

基于累积量（半不变量）的混合 RNA-seq 样本的协方差矩阵的假设检验量的构造

目录

■Part I 累积量（半不变量）理论回顾.....	1
■Part II 问题背景的回顾和转化.....	3
■Part III 回顾 CSnet Estimator 论文中一阶原点矩和二阶原点矩的估计	5
■Part IV 三阶和四阶高阶（混合）原点矩的估计思路	6
■Part V 三阶和四阶高阶（混合）原点矩估计的理论推导：	7
■Part VI: 协方差估计量的方差 $\text{Var}(\sigma_{jj}'(\mathbf{k}))$ 的 MSE 一致估计方法：	10
■Part VII: Max-entry 假设检验量的构造.....	11
■参考文献	11

■Part I 累积量（半不变量）理论回顾

一、 随机变量：

1、设随机变量 ξ_1 的 n 阶矩 m_n 存在，把它的特征函数 $\varphi(t)$ 按 Taylor 级数展开而得重要公式：

$$\varphi(t) = 1 + \sum_{k=1}^n \frac{m_k}{k!} (it)^k + o(t^n). \quad (1)$$

由特征函数性质知： $\log \varphi(t)$ 的前 n 级导数在 0 点存在，令

$$\chi_k = \frac{1}{i^k} \left[\frac{d^k}{dt^k} \log \varphi(t) \right]_{t=0} \quad (k \leq n), \quad (2)$$

称 χ_k 为随机变量 ξ 的 k 阶累积量（半不变量），显然有

$$\log \varphi(t) = \sum_{k=1}^n \frac{\chi_k}{k!} (it)^k + o(t^n). \quad (3)$$

通过比较（1）和（3）可得到累积量 χ_k 与矩 m_n 之间的关系，详见[1]¹.

2、累积量具有下列性质：

a. 如随机变量 ξ_1 、 ξ_2 独立，且它们的 k 阶累积量 $\chi_k^{(1)}$ ， $\chi_k^{(2)}$ 存在，则 $\xi = \xi_1 + \xi_2$

的 k 阶累积量 χ_k 为：

$$\chi_k = \chi_k^{(1)} + \chi_k^{(2)}. \quad (4)$$

b. 在变换 $\eta = \xi + b$ 下，累积量（除一阶外）不变， b 为常数.

二、随机向量：

1、设 $\xi = (\xi_1, \dots, \xi_k)$ 是随机向量，函数 $\varphi_\xi(\mathbf{t}) = \mathbf{E}e^{i(\mathbf{t}, \xi)}$ ， $\mathbf{t} = (t_1, \dots, t_k)$ ，是 ξ 的特征函数.

令 $\mathbf{E}|\xi_i|^n < \infty, i = 1, \dots, k$ ，则将 $\xi = (\xi_1, \dots, \xi_k)$ 展开为 Taylor 级数，得：

$$\varphi_\xi(t_1, \dots, t_k) = \sum_{v_1 + \dots + v_k \leq n} \frac{i^{v_1 + \dots + v_k}}{v_1! \dots v_k!} m_\xi^{(v_1, \dots, v_k)} t_1^{v_1} \dots t_k^{v_k} + o(|t|^n), \quad (5)$$

其中 $|t| = |t_1| + \dots + |t_k|$ ， $m_\xi^{(v_1, \dots, v_k)} = \mathbf{E}\xi_1^{v_1} \dots \xi_k^{v_k}$ 是 $v = (v_1, \dots, v_k)$ 阶（混合）矩.

由特征函数性质知， $\ln \varphi_\xi(t_1, \dots, t_k)$ 可以在 $(0, \dots, 0)$ 表示为 Taylor 公式：

$$\ln \varphi_\xi(t_1, \dots, t_k) = \sum_{v_1 + \dots + v_k \leq n} \frac{i^{v_1 + \dots + v_k}}{v_1! \dots v_k!} s_\xi^{(v_1, \dots, v_k)} t_1^{v_1} \dots t_k^{v_k} + o(|t|^n), \quad (6)$$

其中 $s_\xi^{(v_1, \dots, v_k)} = \frac{1}{i^{v_1 + \dots + v_k}} \left[\frac{\partial^{v_1 + \dots + v_k}}{\partial t_1^{v_1} \dots \partial t_k^{v_k}} \ln \varphi_\xi(t_1, \dots, t_k) \right]_{t_1 = \dots = t_k = 0}$.

我们称 $s_\xi^{(v_1, \dots, v_k)}$ 为随机向量 $\xi = (\xi_1, \dots, \xi_k)$ 的 $v = (v_1, \dots, v_k)$ 阶（混合）累积量（半不变量）.

