

PP 检验的渐近功效¹⁾

张健成平

(中国科学院系统科学研究所)

引言

投影寻踪 (Projection Pursuit, 简称 PP) 是一种处理高维数据的统计方法, 近年来借助于计算机的发展, 它的理论得到了迅速的发展。假如 Q 为一能反映统计性质的指标, 对每个 P 维方向 a , 计算出 $Q(a^T X_1, \dots, a^T X_n)$, 从中找出使 $Q(a^T X_1, \dots, a^T X_n)$ 最大的方向 a_0 , 通过研究数据 $a_0^T X_1, \dots, a_0^T X_n$ 的性质, 来了解原数据 X_1, \dots, X_n 的性质, 这就是所谓数值 PP 的基本想法。如果用一维检验或估计统计量 Q 作指标, 则可得到多维检验或估计 $\sup_{\|a\|=1} Q(a^T X_1, \dots, a^T X_n)$ 。见文献 [2,3,5—7]。在上述文章中这些检验和估计的渐近性质得到了部分研究, 但当对立假设成立时, 有关这些检验统计量的渐近分布问题所得结果尚不多, 本文系统地研究了这类问题。我们部分地采用 [4] 中的记号: $Pf(X)$ 记 $f(X)$ 的期望 ($X \sim P$), P_n 为经验分布, “ \rightsquigarrow ” 记依分布收敛。

§ 1. 主要结果与证明

定理. 设有概率分布 P , $X \sim P$, 分布泛函 $V(t) \triangleq V(t, P)$ 及样本泛函 $V_n(t) \triangleq V_n(t, P_n)$, $t \in \pi$, 其中 P_n 是经验分布, $\pi \subseteq \mathcal{G}$, (\mathcal{G}, d) 是一距离空间, $V(t)$ 关于距离 d 在 π 上连续, $B = \{t \in \pi : V(t) = \sup_{s \in \pi} V(s)\}$ 非空。令

$$S_n(t) \triangleq \sqrt{n}(V_n(t) - V(t)), n \geq 1, t \in \pi. \quad (*)$$

若存在过程 $\{S(t); t \in \pi\}$ 满足

$$\sup_{t \in \pi} |S_n(t) - S(t)| \rightarrow 0 \text{ a.s. } (P) \quad (*)$$

且 $\{S(t); t \in \pi\}$ 的样本轨迹在 π 上一致连续, 一致有界。则

$$\sqrt{n}(\sup_{t \in \pi} V_n(t) - \sup_{t \in \pi} V(t)) \rightarrow \sup_{t \in B} S(t), \text{ a.s. } (P).$$

注. 1) 对 \inf 也有同样的结论。

2) 条件 $S(t)$ 在 π 上一致连续, 可改成存在 $\delta > 0$, 使 $S(t)$ 在 $B^\delta = \{t \in \pi : d(t, B) < \delta\}$ 内一致连续, 则定理仍然成立。

证. 因为 $V(t)$ 连续, 所以 B 是闭集, 令 $B^c \triangleq \pi - B$ 。现令 $h = \sup_{t \in \pi} V(t)$, 则

1989年7月5日收到。

1) 献给系统科学研究所建所十周年 (1979—1989)。

$$\begin{aligned}
& \sqrt{n} (\sup_{t \in \pi} V_n(t) - \sup_{t \in \pi} V(t)) \\
&= \sqrt{n} \sup_{t \in \pi} (V_n(t) - h) \\
&= \sqrt{n} \max \left\{ \sup_{t \in B} (V_n(t) - h), \sup_{t \in B^c} (V_n(t) - h) \right\} \\
&= \sqrt{n} \max \left\{ \sup_{t \in B} (V_n(t) - V(t)), \sup_{t \in B^c} (V_n(t) - V(t)) \right\} \\
&= \max \left\{ \sup_{t \in B} S_n(t), \sup_{t \in B^c} (S_n(t) + \sqrt{n} (V(t) - h)) \right\},
\end{aligned}$$

现任意固定使(*)式收敛性成立的样本点。

由 $S(t)$ 在 π 上一致连续性知

$\forall \varepsilon > 0, \exists \delta > 0$, 使当 $t_1, t_2 \in \pi, d(t_1, t_2) < \delta$ 时, 有 $|S(t_1) - S(t_2)| < \varepsilon$.

