 1$)。我们称

$$u_1 = \mathbf{a}'_1 \mathbf{x}, \quad v_1 = \mathbf{b}'_1 \mathbf{y}$$

为**第一对典型相关变量**, 称 ρ_1 为**第一个典型相关系数**

- 第一对典型相关变量 u_1, v_1 提取了原始变量 \mathbf{x} 与 \mathbf{y} 之间相关的主要部分, 如果这一部分还显得不够, 可以在剩余相关中再求出**第二对典型相关变量** $u_2 = \mathbf{a}'_2 \mathbf{x}, v_2 = \mathbf{b}'_2 \mathbf{y}$, 也就是 \mathbf{a}, \mathbf{b} 应满足标准化条件且应使得第二对典型相关变量不包括第一对典型相关

变量所含的信息，即

$$\rho(u_2, u_1) = \rho(\mathbf{a}'\mathbf{x}, \mathbf{a}_1'\mathbf{x}) = \text{Cov}(\mathbf{a}'\mathbf{x}, \mathbf{a}_1'\mathbf{x}) = \mathbf{a}'\Sigma_{11}\mathbf{a}_1 = 0$$

$$\rho(v_2, v_1) = \rho(\mathbf{b}'\mathbf{y}, \mathbf{b}_1'\mathbf{y}) = \text{Cov}(\mathbf{b}'\mathbf{y}, \mathbf{b}_1'\mathbf{y}) = \mathbf{b}'\Sigma_{22}\mathbf{b}_1 = 0$$

在这些约束条件下使得

$$\rho(u_2, v_2) = \rho(\mathbf{a}'\mathbf{x}, \mathbf{b}'\mathbf{y}) = \mathbf{a}'\Sigma_{12}\mathbf{b}$$

达到最大。

一般地，第 i ($1 < i \leq m$) 对典型相关变量 $u_i = \mathbf{a}'\mathbf{x}, v_i = \mathbf{b}'\mathbf{y}$ 是指，找出 $\mathbf{a} \in R^p, \mathbf{b} \in R^q$ ，在约束条件

$$\mathbf{a}'\Sigma_{11}\mathbf{a} = 1, \quad \mathbf{b}'\Sigma_{22}\mathbf{b} = 1$$

$$\mathbf{a}'\Sigma_{11}\mathbf{a}_k = 0, \quad \mathbf{b}'\Sigma_{22}\mathbf{b}_k = 0, \quad k = 1, 2, \dots, i-1$$

下，使得

$$\rho(u_i, v_i) = \rho(\mathbf{a}'\mathbf{x}, \mathbf{b}'\mathbf{y}) = \mathbf{a}'\Sigma_{12}\mathbf{b}$$

达到最大。当取 $\mathbf{a} = \mathbf{a}_i, \mathbf{b} = \mathbf{b}_i$ 时， $\rho(u_i, v_i)$ 达到最大值 ρ_i ，称它为第 i 个典型相关系数，称 $\mathbf{a}_i, \mathbf{b}_i$ 为第 i 对典型系数。

二、典型相关变量的性质

- 1.同一组的典型变量互不相关
- 2.不同组的典型变量之间的相关性
- 3.原始变量与典型变量之间的相关系数
- 4.简单相关、复相关和典型相关之间的关系

1.同一组的典型变量互不相关

· 设 x, y 的第 i 对典型变量为

$$u_i = \mathbf{a}_i' \mathbf{x}, \quad v_i = \mathbf{b}_i' \mathbf{y}, \quad i=1, 2, \dots, m$$

则有

$$V(u_i) = \mathbf{a}_i' \boldsymbol{\Sigma}_{11} \mathbf{a}_i = 1, \quad V(v_i) = \mathbf{b}_i' \boldsymbol{\Sigma}_{22} \mathbf{b}_i = 1, \quad i=1, 2, \dots, m$$

$$\rho(u_i, u_j) = \text{Cov}(u_i, u_j) = \mathbf{a}_i' \boldsymbol{\Sigma}_{11} \mathbf{a}_j = 0, \quad 1 \leq i \neq j \leq m$$

$$\rho(v_i, v_j) = \text{Cov}(v_i, v_j) = \mathbf{b}_i' \boldsymbol{\Sigma}_{22} \mathbf{b}_j = 0, \quad 1 \leq i \neq j \leq m$$