2、累积量的性质：

a. 如随机向量 ξ 、 η 独立，且它们的 (v_1, \dots, v_k) 阶累积量 $s_\xi^{(v_1, \dots, v_k)}$ ， $s_\eta^{(v_1, \dots, v_k)}$ 存在，则 $\xi + \eta$ 的 (v_1, \dots, v_k) 阶累积量为：

¹ [3]王梓坤.概率论基础及其应用.第3版.北京师范大学出版社.2007

$$s_{\xi+\eta}^{(v_1, \dots, v_k)} = s_{\xi}^{(v_1, \dots, v_k)} + s_{\eta}^{(v_1, \dots, v_k)}. \quad (7)$$

b. 累积量与矩的关系（实际应用主要是计算公式（10））

为简化记号，对具有非负整数分量的向量 $v = (v_1, \dots, v_k)$ ，令

$$v! = v_1! \cdots v_k!, |v| = v_1 + \cdots + v_k, t^v = t_1^{v_1} \cdots t_k^{v_k}, s_{\xi}^{(v)} = s_{\xi}^{(v_1, \dots, v_k)}, m_{\xi}^{(v)} = m_{\xi}^{(v_1, \dots, v_k)}. \quad (8)$$

设 $\xi = (\xi_1, \dots, \xi_k)$ 是随机向量， $\mathbf{E}|\xi_i|^n < \infty, i = 1, \dots, k$. 则有：

$$m_{\xi}^{(v)} = m_{\xi}^{(v_1, \dots, v_k)} = \sum_{\lambda(1) + \cdots + \lambda(q) = v} \frac{v_1! \cdots v_k!}{\lambda(1)! \cdots \lambda(q)!} \prod_{p=1}^q s_{\xi}^{(\lambda(p))}, \quad (9)$$

$$s_{\xi}^{(v)} = s_{\xi}^{(v_1, \dots, v_k)} = \sum_{\lambda(1) + \cdots + \lambda(q) = v} (-1)^{q-1} (q-1)! \frac{v_1! \cdots v_k!}{\lambda(1)! \cdots \lambda(q)!} \prod_{p=1}^q m_{\xi}^{(\lambda(p))}, \quad (10)$$

其中 $\sum_{\lambda(1) + \cdots + \lambda(q) = v}$ 表示对于一切满足 $|\lambda(p)| > 0, \lambda(1) + \cdots + \lambda(q) = v$ 的”有序非负整数

分量的向量 $\lambda(p)$ ”求和。

更一般的形式，

$$m_{\xi}^{(v)} = m_{\xi}^{(v_1, \dots, v_k)} = \sum_{\{r_1 \lambda(1) + \cdots + r_x \lambda(x) = v\}} \frac{1}{r_1! \cdots r_x!} \frac{v_1! \cdots v_k!}{(\lambda(1)!)^{r_1} \cdots (\lambda(x)!)^{r_x}} \prod_{j=1}^x [s_{\xi}^{(\lambda(j))}]^{r_j}, \quad (11)$$

$$s_{\xi}^{(v)} = s_{\xi}^{(v_1, \dots, v_k)} = \sum_{\{r_1 \lambda(1) + \cdots + r_x \lambda(x) = v\}} \frac{(-1)^{q-1} (q-1)!}{r_1! \cdots r_x!} \frac{v_1! \cdots v_k!}{(\lambda(1)!)^{r_1} \cdots (\lambda(x)!)^{r_x}} \prod_{j=1}^x [m_{\xi}^{(\lambda(j))}]^{r_j}, \quad (12)$$

其中 $\sum_{\{r_1 \lambda(1) + \cdots + r_x \lambda(x) = v\}}$ 表示对于一切满足 $|\lambda(j)| > 0, r_1 \lambda(1) + \cdots + r_x \lambda(x) = v$ 的”有序正

整数 r_j 的数组”求和。

■ Part II 问题背景的回顾和转化

仍设 n 维混合 RNA-序列数据样本 across p genes: $x_1, \dots, x_n \in \mathbb{R}^p$. 并考虑 K 种细胞的混合采样模型：

$$\mathbf{x}_i = \sum_{k=1}^K \pi_{ik} \mathbf{x}_i^{(k)}, \quad (1)$$

其中 $\mathbf{x}_i^{(k)}$ ($k \in K$) 之间独立:

$$\mathbb{E}(\mathbf{x}_i) = \sum_{k=1}^K \pi_{ik} \boldsymbol{\mu}^{(k)}, \quad \text{Cov}(\mathbf{x}_i) = \sum_{k=1}^K \pi_{ik}^2 \boldsymbol{\Sigma}^{(k)}. \quad (2)$$

