



Grouping, adjusting and Pooling (GAP) Procedure.

1. 改进动机: ① 本文是提出一种利用稀疏信息去进行两样本的多重检验. 其中第一步为构建协变量序列, 第二步为利用辅助协变量进行推断 (三步算法, grouping, adjusting, pooling) 其中本文 GAP 的显著优点是它能够处理各种相关结构 (高维情形)

② 对于两样本检验, $\beta_1 = \beta_2$, 传统方法为构建统计量 $\{T_1, \dots, T_m\}$, 选取阈值去控制 FDR. 这种方法忽视了 β_1, β_2 本身稀疏的信息, 导致了信息上的损失. 而类似的改进方法有 CARS procedure. 其通过构建辅助协变量 $\{S_1, \dots, S_m\}$ 去利用稀疏信息, 进而进行综合推断. 但 CARS 不能应用在相关结构的检验问题中. 其只能用在独立结构的检验问题中.

③ 思想分两步. 一为构建 (T_i, S_i) 可以准确地获取样本信息, 二是有效地整合 T_i 和 S_i 的信息. (类似 US 方法). 算法中的三步: GAP. grouping 步将根据 S_i 将假设分为 K 组. 分成了稀疏水平不同的各个组别; adjusting 步通过调整 p 值将结构信息考虑进来. 结构信息由各组别反映; pooling 步综合考虑调整后的所有 p 值, 选取一个阈值去控制 FDR

④ 根据 $\{S_1, \dots, S_m\}$ 去进行分组会造成信息损失. 组别越少, 分组变量越离散, 信息损失越多. 组别越多, 计算越复杂, 计算成本越高. 实际选择 $K=3$ 或 4 .

⑤ 贡献: (a) 处理了相依结构的多重检验 FDR 问题. (b) 存在的方法分组依赖先验. GAP 由原样本构造 S_i 进行分组. (c) 不恰当的分组会扭曲原假设下 p 值的分布. GAP 通过条件独立性的准则确保合适的分组以及 FDR 的渐近控制. (d) 新颖的 p 值加权方法, 使得其适用于各种相依结构

⑥ 已知组别的多重检验方法: pool 和 separate. 前者是忽视了组别标签且直接应用 BH 方法在所有检验上, 经证明 pool 是不有效的, 甚至可能 invalid. 后者是在各组别内单独应用 BH 方法最后将所有拒绝的假设并到一起, 如 US 方法.

⑦ 经证明, 适当的权重可以控制 FDR, 但 power 会受当权重信息的影响. GAP 的权重是适当的.

⑧ 对于 non-null 比例的估计: 用 $\hat{\pi}_i = \frac{\# \{P_i < \lambda\}}{m\lambda}$ 来估计 non-null 比例, 进而使用 $\hat{\pi}_i = (\varepsilon \vee \hat{\pi}_i^*) \wedge (1 - \varepsilon)$ 使得其比例被约束在 $(\varepsilon, 1 - \varepsilon)$. $\varepsilon = 0.05$. λ 需要调整.

2. 方法步骤 Step 1: (Grouping) ① 计算 (T_i, S_i)

② 定义 $\lambda_0 = -\infty, -4\sqrt{\log m} \leq \lambda_1 < \lambda_2 < \dots < \lambda_{K-1} \leq 4\sqrt{\log m}, \lambda_K = +\infty$.

令 $\Lambda = \{\lambda_1, \dots, \lambda_K\}$ 是 $\mathcal{X} = \{\frac{1}{\sqrt{N}} \sqrt{\log m}, l = -4N, \dots, 4N\}$, N 充分大的一个子集, 根据 S_i 进行分组. $C_l = \{1 \leq i \leq m: \lambda_{l-1} < S_i \leq \lambda_l\}, 1 \leq l \leq K$.

Step 2: (Adjusting) 定义 $m_l = |C_l|$, 计算调整 p 值 $p_i^w = \min\{\frac{p_i}{w_i}, 1\}, i \in C_l$.