2.不同组的典型变量之间的相关性

$$\rho(u_i, v_j) = \rho_{ij}, \quad i=1, 2, \dots, m$$

$$\rho(u_i, v_j) = \text{Cov}(u_i, v_j) = \text{Cov}(\mathbf{a}'_i \mathbf{x}, \mathbf{b}'_j \mathbf{y}) = \mathbf{a}'_i \text{Cov}(\mathbf{x}, \mathbf{y}) \mathbf{b}_j$$

$$= \mathbf{a}'_i \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} \mathbf{b}_j = \rho_{ij} \mathbf{a}'_i \mathbf{a}_j = 0, \quad 1 \leq i \neq j \leq m$$

记 $\mathbf{u} = (u_1, u_2, \dots, u_m)'$, $\mathbf{v} = (v_1, v_2, \dots, v_m)'$, 则上述两个性质可用矩阵表示为

$$V(\mathbf{u}) = \mathbf{I}_m, \quad V(\mathbf{v}) = \mathbf{I}_m, \quad \text{Cov}(\mathbf{u}, \mathbf{v}) = \mathbf{A}$$

或

$$V \begin{pmatrix} \mathbf{u} \\ \mathbf{v} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_m & \mathbf{A} \\ \mathbf{A} & \mathbf{I}_m \end{pmatrix}$$

其中 $\mathbf{A} = \text{diag}(\rho_1, \rho_2, \dots, \rho_m)$ 。

3.原始变量与典型变量之间的相关系数

记

$$A=(a_1, a_2, \dots, a_m)=(a_{ij})_{p \times m}$$

$$B=(b_1, b_2, \dots, b_m)=(b_{ij})_{q \times m}$$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} = \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{1p} & \sigma_{1,p+1} & \cdots & \sigma_{1,p+q} \\ \vdots & & \vdots & \vdots & & \vdots \\ \sigma_{p1} & \cdots & \sigma_{pp} & \sigma_{p,p+1} & \cdots & \sigma_{p,p+q} \\ \sigma_{p+1,1} & \cdots & \sigma_{p+1,p} & \sigma_{p+1,p+1} & \cdots & \sigma_{p+1,p+q} \\ \vdots & & \vdots & \vdots & & \vdots \\ \sigma_{p+q,1} & \cdots & \sigma_{p+q,p} & \sigma_{p+q,p+1} & \cdots & \sigma_{p+q,p+q} \end{pmatrix}$$

则

$$\text{Cov}(\mathbf{x}, \mathbf{u}) = \text{Cov}(\mathbf{x}, \mathbf{A}'\mathbf{x}) = \boldsymbol{\Sigma}_{11}\mathbf{A}$$

$$\text{Cov}(\mathbf{x}, \mathbf{v}) = \text{Cov}(\mathbf{x}, \mathbf{B}'\mathbf{y}) = \boldsymbol{\Sigma}_{12}\mathbf{B}$$

$$\text{Cov}(\mathbf{y}, \mathbf{u}) = \text{Cov}(\mathbf{y}, \mathbf{A}'\mathbf{x}) = \boldsymbol{\Sigma}_{21}\mathbf{A}$$

$$\text{Cov}(\mathbf{y}, \mathbf{v}) = \text{Cov}(\mathbf{y}, \mathbf{B}'\mathbf{y}) = \boldsymbol{\Sigma}_{22}\mathbf{B}$$

上述四个等式也可表达为

$$\text{Cov}(x_i, u_j) = (\sigma_{i1}, \dots, \sigma_{ip}) \begin{pmatrix} a_{1j} \\ \vdots \\ a_{pj} \end{pmatrix} = \sum_{k=1}^p \sigma_{ik} a_{kj}$$

$$\text{Cov}(x_i, v_j) = (\sigma_{i,p+1}, \dots, \sigma_{i,p+q}) \begin{pmatrix} b_{1j} \\ \vdots \\ b_{qj} \end{pmatrix} = \sum_{k=1}^q \sigma_{i,p+k} b_{kj}$$

$$\text{Cov}(y_i, u_j) = (\sigma_{p+i,1}, \dots, \sigma_{p+i,p}) \begin{pmatrix} a_{1j} \\ \vdots \\ a_{pj} \end{pmatrix} = \sum_{k=1}^p \sigma_{p+i,k} a_{kj}$$

$$\text{Cov}(y_i, v_j) = (\sigma_{p+i,p+1}, \dots, \sigma_{p+i,p+q}) \begin{pmatrix} b_{1j} \\ \vdots \\ b_{qj} \end{pmatrix} = \sum_{k=1}^q \sigma_{p+i,p+k} b_{kj}$$