累积量观点: 我们换个观点看 (2), 其中 (2) 能够拆开的原因从累积量理论的观点看, 累积量满足对独立随机变量和的可加性, 且因为 1 阶单随机变量累积量恰好是数学期望, 2 阶单随机变量的累积量恰好是方差, (1, 1) 阶双随机变量的累积量恰好是协方差, 即如下三个定理:

Theorem 1.1 (将会涉及到的 4 阶以下累积量)

我们设 $s_{\xi}^{(v_1, \dots, v_k)}$ 为随机向量 $\xi = (\xi_1, \dots, \xi_k)$ 的 $v = (v_1, \dots, v_k)$ 阶 (混合) 累积量 (半不变量), 现考虑 $\xi = (X, Y)$, 则:

$$(1) \quad s_{\xi}^{(1,0)} = EX \quad (2) \quad s_{\xi}^{(0,1)} = EY$$

$$(3) \quad s_{\xi}^{(1,1)} = \text{Cov}(X, Y) \quad (4) \quad s_{\xi}^{(2,0)} = D(X), \quad s_{\xi}^{(0,2)} = D(Y)$$

$$(5) \quad s_X^{(3)} = E(X - EX)^3 \text{ (仍恰为 3 阶中心矩)} = m_3 - 3m_1m_2 + 2m_1^3$$

$$(6) \quad s_X^{(4)} = m_4 - 4m_1m_3 - 3m_2^2 + 8m_1^2m_2 - 6m_1^4 \text{ (不是 4 阶中心矩了)}$$

$$(7) \quad s_{X,Y}^{(2,1)} = EX^2Y + 4(EX)^2EY - EX^2EY - (EXY)EX$$

$$(8) \quad s_{X,Y}^{(2,2)} = EX^2Y^2 - 2EXY^2EX - 2EX^2Y EY -$$

$$EX^2EY^2 - 4(EXY)^2 + 2(EX)^2EY^2 + 2EX^2(EY)^2 - 24(EX)^2(EY)^2$$

Theorem 1.2 (累积量的独立可加性)

如随机向量 ξ 、 η 独立, 且它们的 (v_1, \dots, v_k) 阶累积量 $s_{\xi}^{(v_1, \dots, v_k)}$, $s_{\eta}^{(v_1, \dots, v_k)}$ 存在,

则 $\xi + \eta$ 的 (v_1, \dots, v_k) 阶累积量为:

$$s_{\xi+\eta}^{(v_1, \dots, v_k)} = s_{\xi}^{(v_1, \dots, v_k)} + s_{\eta}^{(v_1, \dots, v_k)}.$$

Theorem 1.3 (累积量的常数齐次性)

对随机向量 ξ ，且它的 (v_1, \dots, v_k) 阶累积量 $s_\xi^{(v_1, \dots, v_k)}$ 存在，则 $c\xi$ 的 (v_1, \dots, v_k) 阶累积量为

$$c^{\sum v_i} s_\xi^{(v_1, \dots, v_k)}$$

(注：由于 K 种细胞采样的混合比例 $(\pi_{i1}, \pi_{i2}, \dots, \pi_{iK})$ 不同，不同实验所得的 $x_i \in \mathbb{R}^p$ 间不一定 i.i.d，仅仅满足独立)

然而，不幸的是，当累积量的阶高于 3 阶时，唯有单个随机变量的 3 阶中心矩仍保持恰好是一个累积量。

为简化表达，不失一般性，当 $p > 1$ 时，直接用不同的字母 X, Y 表示这 p genes 中的任意某两个不同的基因表达数据 $X = (X_1, X_2, \dots, X_n)$, $Y = (Y_1, Y_2, \dots, Y_n)$ 。

即用 $X_1, X_2, \dots, X_n \in \mathbb{R}$; $Y_1, Y_2, \dots, Y_n \in \mathbb{R}$ 分别表示 X 基因数据分量和 Y 基因数据分量所分别对应采集到的 n 个混合数据。

■问题的形式总结：

依据混合原点矩的线性和不同实验次数间的独立性两条性质，对 $\text{Var}(\bar{\sigma}_{jj'}^{(k)})$ 可以进行直接的拆分计算，对它的估计可以转化为以下两种情形：

(1) 当 $j \neq j'$ 时，不妨设 j 对应的基因为 X , j' 对应的基因为 Y

归结为对广义二次型方差 $\text{Var}(X^T A Y)$ 的估计，最终归结为对 4 阶以下混合原点矩 $E(X_i X_j Y_k Y_l)$ 的估计和讨论，（其中脚标 $i, j, k, l \in [n]$ 注意到脚标不同时，即不同次实验的数据），可以根据不同次实验间的独立性拆开。

