



Covariate-assisted ranking and screening - CARS procedure.

1. 改进目标：本文介绍了一种两样本的多重检验方法，构造辅助统计量以包含更多信息，提高 Power。传统方法一般将样本提炼为一向量（如 p 值）进而选择 cutoff 进行比较，但第一步中的提炼样本信息可能会造成显著的信息损失从而导出非最优 (suboptimal) 的方法。CARS 通过 **辅助协变量序列提取结构信息**（本文考虑均值两样本稀疏情形）

① CARS 分两步，构造 primary test statistics 和 auxiliary covariate，进而将二者结合构造多重检验方法。思想是通过辅助统计量改变各假设的重要性，不再 exchangeable，故关键的一步为利用辅助统计量去反映结构信息

② 相关的两个技巧为 screening 以及 grouping，前者将问题视为 unequal，通过 screen-and-clean 的方法。先将 x_i 和 y_i 的支撑筛选出来，进而再做检验。而后者为将检验分组，divide-and-test。如 Liu (2014) 的 US 方法，而 Screen 方法可理解为 narrow 了我们本来要考虑的，以提高 test 的准确性。但样本的分割往往导致 power loss。这两个方法都涉及到将一个连续变量分成几部分，这导致无法充分利用辅助统计量的信息（在一个离散点切割，本应概率为 0）

③ 本文主要关注 Gaussian 情形，数值表现 non-Gaussian 亦可。设 T_{12}, T_{21} 为统计量，包含了以下

$\theta_{12} = I(u_{x1} \geq u_{y1})$ 和 $\theta_{21} = I(u_{x1} \neq 0 \text{ 或 } u_{y1} \neq 0)$ 的信息，显然有 $\theta_{12} = 0 \Rightarrow \theta_{21} = 0$ 。

(1) 构造统计量：首先， T_{12} 要获得 θ_{12} 的信息，形如 $\bar{x}_1 - \bar{y}_1$ 。其次，获得 θ_{21} 的信息，形如 $\bar{x}_1 + K\bar{y}_1$
 $K = \frac{n_2 \sigma_{y1}}{n_1 \sigma_{x1}}$ ，使 T_{12}, T_{21} 零相关。最后进行标准化即可

$$(T_{12}, T_{21}) = \sqrt{\frac{n_1 n_2}{n}} \left(\frac{\bar{x}_1 - \bar{y}_1}{\sigma_{x1}}, \frac{\bar{x}_1 + K\bar{y}_1}{\sqrt{K} \sigma_{y1}} \right), \quad \sigma_{x1}^2 = \frac{n_2}{n_1} \sigma_{y1}^2 + \frac{n_1}{n} \sigma_{x1}^2$$

(2) 构造二元模型。考虑 T_{12}, T_{21} 的 jpdf: $f(t_{12}, t_{21}) = \sum_{(y, k) \in \{0, 1\}^2} \pi_{y,k} f(t_{12}, t_{21} | \theta_{12} = y, \theta_{21} = k)$
 $\pi_{y,k} = P(\theta_{12} = y, \theta_{21} = k)$ 。令 $\pi_y = P(\theta_{12} = y)$ ，根据 $\pi_{10} = 0$ ，以及零相关（正态独立），可简化该式。见 Prop 1。

(3) 多重检验。利用了 $mFDR = FDR + o(1)$ 。Power 定 X 有很多，ETP (本文)，average power，missed discovery rate，FNP, FNR 等。令 $T_{OR}^1(t_1, t_2) = P(\theta_{12} = 0 | T_{12} = t_1, T_{21} = t_2) = \frac{q^*(t_2) f_{10}(t_1)}{f(t_1, t_2)}$
其中 $f_{10}(t_1) = f(t_1 | \theta_{12} = 0)$ ， $q^*(t_2) = (1 - \pi_1) f(t_2 | \theta_{12} = 0)$