$$i=1,2,\dots,q, \quad j=1,2,\dots,m$$

所以

$$\rho(x_i, u_j) = \sum_{k=1}^p \sigma_{ik} a_{kj} / \sqrt{\sigma_{ii}}, \quad \rho(x_i, v_j) = \sum_{k=1}^q \sigma_{i,p+k} b_{kj} / \sqrt{\sigma_{ii}}$$

$$\rho(y_i, u_j) = \sum_{k=1}^p \sigma_{p+i,k} a_{kj} / \sqrt{\sigma_{p+i,p+i}}, \quad \rho(y_i, v_j) = \sum_{k=1}^q \sigma_{p+i,p+k} b_{kj} / \sqrt{\sigma_{p+i,p+i}}$$

$$i = 1, 2, \dots, p, \quad j = 1, 2, \dots, m$$

4. 简单相关、复相关和典型相关之间的关系

- 当 $p=q=1$ 时， x 与 y 之间的（惟一）典型相关就是它们之间的简单相关；当 $p=1$ 或 $q=1$ 时， x 与 y 之间的（惟一）典型相关就是它们之间的复相关。可见，复相关是典型相关的一个特例，而简单相关是复相关的一个特例。
- 第一个典型相关系数至少同 x （或 y ）的任一分量与 y （或 x ）的复相关系数一样大，即使所有这些复相关系数都较小，第一个典型相关系数仍可能很大；同样，从复相关的定义也可以看出，当 $p=1$ （或 $q=1$ ）时， x （或 y ）与 y （或 x ）之间的复相关系数也不会小于 x （或 y ）与 y （或 x ）的任一分量之间的相关系数，即使所有这些相关系数都较小，复相关系数仍可能很大。

三、从相关矩阵出发计算典型相关

- 有时， x 和 y 的各分量的单位不全相同，我们希望在对各分量作标准化变换之后再作典型相关分析。

- 记 $\mu_1 = E(x)$, $\mu_2 = E(y)$, $D_1 = \text{diag}(\sqrt{\sigma_{11}}, \dots, \sqrt{\sigma_{pp}})$

$$D_2 = \text{diag}(\sqrt{\sigma_{p+1,p+1}}, \dots, \sqrt{\sigma_{p+q,p+q}}), \quad R = \begin{pmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{pmatrix} \text{为} \begin{pmatrix} x \\ y \end{pmatrix} \text{的}$$

相关矩阵。对 x 和 y 的各分量作标准化变换，即令

$$x^* = D_1^{-1}(x - \mu_1), \quad y^* = D_2^{-1}(y - \mu_2)$$

现在来求 x^* 和 y^* 的典型相关变量， $i=1, 2, \dots, m$ 。

$$V(\mathbf{x}^*) = \mathbf{D}_1^{-1} V(\mathbf{x}) \mathbf{D}_1^{-1} = \mathbf{D}_1^{-1} \boldsymbol{\Sigma}_{11} \mathbf{D}_1^{-1} = \mathbf{R}_{11}$$

$$V(\mathbf{y}^*) = \mathbf{D}_2^{-1} V(\mathbf{y}) \mathbf{D}_2^{-1} = \mathbf{D}_2^{-1} \boldsymbol{\Sigma}_{22} \mathbf{D}_2^{-1} = \mathbf{R}_{22}$$

$$\text{Cov}(\mathbf{x}^*, \mathbf{y}^*) = \mathbf{D}_1^{-1} \text{Cov}(\mathbf{x}, \mathbf{y}) \mathbf{D}_2^{-1} = \mathbf{D}_1^{-1} \boldsymbol{\Sigma}_{12} \mathbf{D}_2^{-1} = \mathbf{R}_{12}$$

$$\text{Cov}(\mathbf{y}^*, \mathbf{x}^*) = \mathbf{D}_2^{-1} \text{Cov}(\mathbf{y}, \mathbf{x}) \mathbf{D}_1^{-1} = \mathbf{D}_2^{-1} \boldsymbol{\Sigma}_{21} \mathbf{D}_1^{-1} = \mathbf{R}_{21}$$

于是

$$\begin{aligned} \mathbf{R}_{11}^{-1} \mathbf{R}_{12} \mathbf{R}_{22}^{-1} \mathbf{R}_{21} &= (\mathbf{D}_1^{-1} \boldsymbol{\Sigma}_{11} \mathbf{D}_1^{-1})^{-1} \mathbf{D}_1^{-1} \boldsymbol{\Sigma}_{12} \mathbf{D}_2^{-1} (\mathbf{D}_2^{-1} \boldsymbol{\Sigma}_{22} \mathbf{D}_2^{-1})^{-1} \mathbf{D}_2^{-1} \boldsymbol{\Sigma}_{21} \mathbf{D}_1^{-1} \\ &= \mathbf{D}_1 \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \mathbf{D}_1^{-1} \end{aligned}$$

