



Uncorrelated Screening-based FDR control procedure

1. 改进目标: ① 本文是针对均值检验提出的, 主要改进是利用假设已知的稀疏信息, 包含入统计量中, 以此获得检验更高的 Power. 由稀疏假设, 我们更多关注两内值的支撑 S_1, U_{S_2} . 故思想为第一步筛选 (Screen) S_1, U_{S_2} , 第二步再进行 FDR 检验

② 在以往的 screening 方法, 要求 screening 步骤和检验步骤独立, 本文提出的方法是 screen 统计量与 test 统计量两两进不相关, 以前保证独立通常还需要额外的数据集或对样本进行分割, (这可能会导致 information loss) 且十分依赖于分割方法

③ 本文还论证了 $power_{US} \geq power_{BH}$ 的结果, 还对信号大小 (H_1) 进行了分析证明了 US 方法 power 收敛到 1 的信号范围宽于 BH 方法. NOTE: US 方法可以没有稀疏假设, 此时至少也和 BH 方法一样 powerful. 但稀疏情形可以更优.

2. 方法步骤 (以 $\delta_{i1}^2 \neq \delta_{i2}^2$ 为例)

Step 1: 计算 $T_i = \frac{\bar{X}_{i1} - \bar{X}_{i2}}{\sqrt{\frac{\delta_{i1}^2}{n_1} + \frac{\delta_{i2}^2}{n_2}}}$, $\hat{\delta}_{iy} = \frac{1}{n_y - 1} \sum_{j=1}^{n_y} (X_{ij,y} - \bar{X}_{i,y})^2$, $y=1,2$ 和 $S_i = \sqrt{\frac{n_1}{\delta_{i1}^2 (1 + \frac{n_2 \delta_{i1}^2}{n_1 \delta_{i2}^2})}} (\bar{X}_{i1} + \frac{n_2 \delta_{i1}^2}{n_1 \delta_{i2}^2} \bar{X}_{i2})$

NOTE: 用 S_i 作为 screening 统计量, S_i 较小较大分成两组. 一般情况下, S_1, U_{S_2} 在小组中. 大组中为 (S_1, U_{S_2}) 的元素, 但可能小组中有 U_{S_1}, U_{S_2} 异号的, 这没关系. 分成两组后, 异号的和 S_1, U_{S_2} 在一起更易检验出来. 记 $\hat{R}_{i,\lambda} = \sum_{j=1}^m I(|S_{ij}| > \lambda)$, $\hat{m}_{i,\lambda} = \sum_{j=1}^m I(|S_{ij}| > \lambda)$, $y=1,2$.

Step 2: 参数 λ 的选取 $\hat{\lambda} = \frac{1}{\sqrt{m}} \log m$, $\hat{t} = \arg \max_{1 \leq k \leq 4N} |R_{k,\lambda}|$

NOTE: 相当于用网格法寻找最合适的 λ . 计算 $|R_{k,\lambda}|$ 时要用到第三步的阈值.

Step 3: 计算 threshold $\hat{t}_{1,\lambda} = \inf \{ t \geq 0, \frac{\hat{m}_{1,\lambda} (2 - \Phi(t))}{\max(1, \sum_{i=1}^m I(|S_{i1}| > \lambda, |T_{i1}| \geq t))} \leq \alpha \}$
 $\hat{t}_{2,\lambda} = \inf \{ t \geq 0, \frac{\hat{m}_{2,\lambda} (2 - \Phi(t))}{\max(1, \sum_{i=1}^m I(|S_{i2}| < \lambda, |T_{i2}| \geq t))} \leq \alpha \}$

NOTE: $R_{\lambda} = \{1 \leq i \leq m, |S_{i1}| > \lambda, |T_{i1}| \geq \hat{t}_{1,\lambda}\} \cup \{1 \leq i \leq m, |S_{i2}| < \lambda, |T_{i2}| \geq \hat{t}_{2,\lambda}\}$, 相当于筛选后两组各自使用 BH 方法检验.

Step 4: 对给定的 $\alpha \in (0,1)$, 当 $i \in R_{\lambda}$ 时, 拒绝 H_{i0} , $i=1, \dots, m$.

NOTE: 调参步中只在 $0 < \lambda < 4\sqrt{\log m}$ 考虑是因为 $P(|S_{i1}| > 4\sqrt{\log m}) \rightarrow 0$.