因为 B 是闭集, 所以任给 $t \in B^c \cap \{t: d(t, B) < \delta\}$, 存在 $b_t \in B$, 使 $d(b_t, t) = d(t, B) < \delta$, 从而

$$S(t) < s + S(b_t) \leq s + \sup_{t \in B} S(t), \text{ 也即}$$

$$\begin{aligned}
& \sup_{t \in B^c \cap \{t: d(t, B) < \delta\}} (S_n(t) - S(t) + S(t) + \sqrt{n} (V(t) - h)) \\
& \leq \sup_{t \in \pi} |S_n(t) - S(t)| + \varepsilon + \sup_{t \in B} S(t) \rightarrow s + \sup_{t \in B} S(t).
\end{aligned}$$

而

$$\begin{aligned}
& \sup_{t \in B^c \cap \{t: d(t, B) > \delta\}} (S_n(t) - S(t) + S(t) + \sqrt{n} (V(t) - h)) \\
& \leq \sup_{t \in B^c} \sup_{t \in \pi} |S_n(t) - S(t)| + \sup_{t \in B^c} \sup_{t \in \pi} S(t) + \sqrt{n} (V(t) - h); \\
& t \in B^c \cap \{t: d(t, B) \geq \delta\} \xrightarrow{n \rightarrow +\infty} -\infty.
\end{aligned}$$

由 s 的任意性知:

$$\limsup_{n \rightarrow \infty} \sup_{t \in B^c} (S_n(t) + \sqrt{n} (V(t) - h)) \leq \sup_{t \in B} S(t),$$

显然, $\sup_{t \in B} S_n(t) \rightarrow \sup_{t \in B} S(t)$. 所以, $\sqrt{n} (\sup_{t \in \pi} V_n(t) - \sup_{t \in \pi} V(t)) \rightarrow \sup_{t \in B} S(t)$.

推论 1. M -型 PP 指标其最大特征根的极限分布.

设 $\phi(z_1, z_2)$ 关于 z_1 单调减(或是 z_1 的偶函数且在 $(0, \infty)$ 上单调减). T^* 是 $\int \phi(a^T X, z) dP$ $= 0$ 的最小解, T_*^* 是 $\int \phi(a^T X, z) dP_n = 0$ 的最小解, P_n 是经验分布.

设 $N_{\varepsilon_0} = \{(b, z); \exists a \in S_p, \|b - a\| \leq \varepsilon_0, |z - T^*| \leq \varepsilon_0\}$, $S_p = \{a \in \mathbb{R}^p; \|a\| = 1\}$, $\varepsilon_0 \geq 0$.

假定 $\frac{\partial \phi(z_1, z_2)}{\partial z_2}$ 二元连续, $P \frac{\partial \phi(a^T X, T^*)}{\partial z_2} \neq 0$, 且存在 $M_i(X), 0 \leq i \leq 4$, 满足:

$$\begin{aligned} \sup_{N_{t_0}} |\phi(a^*X, t^*)| &\leq M_0(X), \quad \sup_{N_{t_0}} \left| \frac{\partial \phi(a^*X, t^*)}{\partial t_1} \right| \leq M_1(X), \\ \sup_{N_{t_0}} \left| \frac{\partial \phi(a^*X, t^*)}{\partial t_2} \right| &\leq M_2(X), \quad \sup_{N_{t_0}} \left| \frac{\partial^2 \phi(a^*X, t^*)}{\partial t_1 \partial t_2} \right| \leq M_3(X), \\ \sup_{N_{t_0}} \left| \frac{\partial^3 \phi(a^*X, t^*)}{\partial t_1^2} \right| &\leq M_4(X), \text{ 且} \end{aligned}$$

$PM_0^2(X) < +\infty$, $PM_0^2(X)\|X\|^{\delta} < +\infty$ (某 $\delta > 0$), $PM_1(X)\|X\| < +\infty$,
 $PM_2^2(X) < +\infty$, $PM_2^2(X) < +\infty$, $PM_3(X)\|X\| < +\infty$, $PM_4(X) < +\infty$.