因为

$$\boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \mathbf{a}_i = \rho_i^2 \mathbf{a}_i$$

$$\mathbf{D}_1 \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \mathbf{D}_1^{-1} (\mathbf{D}_1 \mathbf{a}_i) = \rho_i^2 (\mathbf{D}_1 \mathbf{a}_i)$$

所以

$$\mathbf{R}_{11}^{-1} \mathbf{R}_{12} \mathbf{R}_{22}^{-1} \mathbf{R}_{21} \mathbf{a}_i^* = \rho_i^2 \mathbf{a}_i^*$$

式中 $\mathbf{a}_i^* = \mathbf{D}_1 \mathbf{a}_i$, 有 $\mathbf{a}_i^{*'} \mathbf{R}_{11} \mathbf{a}_i^* = \mathbf{a}_i' \mathbf{D}_1 \mathbf{R}_{11} \mathbf{D}_1 \mathbf{a}_i = \mathbf{a}_i' \boldsymbol{\Sigma}_{11} \mathbf{a}_i = 1$ 。同理

$$\mathbf{R}_{22}^{-1} \mathbf{R}_{21} \mathbf{R}_{11}^{-1} \mathbf{R}_{12} \mathbf{b}_i^* = \rho_i^2 \mathbf{b}_i^*$$

式中 $\mathbf{b}_i^* = \mathbf{D}_2 \mathbf{b}_i$, 有 $\mathbf{b}_i^{*'} \mathbf{R}_{22} \mathbf{b}_i^* = \mathbf{b}_i' \mathbf{D}_2 \mathbf{R}_{22} \mathbf{D}_2 \mathbf{b}_i = \mathbf{b}_i' \boldsymbol{\Sigma}_{22} \mathbf{b}_i = 1$ 。由此可见, $\mathbf{a}_i^*, \mathbf{b}_i^*$ 为 \mathbf{x}^* 和 \mathbf{y}^* 的第 i 对典型系数, 其第 i 个典型相关系数仍为 ρ_i , 在标准化变换下具有不变性, 这一点与主成分分析有所不同。

- \mathbf{x}^* 和 \mathbf{y}^* 的第 i 对典型变量 $u_i^* = \mathbf{a}_i^{*'} \mathbf{x}^*$, $v_i^* = \mathbf{b}_i^{*'} \mathbf{y}^*$ 具有零均值, 且与 \mathbf{x} 和 \mathbf{y} 的第 i 对典型变量 $u_i = \mathbf{a}_i' \mathbf{x}$, $v_i = \mathbf{b}_i' \mathbf{y}$ 只相差一个常数。

- 例10.2.1 设 \mathbf{x}, \mathbf{y} 有如下相关矩阵:

$$\mathbf{R}_{11} = \begin{pmatrix} 1 & \alpha \\ \alpha & 1 \end{pmatrix}, \quad \mathbf{R}_{22} = \begin{pmatrix} 1 & \gamma \\ \gamma & 1 \end{pmatrix}, \quad \mathbf{R}_{12} = \begin{pmatrix} \beta & \beta \\ \beta & \beta \end{pmatrix} = \beta \mathbf{1} \mathbf{1}'$$

这里 $|\alpha| < 1$, $|\gamma| < 1$, 可以保证 $\mathbf{R}_{11}^{-1}, \mathbf{R}_{22}^{-1}$ 存在。

$$\begin{aligned}
 \mathbf{R}_{11}^{-1} \mathbf{R}_{12} \mathbf{R}_{22}^{-1} \mathbf{R}_{21} &= \frac{1}{1-\alpha^2} \begin{pmatrix} 1 & -\alpha \\ -\alpha & 1 \end{pmatrix} \beta \mathbf{1} \mathbf{1}' \bullet \frac{1}{1-\gamma^2} \begin{pmatrix} 1 & -\gamma \\ -\gamma & 1 \end{pmatrix} \beta \mathbf{1} \mathbf{1}' \\
 &= \frac{\beta^2}{(1-\alpha^2)(1-\gamma^2)} \begin{pmatrix} 1-\alpha \\ 1-\alpha \end{pmatrix} (1-\gamma, 1-\gamma) \mathbf{1} \mathbf{1}' \\
 &= \frac{\beta^2}{(1+\alpha)(1+\gamma)} \mathbf{1} \mathbf{1}' \mathbf{1} \mathbf{1}' = \frac{2\beta^2}{(1+\alpha)(1+\gamma)} \mathbf{1} \mathbf{1}'
 \end{aligned}$$