④ 注：oracle statistic T_{OR}^1 可看成是后验概率， H_{10} 为真，给定 T_{12}, T_{21} 的概率，这可以作为对 H_{10} 是否为 null 的显著性水平。这种将 T_{12}, T_{21} 通过 pooling 成 T_{OR}^1 的方法没有信息损失。若考虑最坏的情况， T_{21} 无任何信息： $n_x = n_y$ ， $\sigma_{x1}^2 = \sigma_{y1}^2$ ， $u_{x1} = -u_{y1}$ ， T_{OR}^1 化为 Lfdr 统计量（即单变量模型下的 Adaptive z-value 方法）。CARS 也可以分为两步：ranking 和 thresholding。

⑤ 统计量 $T_{OR} = \frac{q^*(t_2) T_{OR}^1(t_1)}{f(t_1, t_2)}$ 的估计。 $f_{10}(t_1)$ 是当 $\theta_{12} = 0$ 时的 mpdf，假设是已知的。 $f(t_1, t_2)$ 是 gpdf 结合前文所提到的化简方式，利用 **非参的核估计** 方法进行估计，重点是对 $q^*(t_2)$ 进行估计
 $q^*(t_2) = (1 - \pi_1) f(t_2 | \theta_{12} = 0)$

(1) 方法一。计数 $Q^*(t_2, h) = \#\{i : T_{21} \in [t_2 - \frac{1}{h}, t_2 + \frac{1}{h}], \theta_{12} = 0\} / m$ ， $q^*(t_2) = \lim_{h \rightarrow 0} \frac{E(Q^*(t_2, h))}{h}$ ，但由于 θ_{12} 未知无法估计。文章利用了 Screening 方法对 null 进行了筛选，设 p_i 为 T_{21} 对应的 p 值，假设 p 值很大的认定为 null，即 $p_i > \tau$ ，则 $\theta_{12} = 0$ 。由此估计 $Q^*(t_2, h) = \#\{i : T_{21} \in [t_2 - \frac{1}{h}, t_2 + \frac{1}{h}], p_i > \tau\}$
估计 $q^*(t_2) = \lim_{h \rightarrow 0} \frac{E(Q^*(t_2, h))}{h} = \frac{\int_{t_2 - \tau, p_i > \tau} f(t_1, t_2) dt_1}{1 - \tau}$ ，综上 $f(t_1, t_2) = \frac{1}{m} \sum_{i=1}^m K_{h1}(t_1 - t_{1i}) K_{h2}(t_2 - t_{2i})$



其中, $K(t)$ 为核函数, h_1, h_2 为带宽, 设 $\tau(\tau) = \{i: p_i > \tau\}$. 有 $\hat{q}^*(t_2) = \frac{\sum_{i \in \tau(\tau)} K_{h_2}(t_2 - t_{2i})}{m(1-\tau)}$

则 $\hat{\tau}_{OR}(t_1, t_2) = \frac{\hat{q}^*(t_2) f_{10}(t_1)}{f(t_1, t_2)} \wedge 1$

(2) 方法二: 另一种 $q^*(t_2)$ 的估计, 是一种相合估计, 且由最优化理论中得到, 如下

$\hat{q}^*(t_2) = \frac{1}{k} \sum_{j=1}^k \frac{(\hat{z}_2 - \hat{z}_1)(\tau_j - \tau_k) K_{h_2}(\tau_j - \tau_k) \hat{q}(\tau_j)}{\hat{z}_2 \hat{z}_0 - \hat{z}_1^2}$, 其中, $\hat{z}_\tau = \frac{1}{k} \sum_{j=1}^k (\tau_j - 1)^r K_{h_1}(\tau_j - \tau_k)$

τ_0, \dots, τ_k 为 $(0, 1)$ 上的分点, $\hat{q}(\tau_j)$ 为方法一的估计. 此法为获得最小的偏差, 在不增加过多方差的前提下.