(2) 当 $j = j'$ 时，不妨设 j, j' 对应的基因都为 X ：

归结为对二次型方差 $\text{Var}(X^T A X)$ 的估计，最终归结为对 4 阶以下混合原点矩 $E(X_i X_j X_k X_l)$ 的估计和讨论。

■最终转化为要去估计的各阶（混合）原点中心矩包括：

- (I) 所有可能的一阶原点矩： $E(X_s)$ ($s \in [n]$); $E(Y_s)$ ($s \in [n]$)
- (II) 所有可能的二阶(混合)原点矩： $E(X_s^2)$ ($s \in [n]$); $E(Y_s^2)$ ($s \in [n]$)
 $E(X_s Y_s)$ ($s \in [n]$)
- (III) 所有涉及到的三阶(混合)原点矩： $E(X_s^2 Y_s)$ ($s \in [n]$); $E(X_s Y_s^2)$ ($s \in [n]$); $E(X_s^3)$ ($s \in [n]$); $E(Y_s^3)$ ($s \in [n]$)
- (IV) 所有涉及到的四阶(混合)原点矩： $E(X_s^2 Y_s^2)$ ($s \in [n]$); $E(X_s^4)$ ($s \in [n]$);

■Part III 回顾 CSnet Estimator 论文中一阶原点矩和二阶原点矩的估计

■(I) 一阶原点矩的估计：

其中 (I), (II) 的估计已在 CSnet Estimator [\[1.SZZ21\]](#) 论文中完成，对 (I) 中的一阶原点矩，回顾：

Denote $\mathbf{y}_j = (x_{1j}, \dots, x_{nj})$ and $\mathbf{D} = (\pi_{ik})_{n \times K}$. Equation (3) entails estimation of the cell-type-specific mean $\boldsymbol{\mu}^{(k)}$ via

$$\hat{\mu}_j^{(k)} = \left[(\mathbf{D}^\top \mathbf{D})^{-1} \mathbf{D}^\top \mathbf{y}_j \right]_k, \quad j \in [p], \quad k \in [K], \quad (5)$$

$$\mathbb{E}(\mathbf{x}_i) = \sum_{k=1}^K \pi_{ik} \boldsymbol{\mu}^{(k)}, \quad \text{Cov}(\mathbf{x}_i) = \sum_{k=1}^K \pi_{ik}^2 \boldsymbol{\Sigma}^{(k)}. \quad (2)$$

设 $\hat{\mu}_{ij}$ 为参数 $\mathbb{E}(x_{ij})$ 的估计，则由 (5) 和 (2)，直接有估计：

$$\hat{\mu}_{ij} = \sum_{k=1}^K \pi_{ik} \hat{\mu}_j^{(k)} \quad (\text{A})$$

■ (II) 二阶(混合)原点矩的估计：

首先回顾：

where $[\mathbf{x}]_k$ is the k th entry in $\mathbf{x} \in \mathbb{R}^K$. Let $\hat{z}_{ij} = x_{ij} - \sum_{k=1}^K \pi_{ik} \hat{\mu}_j^{(k)}$ and $\hat{\mathbf{z}}_j = (\hat{z}_{1j}, \dots, \hat{z}_{nj}) \in \mathbb{R}^n$. Denoting $\mathbf{H} = (\pi_{ik}^2)_{n \times K}$, equation (4) entails estimation of the cell-type-specific covariance $\sigma_{jj'}^{(k)}$ via

$$\hat{\sigma}_{jj'}^{(k)} = \left[(\mathbf{H}^\top \mathbf{H})^{-1} \mathbf{H}^\top (\hat{\mathbf{z}}_j \circ \hat{\mathbf{z}}_{j'}) \right]_k, \quad j, j' \in [p], \quad k \in [K], \quad (6)$$

首先设 $\hat{\sigma}_{ijj'}$ 为参数 $\mathbb{E}(x_{ij} - \mathbb{E}x_{ij})(x_{ij'} - \mathbb{E}x_{ij'})$, 即 $\text{Cov}(x_{ij}, x_{ij'})$ 的估计量，则有：

$$\hat{\sigma}_{ijj'} = \sum_{k=1}^K \pi_{ik}^2 \hat{\sigma}_{jj'}^{(k)} \quad (\text{B})$$

再设 $\hat{\beta}_{ijj'}$ 为参数 $\mathbb{E}x_{ij}x_{ij'}$ 的估计。由恒等式 $\text{Cov}(X, Y) = \mathbb{E}XY - \mathbb{E}X\mathbb{E}Y$ 和 (A)，则有：

$$\hat{\beta}_{ijj'} = \hat{\sigma}_{ijj'} + \hat{\mu}_{ij}\hat{\mu}_{ij'} \quad (\text{C})$$