权重的计算: ① initial adjusting: 估计各组 non-null 比例 $\hat{\pi}_l$, 则各组权重为 $w_l = \left\{ \frac{K}{l-1} \frac{m_l \hat{\pi}_l}{1 - \hat{\pi}_l} \right\}^{-1} \frac{m_l \hat{\pi}_l}{(1 - \hat{\pi}_l)}, 1 \leq l \leq K$. 有 $p_i^w = \min\{\frac{p_i}{w_l}, 1\}, i \in C_l$

② further refining: 寻求所有 Λ 的可能性, 去寻求最优分组. $\Lambda = \phi$. $\Lambda^* = \arg\max \{\Lambda: k\}$. 即为不分组. 组合所有 p_i^w 进行 BH procedure. 记 $k = \max\{l: \hat{\pi}_l \leq \frac{\alpha}{m}\}$.
→ 最优权重 w_i 基于最优分组计算而来. 即便得拒绝数最多的分组



Step 3: (Pooling) 基于最优分组计算出最优权重, 得出调整后的 p 值 p_i^w .
对所有调整后的 p 值进行 BH procedure. 对于所有 $k = \max\{i: p_{(i)}^w \leq \frac{\alpha}{m}\}$, 其中 $p_{(i)}^w$ 为 $p_i^w = \min(\frac{p_i}{w_i}, 1)$, $i=1, \dots, m$, $i \in G_k$ 的次序统计量, 拒绝 H_{k1}, \dots, H_{kk} .

3. 理论假设

NOTE: GAP 的 FDP 定义为 $FDP_{GAP} = \frac{\sum_{i \in H_0} I(p_i^w \leq p_{(k^w)}^w)}{\sum_{i=1}^m I(p_i^w \leq p_{(k^w)}^w)} \vee 1$, 其中 $k^w = \max\{i: p_{(i)}^w \leq \frac{\alpha}{m}\}$

而 $FDR = E(FDP_{GAP})$. 以下为理论证明必要的假设. 令 A_2 为 H_0 的补集, $|A_2| = o(m^v)$, $\forall v > 0$.
定义 $\tilde{H}_0 = H_0 \setminus A_2$, $n = n_1 + n_2$

(A1) 渐近正态性: 对检验统计量 $\{T_i, i \in \tilde{H}_0\}$, 存在两个 iid 随机样本 $\{Z_{k,i}, i=1, \dots, n_1\}$, $\{Z_{k,i}, i=n_1+1, \dots, n_1+n_2\}$, 且 $E(Z_{k,i}) = 0$, $E(e^{kZ_{k,i}}) < \infty$ 对某个 $k > 0$ 使得 $M > 0$, $\exists b_m, b_m = o(\log m^{\frac{1}{2}})$

$$P_{H_0,i}(|T_i - \frac{\sum_{k=1}^n Z_{k,i}}{\sqrt{\text{Var}(\sum_{k=1}^n Z_{k,i})}}| \geq b_m) = O(m^{-M})$$

(A2) 弱相依性: 定义 $(\beta_{j,1}) = R_1 = \text{corr}(Z_k)$, $k \in \{1, \dots, n_1\}$, $(\beta_{j,2}) = R_2 = \text{corr}(Z_k)$, $k \in \{n_1+1, \dots, n\}$

$Z_k = (Z_{k,1}, \dots, Z_{k,m})$, 有 $\max_{1 \leq i, j \leq m} |\beta_{j,a}| \leq \rho_a < 1$, $a=1, 2$. 进一步地, 存在 $r > 0$ 使得 $\max_{1 \leq i \leq m} |\tau_i(r)| \leq \tau_i(r) = \{j: 1 \leq j \leq m, |\beta_{j,a}| \geq (\log m)^{-r}\}$, $d=1$ 或 2 .

