

## 第四讲. 相关系数 (续)

2014.2.27

1

### 非正态情形下的大样本相关性检验和置信区间

非正态假设下, 或即使假设正态分布但  $\rho$  非零情形下, 样本相关系数分布如何?

#### 样本相关系数的大样本分布

性质1. 样本:  $(x_1, y_1), \dots, (x_n, y_n)$  iid, 设  $\rho$  为总体相关系数,  $r$  为样本相关系数, 则(可参看有关数理统计教程)

$$\sqrt{n}(r - \rho) \xrightarrow{d} N(0, (1 - \rho^2)^2), \text{ 当 } n \rightarrow \infty$$

2

#### (1). 假设检验

相关性的大样本检验 (不假设联合正态分布):

$H_0$  成立时, 当  $n \rightarrow \infty$  时,

$$z = \sqrt{n} \times r \xrightarrow{d} N(0, 1)$$

则  $p$  值  $\approx P(|N(0, 1)| > |z|) = 2(1 - \Phi(|z|))$ ,  $\Phi$  是  $N(0, 1)$  累积分布函数。

注: 通常使用  $z = \sqrt{n-2} \times r \sim N(0, 1)$

3

基于性质1, 可以如下构造  $\rho$  的95%置信区间:  $\left\{ \rho: \left| \frac{\sqrt{n}(r - \rho)}{1 - \rho^2} \right| \leq 1.96 \right\}$ ,  
但通常使用相关系数的Fisher变换构造置信区间。

#### (2). 置信区间

$\text{atanh}(r) = \frac{1}{2} \log\left(\frac{1+r}{1-r}\right)$  称为相关系数  $r$  的Fisher变换。  
方差稳定化

基于Fisher变换的95%置信区间:

$$\left\{ \rho: \left| \sqrt{n} \left\{ \frac{1}{2} \log\left(\frac{1+r}{1-r}\right) - \frac{1}{2} \log\left(\frac{1+\rho}{1-\rho}\right) \right\} \right| \leq 1.96 \right\}$$

该区间构造依赖于下述事实:

4

性质2. 设  $(x_1, y_1), \dots, (x_n, y_n)$  iid, 设  $\rho = \text{cor}(x, y)$ , 则当  $n \rightarrow \infty$  时

$$\sqrt{n} \left\{ \frac{1}{2} \log \left( \frac{1+r}{1-r} \right) - \frac{1}{2} \log \left( \frac{1+\rho}{1-\rho} \right) \right\} \xrightarrow{d} N(0, 1)$$

证明: 根据性质1中r的极限分布和delta方法

注: 基于Fisher变换比直接基于r的渐近分布所构造的置信区间具有更好的精确度。

5

## 非参数相关性检验/置信区间

$$H_0: \rho = 0 \leftrightarrow H_1: \rho \neq 0$$

**(1) 置换检验(permutation test):** 不假设正态分布、且/或样本量较小的情形下, 可使用如下置换方法

For  $k = 1, 2, \dots, N$  ( $N$ : 置换次数)

(1) 置换  $x_1, x_2, \dots, x_n$ , 得序列  $x_{i_1}, \dots, x_{i_n}$

(2) 计算置换后的样本对  $(x_{i_1}, y_1), \dots, (x_{i_n}, y_n)$  的相关系数  $r_k$

$$p\text{值} = \frac{\#\{k: |r_k| \geq |r|\}}{N}$$

6

## (2). 自助法 (bootstrap) 构造置信区间

无联合正态分布假设下,  $\rho = \text{corr}(x, y)$  的95%置信区间可由自助法 (bootstrap) 构建:

从  $(x_1, y_1), \dots, (x_n, y_n)$  中有放回地抽取  $n$  对, 计算其样本相关系数  $r^*$ 。

重复  $N$  次, 得到  $N$  个  $r^*$ 。

计算  $N$  个  $r^*$  的 2.5% 和 97.5% 分位点, 分别作为置信区间的下界和上界。

7

## 3. Spearman's rho, Kendall's tau

### ■ Spearman's rho

样本:  $(x_1, y_1), \dots, (x_n, y_n)$ . 记  $R_i = x_i$  在所有  $x$  中的秩 (排名);  $S_i = y_i$  在所有  $y$  中的秩