■ Part IV 三阶和四阶高阶（混合）原点矩的估计思路

■ (III) 和 (IV) 中所涉及高阶混合矩的估计：

(I)，(II) 的估计能够通过 (2) 和对应的回归方程 (3)，(4) 得到，虽然是一阶矩和二阶矩的情况，但可以看成之前提过的 **Theorem1.1** 和 **Theorem1.2**，以及 **Theorem 1.3** 的特殊情形，这里 $z_{ij} = x_{ij} - \mathbb{E}(x_{ij})$ 。

$$\mathbb{E}(\mathbf{x}_i) = \sum_{k=1}^K \pi_{ik} \boldsymbol{\mu}^{(k)}, \quad \text{Cov}(\mathbf{x}_i) = \sum_{k=1}^K \pi_{ik}^2 \boldsymbol{\Sigma}^{(k)}. \quad (2)$$

Letting $\boldsymbol{\mu}^{(k)} = (\mu_1^{(k)}, \dots, \mu_p^{(k)})$ and $\boldsymbol{\Sigma}^{(k)} = (\sigma_{jj'}^{(k)})_{p \times p}$, (1) and (2) together imply

$$x_{ij} = \sum_{k=1}^K \pi_{ik} \mu_j^{(k)} + z_{ij}, \quad j \in [p], \quad (3)$$

$$z_{ij} z_{ij'} = \sum_{k=1}^K \pi_{ik}^2 \sigma_{jj'}^{(k)} + \epsilon_{ijj'}, \quad j, j' \in [p], \quad (4)$$

where $\mathbb{E}(z_{ij}) = 0$ and $\mathbb{E}(\epsilon_{ijj'}) = 0$. This formulation facilitates an efficient least squares

所以一个自然的想法是：

想估计 (III), (IV) 中所涉及的更高阶的（混合）原点矩的估计，可以考虑先用同样的回归方法估计相应的整体的高阶累积量，然后代入之前已估计出来的低阶的（混合）原点矩，最终得到，所想要的高阶（混合）原点矩的估计。

■Part V 三阶和四阶原点矩估计的理论推导：

本节中，我们将以 $\mathbf{E}(\mathbf{x}_{ij}^3)$ ($i \in [n], j \in [p]$) 和 $\mathbf{E}(\mathbf{x}_{ij}^4)$ ($i \in [n], j \in [p]$) 为例子，推导它们的 MSE 一致估计。

首先进行随机变量 X 的 (3) 阶, (4) 阶累积量 $s_{X,Y}^{(3,0)}, s_{X,Y}^{(4,0)}$, 简写为 $s_X^{(3)}, s_X^{(4)}$ 的计算，始终假设 $m_i = \mathbf{E}(X^i)$, 且相应累积量存在，则：

$$(1) \quad s_X^{(3)} = \mathbf{E}(X - \mathbf{E}X)^3 \text{ (仍恰为 3 阶中心矩)} = m_3 - 3m_1 m_2 + 2m_1^3$$

$$(2) \quad s_X^{(4)} = m_4 - 4m_1 m_3 - 3m_2^2 + 8m_1^2 m_2 - 6m_1^4 \text{ (不是 4 阶中心矩了)}$$

注意：

$$s_X^{(3)} = m_3 - 3m_1 m_2 + 2m_1^3 = \mathbf{E}(X^3 - 3m_1 m_2 + 2m_1^3)$$

$$\begin{aligned} s_X^{(4)} &= m_4 - 4m_1 m_3 - 3m_2^2 + 8m_1^2 m_2 - 6m_1^4 \\ &= \mathbf{E}X^4 - 4m_1 m_3 - 3m_2^2 + 8m_1^2 m_2 - 6m_1^4 \\ &= \mathbf{E}(X^4 - 4m_1 m_3 - 3m_2^2 + 8m_1^2 m_2 - 6m_1^4) \end{aligned}$$

其中 $(X^3 - 3m_1 m_2 + 2m_1^3)$ 以及 $(X^4 - 4m_1 m_3 - 3m_2^2 + 8m_1^2 m_2 - 6m_1^4)$ 地位上

对等于如下回归方程（3）和（4）中的 $z_{ij}z_{ij'} = (x_{ij} - E(x_{ij}))(x_{ij'} - E(x_{ij'}))$ 部分.