(A3) 渐近独立性: T_i 和 S_i 渐近独立, 在 null T, i.e. \forall 常数 $M > 0$, 有

$$P_{H_0,i}(|T_i| \geq t, |S_i| \geq \lambda) = (1+o(1))G(t)P(|N(0,1)+S_i| \geq \lambda) + O(m^{-M})$$

一致地对于 $0 \leq t \leq 4\sqrt{\log m}$, $0 \leq \lambda \leq 4\sqrt{\log m}$, $i \in \tilde{H}_0$, $S_i = E(S_i)$, $\forall 0 \leq j \leq 4N$, 固定的 N , 有

$$P_{H_0,i}(|T_i| \geq t, |S_i| < \lambda) = (1+o(1))G(t)P(|N(0,1)+S_i| < \lambda) + O(m^{-M})$$

一致地对于 $0 \leq t \leq 4\sqrt{\log m}$, $0 \leq \lambda \leq \frac{4}{N}\sqrt{\log m}$, $i \in \tilde{H}_0$, 其中 $\lambda_N = \frac{4}{N}\sqrt{\log m}$.

NOTE: (A1) 比较弱, 只要求了 $T_i \sim AN(0,1)$, (A2) 指明不是所有统计量都强相关, (A3) 由 (T_i, S_i) 的构造可知. 定义 $S_p = \{i: 1 \leq i \leq m, |\beta_{i,1} - \beta_{i,2}| \geq (\frac{(\log m)^{1/r}}{m})^{\frac{1}{2}}\}$

4. 理论结果

定理一: 设 $\rho > 0$, $\delta > 0$, $|S_0| \geq \frac{1}{\sqrt{\alpha}} + \delta)(\log m)^{\frac{3}{2}}$, n_1, n_2 , $m_0 = |H_0| \geq cm$, $\exists c > 0$. 则在 (A1)-(A3) 且 $\log m = o(n^{\frac{1}{2}})$, $c > 5$ 有 $\lim_{m, n \rightarrow \infty} FDP_{GAP} \leq \alpha$, 且 $\lim_{m, n \rightarrow \infty} P(FDR_{GAP} \leq \alpha + \varepsilon) = 1$

定理二: 在定理一的条件下, 有 $\Psi_{GAP} \geq \Psi_{BH} + o(1)$, $m \rightarrow \infty$.

Pf: 由 Thm 1, $P(\frac{FDP_{GAP}}{\alpha'} - 1 \geq \varepsilon) \rightarrow 0$, $c \in (0, 1]$, $P(\frac{FDP_{BH}}{\alpha'} - 1 \geq \varepsilon) \rightarrow 0$, $\alpha' = \frac{m_0}{m}\alpha$

即有 $FDP_{GAP} + o(1) \leq FDP_{BH}$, 即 $\frac{\sum_{H_0} I(p_i \leq \alpha)}{\sum_{H_0} I(p_i \leq \alpha)} \geq \frac{\sum_{H_0} I(p_i^w \leq p_{(k^w)}^w)}{\sum_{H_0} I(p_i^w \leq p_{(k^w)}^w)}$, $\because \Lambda = \emptyset$ 时.

GAP 和 BH 等价. 又由 GAP 的最优组的定义, $\sum_{H_0} I(p_i^w \leq p_{(k^w)}^w) \geq \sum_{H_0} I(p_i \leq p_{(k)}) \leq A \geq B$

要证: $\sum_{H_0} I(p_i^w \leq p_{(k^w)}^w) \geq \sum_{H_0} I(p_i \leq p_{(k)})$, 即证: $(1-\alpha')A - (1-\alpha)B \geq 0$, $c \in (0, 1]$, 此为 c 的减函数, $c=1$ 时取 min. 为 $(1-\alpha')A - (1-\alpha)B \geq 0$. 成立. 进而有

$$\frac{\sum_{H_0} I(p_i^w \leq p_{(k^w)}^w)}{|H_0|} \geq \frac{\sum_{H_0} I(p_i \leq p_{(k)})}{|H_0|} + o(1), \Rightarrow \Psi_{GAP} \geq \Psi_{BH} + o(1)$$