二. 方法步骤

Step 1: 计算检验统计量 T_{1i} 和辅助统计量 T_{2i} , $i=1, \dots, m$

Step 2: 估计 pool 统计量 \hat{T}_{OR}^1 , 并排序 $\hat{T}_{OR}^{(1)} \leq \dots \leq \hat{T}_{OR}^{(m)}$

Step 3: 令 $k = \max\{j: \frac{1}{j} \sum_{i=1}^j \hat{T}_{OR}^{(i)} \leq \alpha\}$, 拒绝 $H_{(1)}, \dots, H_{(k)}$.

NOTE: 这里的 $\frac{1}{j} \sum_{i=1}^j \hat{T}_{OR}^{(i)}$ 是对 mFDR 的估计, 思想与 Adaptive z -value 一样.

三. 理论结果

定理一 (ORACLE) (1) 对 $0 < \lambda \leq 1$, 设 $Q_{OR}(\lambda)$ 为 $\{I(T_{OR} < \lambda): 1 \leq i \leq m\}$ 的 mFDR. 则 $Q_{OR}(\lambda) < \lambda$ 且 $Q_{OR}(\lambda)$ 为 λ 的非降函数. (2) 设选取 $\alpha < \bar{\alpha} = Q_{OR}(1)$, 则 oracle 阈值 $\lambda_{OR} = \sup\{\lambda: Q_{OR}(\lambda) \leq \alpha\}$ 存在且唯一, $Q_{OR}(\lambda_{OR}) = \alpha$. 进而, 定义 $\delta_{OR} = (\delta_{OR}^i, i=1, \dots, m)$, 其中, $\delta_{OR}^i = I(T_{OR} < \lambda_{OR})$, 则 δ_{OR} 是最优的, 即对 \forall FDR $\leq \alpha$ 的检验类 \mathcal{D}_α , $\delta \in \mathcal{D}_\alpha$, $ETP_\delta \leq ETP_{\delta_{OR}}$, 且 \mathcal{D}_α 是包含基于 T_1, T_2 的 mFDR $\leq \alpha$ 的检验.

证明: 首先 $T_{OR}(t_1, t_2) = \frac{P(\Theta_i = 0 | f(t_1, t_2) | \Theta_i = 0)}{f(t_1, t_2)} = \frac{(1 - \pi_1) f(t_1 | \Theta_i = 0) f(t_2 | \Theta_i = 0)}{f(t_1, t_2)} = \frac{q^*(t_2) f_{10}(t_1)}{f(t_1, t_2)}$

(1) ① 证 $\alpha t < t$, 在 Adaptive z -value 文献查知. $Q_{OR}(\lambda) = E(I(T_{OR} < \lambda)) \rightarrow$ 由 Prop 1 条件独立性.

设 $Q_{OR}(t) = \alpha t$. 由 mFDR 定义 ($Q_{OR}(\lambda)$ 定义). $E\{\sum_{i=1}^m (T_{OR}^i - \alpha t) I(T_{OR}^i < t)\} = \alpha A.1)$ 对 (t_1, t_2) 积分. 由 (A.1) 得 $\alpha t < t$. ② 证 $Q_{OR}(t)$ 对 t 单调. 设 $Q_{OR}(t_y) = \alpha y, y=1, 2$. 证 $t_1 < t_2$. 有 $\alpha_1 \leq \alpha_2$. 反证若 $\alpha_1 > \alpha_2$ 则 $(T_{OR} - \alpha_2) I(T_{OR} < t_2) = (T_{OR} - \alpha_2) I(T_{OR} < t_1) + (T_{OR} - \alpha_2) I(t_1 \leq T_{OR} \leq t_2)$

$\geq (T_{OR} - \alpha_1) I(T_{OR} < t_1) + (\alpha_1 - \alpha_2) I(T_{OR} < t_1) + (T_{OR} - \alpha_1) I(t_1 \leq T_{OR} \leq t_2)$ 两边取期望, 有 $E\{(T_{OR} - \alpha_2) I(T_{OR} < t_2)\} > 0$. 这与 (A.1) 矛盾.