Letting $\mu^{(k)} = (\mu_1^{(k)}, \dots, \mu_p^{(k)})$ and $\Sigma^{(k)} = (\sigma_{jj'}^{(k)})_{p \times p}$, (1) and (2) together imply

$$x_{ij} = \sum_{k=1}^K \pi_{ik} \mu_j^{(k)} + z_{ij}, \quad j \in [p], \quad (3)$$

$$z_{ij}z_{ij'} = \sum_{k=1}^K \pi_{ik}^2 \sigma_{jj'}^{(k)} + \epsilon_{ijj'}, \quad j, j' \in [p], \quad (4)$$

where $E(z_{ij}) = 0$ and $E(\epsilon_{ijj'}) = 0$. This formulation facilitates an efficient least squares

下面我们以 $E(x_{ij}^3)$ ($i \in [n], j \in [p]$) 和 $E(x_{ij}^4)$ ($i \in [n], j \in [p]$) 的估计为例推导其基于高阶累积量的回归估计, 类似 (3), (4) 的回归顺序, 要想先估计 4 阶矩, 先要先估计 3 阶矩。

首先设 $w_{X_{ij}}^{(3)} = x_{ij}^3 - 3m_{ij1}m_{ij2} + 2m_{ij1}^3$, ($m_{ijk} = E(x_{ij}^k)$)

而注意到在 (3), (4) 回归方程的估计中已经得到 m_{ij1} , m_{ij2} 的估计, 即一阶原点矩和二阶（混合）原点矩的估计:

$$\hat{\mu}_{ij} = \sum_{k=1}^K \pi_{ik} \hat{\mu}_j^{(k)} \dots (A)$$

$$\hat{\beta}_{ijj'} = \hat{\sigma}_{ijj'} + \hat{\mu}_{ij} \hat{\mu}_{ij'} \dots (C)$$

这里重新记估计量 $\hat{\mu}_{ij}$ 为符号 \hat{m}_{ij1} , $\hat{\beta}_{ijj'}$ 为符号 \hat{m}_{ij2} , 再设

$$h_{X_{ij}}^{(4)} = x_{ij}^4 - 4m_{ij1}m_{ij3} - 3m_{ij2}^2 + 8m_{ij1}^2m_{ij2} - 6m_{ij1}^4 (m_{ijk} = E(X_{ij}^k))$$

$$\text{则} \quad w_{X_{ij}}^{(3)} = \sum_{k=1}^K \pi_{ik}^3 s_{X_{ij}}^{(k)(3)} + \epsilon_{ij}, \quad j \in [p] \quad (5)$$

$$h_{X_{ij}}^{(4)} = \sum_{k=1}^K \pi_{ik}^4 s_{X_{ij}}^{(k)(4)} + \epsilon'_{ij}, \quad j \in [p] \quad (6)$$

其中 $E(\epsilon_{ij}) = 0$, $E(\epsilon'_{ij}) = 0$.

记 $\hat{w}_{X_{ij}}^{(3)} = x_{ij}^3 - 3\hat{m}_{ij1}\hat{m}_{ij2} + 2\hat{m}_{ij1}^3$, 类似有 $\hat{h}_{X_{ij}}^{(4)}$, 但注意 $\hat{h}_{X_{ij}}^{(4)}$ 中的 \hat{m}_{ij3} 需要在 (5) 做完后才有相应的估计量, 就如 (3) 和 (4) 之间的关系, 先有期望才能估计方差。

设 $\hat{w}_{X_j}^{(3)} = (\hat{w}_{X_{1j}}^{(3)}, \hat{w}_{X_{2j}}^{(3)}, \dots, \hat{w}_{X_{nj}}^{(3)}) \quad (n \times 1)$

$$\hat{h}_{X_j}^{(4)} = (\hat{h}_{X_{1j}}^{(4)}, \hat{h}_{X_{2j}}^{(4)}, \dots, \hat{h}_{X_{nj}}^{(4)}) \quad (n \times 1)$$

于是：

设 $W=(\pi_{ik}^3)n \times K$; $V=(\pi_{ik}^4)n \times K$.

则 $s_{Xj}^{(k)(3)}$ 的估计为：

$$\hat{s}_{Xj}^{(k)(3)}=[(W^TW)^{-1}W^T\hat{w}_{Xj}^{(3)}]_k \quad j \in [p], k \in [K]$$

从而：

$$\hat{s}_{Xij}^{(3)} = \sum_{k=1}^K \pi_{ik}^3 \hat{s}_{Xj}^{(k)(3)}$$

Theorem 5.1

(1) 最后首先得到 $E(X_{ij}^3)$ 的一个 MSE 一致估计为：

$$\hat{m}_{ij3} = \hat{s}_{Xij}^{(3)} + 3\hat{m}_{ij1}\hat{m}_{ij2} - 2\hat{m}_{ij1}^3$$