由于 $\mathbf{1} \mathbf{1}'$ 有唯一的非零特征值 $\mathbf{1}' \mathbf{1} = 2$, 故 $\mathbf{R}_{11}^{-1} \mathbf{R}_{12} \mathbf{R}_{22}^{-1} \mathbf{R}_{21}$ 有唯一非零特征值

$$\rho_1^2 = \frac{4\beta^2}{(1+\alpha)(1+\gamma)}$$

- 在约束条件 $\mathbf{a}_1^* \mathbf{R}_{11} \mathbf{a}_1^* = 1$ 下, 相应于特征值 ρ_1^2 的特征向量为 $\mathbf{a}_1^* = [2(1+\alpha)]^{-1/2} \mathbf{1}$ 。同理, 在约束条件 $\mathbf{b}_1 \mathbf{R}_{22} \mathbf{b}_1^* = 1$ 下,

$R_{22}^{-1}R_{21}R_{11}^{-1}R_{12}$ 相应于特征值 ρ_1^2 的特征向量为 $b_1^* = [2(1+\gamma)]^{-1/2} \mathbf{1}$

所以，第一对典型相关变量为

$$a_1^* x^* = [2(1+\alpha)]^{-1/2} \mathbf{1}' x^*, b_1^* y^* = [2(1+\gamma)]^{-1/2} \mathbf{1}' y^*$$

其中 x^* 和 y^* 分别是对 x 和 y 各分量标准化后的向量。第一个典型相关系数为 $\rho_1 = 2|\beta| / [(1+\alpha)(1+\gamma)]^{1/2}$ 。由于 $|\alpha| < 1$, $|\gamma| < 1$, 故 $\rho_1 > |\beta|$, 表明第一个典型相关系数大于两组原始变量之间的相关系数。

§10.3 样本典型相关

· 设数据矩阵为

$$(X \ Y) = \begin{pmatrix} \mathbf{x}'_1 & \mathbf{y}'_1 \\ \vdots & \vdots \\ \mathbf{x}'_n & \mathbf{y}'_n \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1p} & y_{11} & \cdots & y_{1q} \\ \vdots & & \vdots & \vdots & & \vdots \\ x_{n1} & \cdots & x_{np} & y_{n1} & \cdots & y_{nq} \end{pmatrix}$$

则样本协方差矩阵为

$$\mathbf{S} = \begin{pmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{pmatrix}$$

\mathbf{S} 可用来作为 Σ 的估计。当 $n > p+q$ 时, $\mathbf{S}_{11}^{-1} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{21}$ 和 $\mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12}$ 可分别作为 $\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$ 和 $\Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$ 的估计; 它们的非零特征值 $r_1^2 \geq r_2^2 \geq \cdots \geq r_m^2$ 可用来估计 $\rho_1^2 \geq \rho_2^2 \geq \cdots \geq \rho_m^2$;

- 相应的特征向量 $\hat{a}_1, \hat{a}_2, \dots, \hat{a}_m$ 作为 a_1, a_2, \dots, a_m 的估计, $\hat{b}_1, \hat{b}_2, \dots, \hat{b}_m$ 作为 b_1, b_2, \dots, b_m 的估计。 r_j^2 的正平方根 r_j 称为第 j 个样本典型相关系数, $\hat{a}'_j x$ 和 $\hat{b}'_j y$ 称为第 j 对样本典型相关变量, $j=1, 2, \dots, m$ 。将样本 (x_i, y_i) , $i=1, 2, \dots, n$ 经中心化后代入 m 对典型变量, 即令

则称 $u_{ij} = \hat{a}'_j (x_i - \bar{x})$ 为第 i 个样品的第 j 个样本典型变量得分, 称 $v_{ij} = \hat{b}'_j (y_i - \bar{y})$ 为第 i 个样品的第 j 个样本典型变量得分。由约束条件 可得

$$\hat{a}'_j S_{11} \hat{a}_j = 1$$

- 同理可得 $\frac{1}{n-1} \sum_{i=1}^n u_{ij}^2 = \frac{1}{n-1} \hat{a}'_j \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})' \hat{a}_j = \hat{a}'_j S_{11} \hat{a}_j = 1, \quad j = 1, 2, \dots, m$