则 $\sqrt{n}(\sup_{\|\alpha\|=1} T_{\alpha}^* - \sup_{\|\alpha\|=1} T^*) \rightsquigarrow \sup_{\alpha \in B} W(h(\alpha))$, 其中

$$h(\alpha) \triangleq h(\alpha, X) = \frac{\phi(a^*X, T^*)}{-P \frac{\partial \phi(a^*X, T^*)}{\partial t_1}}, \quad \mathcal{F} = \{h(\alpha) : \|\alpha\| = 1\},$$

W 是 \mathcal{F} 上的 P -桥, $B = \{b : \|b\| = 1, T^b = \sup_{\|\alpha\|=1} T_{\alpha}^*\}$. P -桥的定意见 Pollard^[4].

注记. 如果 $\phi(t_1, t_2) = \frac{t_1^2 - t_2^2}{t_1^2 + t_2^2}$, $X \sim$ 椭球等高分布, 则只要 $P\|X\|^{\delta} < +\infty$ (某 $\delta > 0$), 本定理的结论就成立.

证. 由[7]知

$$\sup_{\|\alpha\|=1} |\sqrt{n}(T_{\alpha}^* - T^*) - \sqrt{n}(P_{\alpha} - P)h(\alpha, X)| = o_p(1),$$

注意到

$$\begin{aligned} &\sqrt{n}(\sup_{\|\alpha\|=1} T_{\alpha}^* - \sup_{\|\alpha\|=1} T^*) \\ &= \sqrt{n}\{\sup_{\|\alpha\|=1} V_{\alpha}(a) - \sup_{\|\alpha\|=1} V(a)\} + o_p(1), \end{aligned}$$

其中 $V_{\alpha}(a) \triangleq P_{\alpha}(h(a, X) + T^*)$, $V(a) \triangleq T^*$,

$$S_{\alpha}(a) \triangleq \sqrt{n}(V_{\alpha}(a) - V(a)) = \sqrt{n}(P_{\alpha} - P)h(a, X), \quad \|\alpha\| = 1.$$

由 Dudley-Philipp[1] 或 Pollard [6] 表现定理知, 存在过程 $\{\bar{S}_{\alpha}(a), S(a) : \|\alpha\|=1\}$, 满足

$\{\bar{S}_{\alpha}(a) : \|\alpha\|=1\}$ 与 $\{S_{\alpha}(a) : \|\alpha\|=1\}$ 同分布, $\{\bar{S}(a) : \|\alpha\|=1\} = \{S(h(a)) : \|\alpha\|=1\}$, 后者是 \mathcal{F} 上 P -桥且

$$\sup_{\|\alpha\|=1} |\bar{S}_{\alpha}(a) - S(a)| \rightarrow 0 \text{ a.s. } P.$$

因为 $S(h(a))$ 的样本轨迹在 \mathcal{F} 上一致连续, 其中指标集 \mathcal{F} 的本身度量为 $L_2(P)$ 半模, 而由 $PM_0^2(X) < +\infty$ 知, 当 $\|\alpha - b\| \rightarrow 0$ 时, $P|h(a) - h(b)|^2 \rightarrow 0$.

注意到 $\{\alpha \in \mathbb{R}^p : \|\alpha\|=1\}$ 是闭集, 所以

$S(a)$ 的样本轨迹在 $\{\alpha \in \mathbb{R}^p : \|\alpha\|=1\}$ 上关于欧氏模一致连续, 且一致有界.

综上定理的条件满足, 因此

$$\sqrt{n}(\sup_{\|\alpha\|=1} T_{\alpha}^* - \sup_{\|\alpha\|=1} T^*) \rightsquigarrow \sup_{\alpha \in B} W(h(\alpha)).$$

推论 2. Friedman 判别标准的功效。

设 $L_i \triangleq L_i(a, Z) \triangleq Q_i(2\Phi(a^T Z) - 1)$, $1 \leq i \leq J$.