(2) oracle 阈值 $t_{OR} = \sup\{t \in (0, 1), Q_{OR}(t) \leq \alpha\}$, 要证明 t_{OR} 可取到. $Q_{OR}(\lambda) = \frac{(1-p) Q_{OR}(\lambda)}{(1-p) Q_{OR}(\lambda)}$ 为连续且

又为单调函数. 故有 $Q_{OR}(t_{OR}) = \alpha, \alpha < \bar{\alpha}$. 定义 δ_{OR} . 令 δ^* 为任意决策 mFDR(δ^*) $\leq \alpha$, 则有 $E(\sum_{i=1}^m (T_{OR}^i - \alpha) \delta_{OR}^i) = 0, E(\sum_{i=1}^m (T_{OR}^i - \alpha) \delta_{OR}^i) \leq \alpha$ (A.2) 由此有 $E(\sum_{i=1}^m (\delta_{OR}^i - \delta_{OR}^*) (T_{OR}^i - \alpha)) \geq 0$. 考虑变换, $f(x) = \frac{x - \alpha}{1 - x}$ (个), $\delta_{OR}^i \Leftrightarrow \delta_{OR}^i = I(\frac{T_{OR}^i - \alpha}{1 - T_{OR}^i} < \lambda_{OR})$. $\lambda_{OR} = f(t_{OR})$.

$\therefore \alpha < t_{OR} < 1$, 则有 $\lambda_{OR} > 0$. 注意有 $\begin{cases} \delta_{OR}^i > \delta_{OR}^i, T_{OR}^i - \alpha - \lambda_{OR}(1 - T_{OR}^i) < 0 \\ \delta_{OR}^i < \delta_{OR}^i, T_{OR}^i - \alpha - \lambda_{OR}(1 - T_{OR}^i) > 0 \end{cases}$ 分类讨论即可.

$\forall i, (\delta_{OR}^i - \delta_{OR}^*) (T_{OR}^i - \alpha - \lambda_{OR}(1 - T_{OR}^i)) \leq 0$. 对 i 求和取期望, $E(\sum_{i=1}^m (\delta_{OR}^i - \delta_{OR}^*) (T_{OR}^i - \alpha - \lambda_{OR}(1 - T_{OR}^i))) \leq 0$ 结合 (A.3), (A.4). 有 $\lambda_{OR} E(\sum_{i=1}^m (\delta_{OR}^i - \delta_{OR}^*) (1 - T_{OR}^i)) \geq E(\sum_{i=1}^m (\delta_{OR}^i - \delta_{OR}^*) (T_{OR}^i - \alpha)) \geq 0$, 由 Adaptive z -value 知 $ETP = E(\sum_{i=1}^m \delta_{OR}^i) = E(\sum_{i=1}^m (1 - T_{OR}^i) \delta_{OR}^i)$. 由此得 $\therefore \lambda_{OR} > 0, ETP_{\delta_{OR}} \geq ETP_{\delta^*}$



同濟大學
TONGJI UNIVERSITY

地址：中国上海市四平路1239号 邮编：200092
1239 SIPING ROAD SHANGHAI CHINA 200092
电话 (TEL) : +86 21- 传真 (FAX) : +86 21-
网址 (WEB) : www.tongji.edu.cn

定理=(Data-Driven) 条件 (1) $E\|\hat{q}^T - q^T\|^2 \rightarrow 0$. (1)' $E\|\hat{q}^X - q^X\|^2 \rightarrow 0$

(2) $E\|\hat{f} - f\|^2 = E[\int \int \hat{f}(t_1, t_2) - f(t_1, t_2)^2 dt_1 dt_2] \rightarrow 0$.

(1) 若条件 (1), (2) 成立, 则 mFDR 和 FDR 在 CARS 下控制在 $\alpha + o(1)$.

2) 若条件 (1)', (2) 成立, 用 \hat{q}^X 估计 q^X , 则 CARS 的 FDR 水平为 $\alpha + o(1)$, 进而, 用 ETP_{CARS} 和 ETP_{OR} 表示 CARS 和 oracle 的 ETP, 则有 $ETP_{CARS} / ETP_{OR} = 1 + o(1)$