- 对每个 j , 可画出 $(u_{ij}, v_{ij})_{i=1}^n$ 的散点图, 该图也可用来检查是否有异常值出现。

- 例10.3.1 某康复俱乐部对20名中年人测量了三个生理指标：体重(x_1)、腰围(x_2)、脉搏(x_3)和三个训练指标：引体向上(y_1)、起坐次数(y_2)、跳跃次数(y_3)。其数据列于表10.3.1。

表10.3.1 某康复俱乐部的生理指标和训练指标数据

编号	x_1	x_2	x_3	y_1	y_2	y_3
1	191	36	50	5	162	60
2	189	37	52	2	110	60
3	193	38	58	12	101	101
4	162	35	62	12	105	37
5	189	35	46	13	155	58
6	182	36	56	4	101	42
7	211	38	56	8	101	38

8	167	34	60	6	125	40
9	176	31	74	15	200	40
10	154	33	56	17	251	250
11	169	34	50	17	120	38
12	166	33	52	13	210	115
13	154	34	64	14	215	105
14	247	46	50	1	50	50
15	193	36	46	6	70	31
16	202	37	62	12	210	120
17	176	37	54	4	60	25
18	157	32	52	11	230	80
19	156	33	54	15	225	73
20	138	33	68	2	110	43

$$\hat{\mathbf{R}}_{11} = \begin{pmatrix} 1 & & \\ 0.870 & 1 & \\ -0.366 & -0.353 & 1 \end{pmatrix}, \quad \hat{\mathbf{R}}_{22} = \begin{pmatrix} 1 & & \\ 0.696 & 1 & \\ 0.496 & 0.669 & 1 \end{pmatrix}$$

$$\hat{\mathbf{R}}_{12} = \hat{\mathbf{R}}'_{21} = \begin{pmatrix} -0.390 & -0.493 & -0.226 \\ -0.552 & -0.646 & -0.192 \\ 0.151 & 0.225 & 0.035 \end{pmatrix}$$

$\hat{\mathbf{R}}_{11}^{-1} \hat{\mathbf{R}}_{12} \hat{\mathbf{R}}_{22}^{-1} \hat{\mathbf{R}}_{21}$ 的特征值分别为0.6630、0.0402和0.0053
， 于是

$$r_1=0.797, \quad r_2=0.201, \quad r_3=0.073$$

相应的样本典型变量系数为

$$\hat{\mathbf{a}}_1^* = \begin{pmatrix} -0.775 \\ 1.579 \\ -0.059 \end{pmatrix}, \quad \hat{\mathbf{a}}_2^* = \begin{pmatrix} -1.884 \\ 1.181 \\ -0.231 \end{pmatrix}, \quad \hat{\mathbf{a}}_3^* = \begin{pmatrix} -0.191 \\ 0.506 \\ 1.051 \end{pmatrix}$$

$$\hat{\mathbf{b}}_1^* = \begin{pmatrix} -0.350 \\ -1.054 \\ 0.716 \end{pmatrix}, \quad \hat{\mathbf{b}}_2^* = \begin{pmatrix} -0.376 \\ 0.124 \\ 1.062 \end{pmatrix}, \quad \hat{\mathbf{b}}_3^* = \begin{pmatrix} -1.297 \\ 1.237 \\ -0.419 \end{pmatrix}$$

因此，第一对样本典型变量为

$$u_1^* = -0.775x_1^* + 1.579x_2^* - 0.059x_3^*$$

$$v_1^* = -0.350y_1^* - 1.054y_2^* + 0.716y_3^*$$

· 如果需要，第二对样本典型变量为

$$u_2^* = -1.884x_1^* + 1.181x_2^* - 0.231x_3^*$$

$$v_2^* = -0.376y_1^* + 0.124y_2^* + 1.062y_3^*$$

$$\hat{\mathbf{R}}_{22} = \begin{pmatrix} 1.00 & & & & & & \\ 0.43 & 1.00 & & & & & \\ 0.27 & 0.33 & 1.00 & & & & \\ 0.24 & 0.26 & 0.25 & 1.00 & & & \\ 0.34 & 0.54 & 0.46 & 0.28 & 1.00 & & \\ 0.37 & 0.32 & 0.29 & 0.30 & 0.35 & 1.00 & \\ 0.40 & 0.58 & 0.45 & 0.27 & 0.59 & 0.31 & 1.00 \end{pmatrix}$$

$$\hat{\mathbf{R}}_{12} = \hat{\mathbf{R}}'_{21} = \begin{pmatrix} 0.33 & 0.32 & 0.20 & 0.19 & 0.30 & 0.37 & 0.21 \\ 0.30 & 0.21 & 0.16 & 0.08 & 0.27 & 0.35 & 0.20 \\ 0.31 & 0.23 & 0.14 & 0.07 & 0.24 & 0.37 & 0.18 \\ 0.24 & 0.22 & 0.12 & 0.19 & 0.21 & 0.29 & 0.16 \\ 0.38 & 0.32 & 0.17 & 0.23 & 0.32 & 0.36 & 0.27 \end{pmatrix}$$