同理可得到：

$$\hat{s}_{Xj}^{(k)(4)}=[(V^TV)^{-1}V^T\hat{h}_{Xj}^{(4)}]_k \quad j \in [p], k \in [K]$$

从而：

$$\hat{s}_{Xij}^{(4)} = \sum_{k=1}^K \pi_{ik}^4 \hat{s}_{Xj}^{(k)(4)}$$

(2) 进而得到 $E(X_{ij}^4)$ 的一个 MSE 一致估计为：

$$\hat{m}_{ij4} = \hat{s}_{Xij}^{(4)} + 4\hat{m}_{ij1}\hat{m}_{ij3} - 3\hat{m}_{ij2}^2 - 8\hat{m}_{ij1}^2\hat{m}_{ij2} + 6\hat{m}_{ij1}^4$$

对于 (III), (IV) 中其他的混合 3 阶原点矩, 4 阶原点矩也可类似的给出回归估计。

最后再计算两个要使用的高阶累积量作为本节的结束：

Example 5.1

随机变量 X, Y 的 $s_{X,Y}^{(2,1)}$, $s_{X,Y}^{(2,2)}$ 的计算, 设 $m_i = E(X^i)$, 若相应累积量存在, 则：

$$(1) \quad s_{X,Y}^{(2,1)} = EX^2Y + 4(EX)^2EY - EX^2EY - (EXY)EX$$

$$(2) \quad s_{X,Y}^{(2,2)} = EX^2Y^2 - 2EXY^2EX - 2EX^2Y EY -$$

$$EX^2EY^2 - 4(EXY)^2 + 2(EX)^2EY^2 + 2EX^2(EY)^2 - 24(EX)^2(EY)^2$$

■Part VI: 协方差估计量的方差 $\text{Var}(\bar{\sigma}_{jj'}^{(k)})$ 的 MSE 一致估计方法:

固定 k , 对任意第 k 类细胞, 仍用 X, Y 表示 p 个基因中的某两个基因, 类似于 $E(x_{ij}^3)$ 和 $E(x_{ij}^4)$ ($i \in [n], j \in [p]$) 情形的估计, 原始论文中通过一阶累积量数学期望 $s_{X,Y}^{(1,0)(k)} = EX$ 和 $s_{X,Y}^{(0,1)(k)} = EY$ 给出了 EX 和 EY 的估计, 通过二阶累积量协方差 $s_{X,Y}^{(1,1)(k)} = \text{Cov}(X, Y) = E(X - EX)(Y - EY)$ 给出了 EXY 的估计。方才我们用三阶累积量 $s_{X,Y}^{(3,0)(k)} = E(X - EX)^3$ (仍恰为 3 阶中心矩) $= m_3 - 3m_1m_2 + 2m_1^3$ 给出了 EX^3 的估计, 四阶累积量 $s_{X,Y}^{(4,0)(k)} = m_4 - 4m_1m_3 - 3m_2^2 + 8m_1^2m_2 - 6m_1^4$ (不是 4 阶中心矩) 给出了 EX^4 的估计

我们可以依次通过对 $s_{X,Y}^{(2,1)(k)}, s_{X,Y}^{(2,2)(k)}$, 以及对称的 $s_{X,Y}^{(1,2)(k)}$ 的回归方程 ($k \in [K]$) 去估计相应的累积量, 再代入前面已经被估计出的低阶累积量, 最余给出其他所需要的混合中心矩 $E(X_s^2 Y_s)$ ($s \in [n]$); $E(X_s Y_s^2)$ ($s \in [n]$) $E(X_s^2 Y_s^2)$ ($s \in [n]$) 的估计, 用 $\hat{E}(X_s^m Y_s^l)$ ($m, l = 0, 1, 2 \dots$) 表示相应的估计量, 从而通过这些混合局估计可以如下构造出一种 $\text{Var}(\bar{\sigma}_{jj'}^{(k)})$ 的 MSE 一致估计:

又由于当 $j \neq j'$ 时, 若 $\text{Var}(\bar{\sigma}_{jj'}^{(k)}) = F_{jj'}(E(X_{ij}^1 X_{ij'}^2), E(X_{ij}^2 X_{ij'}^2), E(X_{ij}^2 X_{ij'}^1), E(X_{ij}^1), E(X_{ij'}^2), E(X_{ij'}^1), E(X_{ij'}^2), E(X_{ij} X_{ij'}), i \in [n])$ 其中 $F_{jj'}$ 指的是将标准差估计量 $\bar{\sigma}_{jj'}^{(k)}$ 的方差直接依据混合矩的线性和不同实验间的独立性将其直接计算拆开, 表示为以关于 X_{ij} 和 $X_{ij'}$ ($i \in [n]$) 高阶混合中心矩为自变量的函数。