Spearman's rho

$$\hat{\rho} = r_{RS} = \frac{\sum (R_i - \bar{R})(S_i - \bar{S})}{\sqrt{\sum (R_i - \bar{R})^2 \sum (S_i - \bar{S})^2}}$$

例:  $\mathbf{x} = (23, 67, 0.2, 99)$   
 $\text{Rank}(\mathbf{x}) = (2, 3, 1, 4)$

### ■ Kendall's tau

$$\hat{\tau} = \frac{C - D}{\binom{n}{2}}, \text{ 其中 } C = \#\{(i, j): (x_i - x_j)(y_i - y_j) > 0, i < j\}, D = \binom{n}{2} - C$$

统计推断: 置换或自助法

8

## 4. 随机向量（中文课本第二章）

1. 均值：设  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix}$  是随机向量，其均值定义为  $E(x) = \begin{pmatrix} E(x_1) \\ \vdots \\ E(x_k) \end{pmatrix}$

2. 协方差矩阵

设  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix}$ ,  $y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}$  是两个随机向量，记  $\mu_x = E(x)$ ,  $\mu_y = E(y)$ .

$x$ ,  $y$  的协方差矩阵 ( $k \times m$ ) 定义为：

$$\text{cov}(x, y) = E(x - \mu_x)(y - \mu_y)' = (\text{cov}(x_i, y_j))_{i=1, \dots, k; j=1, \dots, m}$$

3. 方差矩阵 (或方差 - 协方差矩阵)：

$$\text{var}(x) \text{ 或 } \text{cov}(x) = E(x - \mu)(x - \mu)' = \text{cov}(x, x)$$

9

性质：

1. 设  $x, y$  分别为  $n \times 1$  的随机向量， $z, w$  是  $n \times 1$  随机向量，则

$$(1) E(x + y) = E(x) + E(y)$$

$$(2) \text{cov}(x + y, z + w) = \text{cov}(x, z) + \text{cov}(x, w) + \text{cov}(y, z) + \text{cov}(y, w)$$

2. 设  $x$  为  $n \times 1$  随机向量， $A$  为  $m \times n$  常数矩阵，则

$$(1) E(Ax) = AE(x)$$

$$(2) \text{cov}(Ax) = A \text{cov}(x) A'$$

$$\begin{aligned} \text{证明：} (2) \text{cov}(Ax) &= E(Ax - AE(x))(Ax - AE(x))' \\ &= EA(x - E(x))(x - E(x))' A' = A[E(x - E(x))(x - E(x))'] A' \\ &= A \text{cov}(x) A' \end{aligned}$$

3. 设  $x, y$  分别为  $n \times 1$  和  $m \times 1$  随机向量， $A, B$  分别为  $p \times n, q \times n$  常数矩阵，则  $\text{cov}(Ax, By) = A \text{cov}(x, y) B'$

10

4. 设  $x$  为  $n \times 1$  随机向量，则方差 - 协方差矩阵  $\text{cov}(x) \geq 0$  (半正定).

证明：记  $\Sigma = \text{cov}(x)$ ，对于任何常数向量  $c \in R^n$ ,

由性质2(2)， $c' \Sigma c = c' \text{cov}(x) c = \text{cov}(c' x)$

注意  $\text{cov}(c' x)$  是随机变量  $c' x$  的方差，非负，所以  $c' \Sigma c \geq 0$ .

5. 设  $x$  为  $n \times 1$  随机向量， $E(x) = \mu$ ， $\text{cov}(x) = \Sigma$ ， $A$  为  $n \times n$  常数矩阵，则  $E(x' Ax) = \mu' A \mu + \text{tr}(A \Sigma)$

证明： $x' Ax = \text{tr}(x' Ax) = \text{tr}(Axx') \Rightarrow E(x' Ax) = \text{tr}(AE(xx'))$

而  $\Sigma = E(xx') - \mu \mu'$

$$\text{tr}(AE(xx')) = \text{tr}(A[\Sigma + \mu \mu']) = \text{tr}(A \Sigma) + \text{tr}(A \mu \mu') = \text{tr}(A \Sigma) + \mu' A \mu$$