样本典型相关系数和样本典型变量系数列于表10.3.2中。

表10.3.2

典型相关系数和典型变量系数

标准化变量	\hat{a}_1^*	\hat{a}_2^*	\hat{a}_3^*	\hat{a}_4^*	\hat{a}_5^*
x_1^*	0.42	0.34	-0.86	-0.79	0.03
x_2^*	0.20	-0.67	0.44	-0.27	0.98
x_3^*	0.17	-0.85	-0.26	0.47	-0.91
x_4^*	-0.02	0.36	-0.42	1.04	0.52
x_5^*	0.46	0.73	0.98	-0.17	-0.44
r_j	0.55	0.24	0.12	0.07	0.06
标准化变量	\hat{b}_1^*	\hat{b}_2^*	\hat{b}_3^*	\hat{b}_4^*	\hat{b}_5^*
y_1^*	0.43	-0.09	0.49	-0.13	-0.48
y_2^*	0.21	0.44	-0.78	-0.34	-0.75
y_3^*	-0.04	-0.09	-0.48	-0.61	0.35
y_4^*	0.02	0.93	-0.01	0.40	0.31
y_5^*	0.29	-0.10	0.28	-0.45	0.70
y_6^*	0.52	-0.55	-0.41	0.69	0.18
y_7^*	-0.11	-0.03	0.93	0.27	-0.01

第一对样本典型变量为

$$u_1^* = 0.42x_1^* + 0.20x_2^* + 0.17x_3^* - 0.02x_4^* + 0.46x_5^*$$

$$v_1^* = 0.43y_1^* + 0.21y_2^* - 0.04y_3^* + 0.02y_4^* + 0.29y_5^* + 0.52y_6^* - 0.11y_7^*$$

- 根据典型系数，主要代表了用户反馈和自主权这两个变量，三个任务变量显得并不重要；而主要代表了主管满意度和工种满意度变量，其次代表了事业前景满意度和公司地位满意度变量。我们也可从相关系数的角度来解释典型变量，原始变量与第一对典型变量间的样本相关系数列于表10.3.3中。

表10.3.3 原始变量与典型变量的样本相关系数

原始变量 x	样本典型变量		原始变量 y	样本典型变量	
	u_1^*	v_1^*		u_1^*	v_1^*
x_1 : 用户反馈	0.83	0.46	y_1 : 主管满意度	0.42	0.76
x_2 : 任务重要性	0.73	0.40	y_2 : 事业前景满意度	0.36	0.64
x_3 : 任务多样性	0.75	0.42	y_3 : 财政满意度	0.21	0.39
x_4 : 任务特性	0.62	0.34	y_4 : 工作强度满意度	0.21	0.38
x_5 : 自主权	0.86	0.48	y_5 : 公司地位满意度	0.36	0.65
			y_6 : 工种满意度	0.45	0.80
			y_7 : 总体满意度	0.28	0.50

- 所有五个职业特性变量与第一典型变量 u_1^* 有大致相同的相关系数，故 u_1^* 可以解释为职业特性变量，这与基于典型系数的解释不同。 v_1^* 主要代表了主管满意度、事业前景满意度、公司地位满意度和工种满意度， v_1^* 可以解释为职业满意度—公司地位变量，这与基于典型系数的解释基本相一致。第一对典型变量 u_1^* 与 v_1^* 的样本相关系数 $r_1=0.55$ ，可见，职业特性与职业满意度之间有一定程度的相关性。

§10.4 典型相关系数的显著性检验

- 一、全部总体典型相关系数均为零的检验
- 二、部分总体典型相关系数为零的检验

一、全部总体典型相关系数均为零的检验

· 设 $(x', y')' \sim N_{p+q}(\mu, \Sigma), \Sigma > 0$ 。又设 S 为样本协方差矩阵, 且 $n > p+q$ 。