3. 假设条件. (1) $|H_1| = \alpha m$, 即 u_1, u_2 的支撑是稀疏的.

(2) $E[\exp(t|X_{iy} - u_{iy}|/\delta_{iy})] \leq K_1$, 对某 $K_1 > 0, t > 0, \forall 1 \leq i \leq m, y=1,2$. 设 $c_1 \leq \frac{n_1}{n_2} \leq c_2$, 其中 $c_1, c_2 > 0$. 以及 $\min(n_1, n_2) \geq c(\log m)^{\frac{1}{\gamma}}$, $\gamma > 5, c > 0$. 该条件相当于对总体的分布条件, 且此在高维情形中常见.

(3). 设 u_i 为 \mathcal{H}_0 的子集, 满足 $\{X_{g,1}, X_{g,2}, y \in u_i\}$ 与 $(X_{i,1}, X_{i,2})$ 相互独立. 对 $\forall i \in \mathcal{H}_0$. 有 $|u_i| \geq m_0 - K$. 即对任意 $(X_{i,1}, X_{i,2})$, 我们允许 K 个变量强烈相关

筛选
↓
检验
↓
调参
↓
决策



四. 理论结果

定理一: 设有 (1), (2), 且 $|R\lambda| \xrightarrow{P} \infty, m \rightarrow \infty, \forall \varepsilon > 0, P(FDP \leq \alpha + \varepsilon) \rightarrow 1$. 因此有 $\limsup_{m \rightarrow \infty} FDR \leq \alpha$

定理二: 设 $m_1 \rightarrow \infty$, (1)-(3) 成立, 则有 $\text{power}_{BS} \geq \text{power}_{BH} + o(1)$, 对某些 $O(1), m \rightarrow \infty$.

接下来探讨 power, 假设 $\frac{|u_{1,1} - u_{1,2}|}{\sqrt{\sigma_{1,1}^2/m_1 + \sigma_{1,2}^2/m_2}} = \theta \sqrt{\log m}, 1 \in H_1, \theta > 0$. θ 衡量了信号 (H_1) 元素的大小 (强弱). 设 $|H_1| = m^\beta, 0 < \beta < 1$ 衡量信号的多少. (*)

定理三: 设 (2), (3) 成立. 若 $0 < \theta < \sqrt{2(1-\beta)}$, 有 $\text{power}_{BH} \xrightarrow{P} 0$. 若 $\theta > \sqrt{2(1-\beta)}$, 则 $\text{power}_{BH} \xrightarrow{P} 1$.

设两均值向量满足: 对某个 $0 \leq r \leq 1, 0 < h \leq 2, 0 < \rho < 1, K > h + \rho\sqrt{1-\beta}$

(4) $\text{Card}\{i \in H_0: \sqrt{\frac{n_1}{\sigma_{1,1}^2(1 + \frac{n_2 \sigma_{1,1}^2}{n_1 \sigma_{1,2}^2})}} |u_{1,1} + \frac{n_2 \sigma_{1,1}^2}{n_1 \sigma_{1,2}^2} u_{1,2}| \geq h \sqrt{\log m}\} = O(m^r)$ (u_1 和 u_2 的稀疏条件)

(5) $\text{Card}\{i \in H_1: \sqrt{\frac{n_1}{\sigma_{1,1}^2(1 + \frac{n_2 \sigma_{1,1}^2}{n_1 \sigma_{1,2}^2})}} |u_{1,1} + \frac{n_2 \sigma_{1,1}^2}{n_1 \sigma_{1,2}^2} u_{1,2}| \geq K \sqrt{\log m}\} \geq \rho |H_1|$ (确保 S_n 可以筛选)

定理四: 设 (2), (3), (*) 成立.

① 若 $\theta > \sqrt{2(1-\beta)}$, $\text{power}_{BS} \xrightarrow{P} 1, m \rightarrow \infty$

② (4), (5) 成立, 令 $\theta > \sqrt{\max(0, 2r-2\beta)}, N > \frac{10}{m \min(1-\beta, \frac{\theta^2}{4})}$, 有 $P(\text{power}_{BS} \geq \rho - \varepsilon) \rightarrow 1$

$\forall \varepsilon > 0, m \rightarrow \infty$.

③ 有 $P(FDP \leq \alpha + \varepsilon) \rightarrow 1, \forall \varepsilon > 0, m \rightarrow \infty$.