11

## “不相关化”（“正交”化/对角化）

任意随机向量  $x_1, x_2$ ，记  $\Sigma = \text{cov} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$

令  $x_1^\perp = x_1 - \Sigma_{12} \Sigma_{22}^{-1} x_2$ ，则  $x_1^\perp$  与  $x_2$  不相关，且

$$\text{var}(x_1^\perp) = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \stackrel{\text{记为}}{=} \Sigma_{11 \cdot 2}$$

注：

$x_1^\perp = x_1$  中消除掉与  $x_2$  线性相关的部分  $\Sigma_{12} \Sigma_{22}^{-1} x_2$ ，

后者可以认为是  $x_1$  在  $x_2$  上的正交投影

12

总结一下:

$$\begin{pmatrix} \mathbf{x}_1^\perp \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{I} & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \Rightarrow \text{cov} \begin{pmatrix} \mathbf{x}_1^\perp \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} \Sigma_{11\bullet 2} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$$

$\Leftrightarrow$

方差矩阵对角化:

$$\begin{pmatrix} \mathbf{I} & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & \mathbf{I} \end{pmatrix} \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \mathbf{I} & 0 \\ -\Sigma_{22}^{-1}\Sigma_{21} & \mathbf{I} \end{pmatrix} = \begin{pmatrix} \Sigma_{11\bullet 2} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$$

13

引理1(分块矩阵的逆):  $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} > 0$  (正定)

$$\text{则 } \Sigma^{-1} = \begin{pmatrix} \Sigma_{11\bullet 2}^{-1} & -\Sigma_{11\bullet 2}^{-1}\Sigma_{12}\Sigma_{22}^{-1} \\ -\Sigma_{22\bullet 1}^{-1}\Sigma_{21}\Sigma_{11}^{-1} & \Sigma_{22\bullet 1}^{-1} \end{pmatrix}$$

其中  $\Sigma_{11\bullet 2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ ,  $\Sigma_{22\bullet 1} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$ ,

证明:由对角化公式( $\Sigma$ 为n阶,  $\Sigma_{11}$ 为k阶):

$$\begin{pmatrix} \mathbf{I}_k & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & \mathbf{I}_{n-k} \end{pmatrix} \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \mathbf{I}_k & 0 \\ -\Sigma_{22}^{-1}\Sigma_{21} & \mathbf{I}_{n-k} \end{pmatrix} = \begin{pmatrix} \Sigma_{11\bullet 2} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$$

$$\text{得 } \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_k & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & \mathbf{I}_{n-k} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{11\bullet 2} & 0 \\ 0 & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \mathbf{I}_k & 0 \\ -\Sigma_{22}^{-1}\Sigma_{21} & \mathbf{I}_{n-k} \end{pmatrix}^{-1}$$

14

两边求逆  $\Rightarrow$

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{I}_k & 0 \\ -\Sigma_{22}^{-1}\Sigma_{21} & \mathbf{I}_{n-k} \end{pmatrix} \begin{pmatrix} \Sigma_{11\bullet 2}^{-1} & 0 \\ 0 & \Sigma_{22}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{I}_k & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & \mathbf{I}_{n-k} \end{pmatrix}$$

$$= \begin{pmatrix} \Sigma_{11\bullet 2}^{-1} & -\Sigma_{11\bullet 2}^{-1}\Sigma_{12}\Sigma_{22}^{-1} \\ -\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11\bullet 2}^{-1} & \Sigma_{22}^{-1} + \Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11\bullet 2}^{-1}\Sigma_{12}\Sigma_{22}^{-1} \end{pmatrix}$$

注意: 由对称性知右下角的  $\Sigma_{22}^{-1} + \Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11\bullet 2}^{-1}\Sigma_{12}\Sigma_{22}^{-1}$

实际上等于  $\Sigma_{22\bullet 1}^{-1}$ , 同样左下角  $-\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11\bullet 2}^{-1} = -\Sigma_{22\bullet 1}^{-1}\Sigma_{21}\Sigma_{11}^{-1}$

15