· 考虑假设检验问题:

$$H_0: \rho_1 = \rho_2 = \cdots = \rho_m = 0$$

$$H_1: \rho_1, \rho_2, \cdots, \rho_m \text{ 至少有一个不为零}$$

其中 $m = \min\{p, q\}$ 。若检验接受 H_0 , 则认为讨论两组变量之间的相关性没有意义; 若检验拒绝 H_0 , 则认为第一对典型变量是显著的。(10.4.1)式实际上等价于假设检验问题

$$H_0: \Sigma_{12} = \mathbf{0}, \quad H_1: \Sigma_{12} \neq \mathbf{0}$$

H_0 成立表明 x 与 y 互不相关。

检验统计量为

$$\Lambda_1 = \prod_{i=1}^m (1 - r_i^2)$$

对于充分大的 n , 当 H_0 成立时, 统计量

$$Q_1 = - \left[n - \frac{1}{2}(p + q + 3) \right] \ln \Lambda_1 \sim \chi^2(pq)$$

在给定的 α 下, 若 $Q_1 \geq \chi_{\alpha}^2(pq)$, 则拒绝 H_0 , 认为典型变量 u_1 与 v_1 之间的相关性是显著的; 否则, 就认为第一个典型相关系数不显著。

· 例10.4.1 在例10.3.1中, 假设为多元正态数据, 欲检验:

$$H_0: \rho_1 = \rho_2 = \rho_3 = 0, \quad H_1: \rho_1 \neq 0$$

它的似然比统计量为

$$A_1 = (1 - r_1^2)(1 - r_2^2)(1 - r_3^2)$$

$$= (1 - 0.6330)(1 - 0.0402)(1 - 0.0053) = 0.3504$$

$$Q_1 = - \left[20 - \frac{1}{2}(3 + 3 + 3) \right] \ln A_1 = -15.5 \times \ln 0.3504 = 16.255$$

查 χ^2 分布表得, $\chi_{0.10}^2(9) = 14.684$, $\chi_{0.05}^2(9) = 16.919$, 因此在 $\alpha=0.10$ 的显著性水平下, 拒绝原假设 H_0 , 也即认为至少有一个典型相关是显著的。

二、部分总体典型相关系数为零的检验

- 若 $H_0: \rho_1 = \rho_2 = \dots = \rho_m = 0$ 经检验被拒绝, 则应进一步检验假设

$$H_0: \rho_2 = \dots = \rho_m = 0$$

$$H_1: \rho_2, \dots, \rho_m \text{ 至少有一个不为零}$$

若原假设 H_0 被接受, 则认为只有第一对典型变量是有用的;

若原假设 H_0 被拒绝, 则认为第二对典型变量也是有用的。

- 如此进行下去, 直至对某个 k , 假设 $H_0: \rho_{k+1} = \dots = \rho_m = 0$ 被接受, 这时可认为只有前 k 对典型变量是显著的。

- 对于假设检验问题

$$H_0: \rho_{k+1} = \dots = \rho_m = 0$$

$$H_1: \rho_{k+1}, \dots, \rho_m \text{ 至少有一个不为零}$$

其检验统计量为

$$\Lambda_{k+1} = \prod_{i=k+1}^m (1 - r_i^2)$$

对于充分大的 n ，当 H_0 为真时，统计量

$$Q_{k+1} = - \left[n - k - \frac{1}{2}(p + q + 3) + \sum_{i=1}^k r_i^{-2} \right] \ln \Lambda_{k+1}$$

近似服从自由度为 $(p-k)(q-k)$ 的 χ^2 分布。给定显著性水平 α ，若 $Q_{k+1} \geq \chi_{\alpha}^2[(p-k)(q-k)]$ ，则拒绝原假设 H_0 ，认为第 $k+1$ 个典型相关系数 ρ_{k+1} 是显著的，即第 $k+1$ 对典型变量显著相关。

以上的一系列检验实际上是一个序贯检验，检验直到对某个 k 值 H_0 未被拒绝为止。事实上，检验的总显著性水平已不是 α 了，且难以确定。还有，检验的结果易受样本容量大小的影响。因此，检验的结果只宜作为确定典型变量个数的重要参考依据，而不宜作为惟一的依据。通常选择尽可能小的 k 。

· 例10.4.2 在例10.3.1中, 欲进一步检验:

$$H_0: \rho_2 = \rho_3 = 0, H_1: \rho_2 \neq 0$$

检验统计量为

$$A_2 = (1 - r_2^2)(1 - r_3^2) = (1 - 0.0402)(1 - 0.0053) = 0.9547$$

$$Q_2 = - \left[20 - 1 - \frac{1}{2}(3 + 3 + 3) + r_1^{-2} \right] \ln A_2$$

$$= -16.08 \times \ln 0.9547 = 0.745 < 7.779 = \chi_{0.10}^2(4)$$

故接受原假设 H_0 , 即认为第二个典型相关是不显著的。因此, 只有一个典型相关是显著的。