当 $j = j'$ 时, 若 $\text{Var}(\bar{\sigma}_{jj}^{(k)}) = \text{Var}(\bar{\sigma}_{jj}^{(k)}) = F_{jj}(E(X_{ij}^2), E(X_{ij}^4), E(X_{ij}^1), E(X_{ij}^3), i \in [n])$ 其中 Function 指的是将标准差估计量 $\bar{\sigma}_{jj}^{(k)}$ 的方差直接依据混合矩的线性和不同实验间的独立性将其直接计算拆开, 表示为以关于 X_{ij} ($i \in [n]$) 的高阶中心矩为自变量的函数。

Theorem 6.1

若用 $\hat{E}(X_{ij}^s X_{ij'}^t)$ 表示 $E(X_{ij}^s X_{ij'}^t)$ 的用累积量回归所得的 MSE 估计, 用 $\widehat{\text{Var}}(\bar{\sigma}_{jj'}^{(k)})$ 表

示 $\text{Var}(\bar{\sigma}_{jj'}^{(k)})$ 的一个一致估计量，则：

$$\begin{aligned} \widehat{\text{Var}}(\bar{\sigma}_{jj'}^{(k)}) &= \begin{cases} F_{jj'} \left(\widehat{E}(X_{ij}^1 X_{ij'}^2), \widehat{E}(X_{ij}^2 X_{ij'}^2), \widehat{E}(X_{ij}^2 X_{ij'}^1), \widehat{E}(X_{ij}^1), \widehat{E}(X_{ij'}^2), \widehat{E}(X_{ij}^2), \widehat{E}(X_{ij'}^1), \widehat{E}(X_{ij}^1 X_{ij'}^1), i \in [n] \right) & \text{if } j \neq j' \\ F_{jj} \left(\widehat{E}(X_{ij}^2), \widehat{E}(X_{ij}^4), \widehat{E}(X_{ij}^1), \widehat{E}(X_{ij}^3), i \in [n] \right) & \text{if } j = j' \end{cases} \end{aligned}$$

而对于 $\widehat{\text{Var}}(\bar{\sigma}_{jj'}^{(k)})$ ， $F_{jj'}$ 是一些 MSE 估计的直接加和，因为每一项加和的单项是至多两个混合矩的 MSE 的估计的乘积所以 $\widehat{\text{Var}}(\bar{\sigma}_{jj'}^{(k)})$ 仍是一个 MSE 一致估计。

■Part VII: Max-entry 假设检验量的构造

Theorem 7.1

从而第 k 类细胞，我们考虑两组 bulk RNA 数据所得的协方差矩阵的假设检验 $H_0: \Sigma_1^{(k)} = \Sigma_2^{(k)}$ ，可以用如下 Max-entry Testing 进行检验,定义

$$M_{jj'}^{(k)} := \frac{(\bar{\sigma}_{jj'(1)}^{(k)} - \bar{\sigma}_{jj'(2)}^{(k)})^2}{\widehat{\text{Var}}(\bar{\sigma}_{jj'(1)}^{(k)}) + \widehat{\text{Var}}(\bar{\sigma}_{jj'(2)}^{(k)})} \quad (1 \leq j \leq j' \leq p)$$

则我们考虑如下 Max-entry 检验量：

$$M_n^{(k)} := \max_{1 \leq j \leq j' \leq p} M_{jj'}^{(k)} = \max_{1 \leq j \leq j' \leq p} \frac{(\bar{\sigma}_{jj'(1)}^{(k)} - \bar{\sigma}_{jj'(2)}^{(k)})^2}{\widehat{\text{Var}}(\bar{\sigma}_{jj'(1)}^{(k)}) + \widehat{\text{Var}}(\bar{\sigma}_{jj'(2)}^{(k)})}$$

从[2.JASA13]的工作知道它是一个呈现 cumulative 分布的假设检验量。

■参考文献

1. Estimating cell-type-specific gene co-expression networks from bulk gene expression data with an application to Alzheimer's disease, Chang Su, Jingfei Zhang, Hongyu Zhao, [doi: https://doi.org/10.1101/2021.12.21.473558](https://doi.org/10.1101/2021.12.21.473558)
- 2.Cai, Tony, et al. "Two-Sample Covariance Matrix Testing and Support Recovery in High-Dimensional and Sparse Settings." Journal of the American Statistical Association, vol. 108, no. 501, 2013, pp. 265–77. JSTOR, <http://www.jstor.org/stable/23427527>. Accessed 25 Sep. 2022.
- 3.王梓坤.概率论基础及其应用.第3版.北京师范大学出版社.2007