Q_i 是 $[-1, 1]$ 上的勒让德多项式, $1 \leq i \leq J$.

$$\Phi \sim N(0, I_p), \quad Z \sim P.$$

一维投影指标:

$$I_p(a) = \frac{1}{2} \sum_{j=1}^J (2j+1) P_R^2(Q_j(R)), \quad R = (2\Phi(a^T Z) - 1), \quad P_R(Q_j(R)) \text{ 是 } Q_j(R) \text{ 的}$$

期望

$$I_n(a) = \frac{1}{2} \sum_{j=1}^J (2j+1) [P_n L_j(a, Z)]^2;$$

二维投影指标:

$$I_p(a, b) = \frac{1}{4} \sum_{j=1}^J (2j+1) E^2[Q_j(R_1)] + \frac{1}{4} \sum_{k=1}^J (2k+1) E^2[Q_k(R_2)] \\ + \frac{1}{4} \sum_{j=1}^J \sum_{k=1}^{J-j} (2j+1)(2k+1) [E(L_j(a, Z)L_k(b, Z))]^2.$$

其中 $R_1 = 2\Phi(X_1) - 1$, $X_1 = a^T Z$, $R_2 = 2\Phi(X_2) - 1$, $X_2 = b^T Z$, $\|a\| = 1 = \|b\|$, $a^T b = 0$, E 表示取期望.

$$I_n(a, b) = \frac{1}{4} \sum_{j=1}^J (2j+1) [P_n L_j(a, Z)]^2 + \frac{1}{4} \sum_{k=1}^J (2k+1) [P_n L_k(b, Z)]^2 \\ + \frac{1}{4} \sum_{j=1}^J \sum_{k=1}^{J-j} (2j+1)(2k+1) [P_n L_j(a, Z)L_k(b, Z)]^2.$$

当 $P \neq \emptyset$ 时,

$$\sqrt{n} (\sup_{\|a\|=1} I_n(a) - \sup_{\|a\|=1} I_p(a)) \sim \sup_{a \in B_1} \sum_{j=1}^J (2j+1) (E_p L_j(a, Z)) \cdot PL_j(a, Z),$$

其中 $B_1 = \{b : I_p(b) = \sup_{\|b\|=1} I_p(b)\}$;

$$\sqrt{n} (\sup \{I_n(a, b) : \|a\| = 1 = \|b\|, a^T b = 0\} - \sup \{I(a, b) : \|a\| = 1 = \|b\|, \\ a^T b = 0\}) \sim 2 \sup_{(a, b) \in B_{11}} \left\{ \frac{1}{4} \sum_{j=1}^J (2j+1) (E_p L_j(a, Z)) PL_j(a, Z) \right. \\ \left. + \frac{1}{4} \sum_{k=1}^J (2k+1) (E_p L_k(b, Z)) PL_k(b, Z) + \frac{1}{4} \sum_{j+k \leq J} (2j+1) \right. \\ \left. \cdot (2k+1) (E_p L_j(a, Z)L_k(b, Z)) P(L_j(a, Z)L_k(b, Z)) \right\}.$$

其中 $B_{11} = \{(a, b) : a^T b = 0, \|a\| = 1 = \|b\|, I_p(a, b) = \sup \{I_p(c, d) : c^T d = 0, \|c\| = \|d\| = 1\}\}$, E_p 是 $\mathcal{F} = \{L_i(a, Z), L_i(c, Z)L_k(d, Z) : i+k \leq J, 1 \leq i \leq J, c^T d = 0, \|a\| = \|c\| = \|d\| = 1\}$ 上的 P -桥.

证. 取 $V_n(a) = I_n(a)$, $V(a) = I_p(a)$.

因为 $S_n(a) \triangleq \sqrt{n} (I_n(a) - I_p(a))$

$$= \sum_{j=1}^J (2j+1) [\sqrt{n} (P_s - P) L_i(a, Z)] PL_i(a, Z) + o_p(1).$$

由 Dudley-Philipp^[3] 及 Pollard^[4] 表现定理知,

$$\begin{aligned} & \text{存在 } \left\{ \bar{S}_s(a), \sum_{j=1}^J (2j+1) (E_p L_i(a, Z)) PL_i(a, Z); \|a\| = 1 \right\} \text{ 满足 } \bar{S}_s(a) \text{ 与 } S_s(a) \\ & - o_p(1) = \sum_{j=1}^J (2j+1) [\sqrt{n} (P_s - P) L_i(a, Z)] PL_i(a, Z) \text{ 同分布且} \\ & \sup_{\|a\|=1} \left| \bar{S}_s(a) - \sum_{j=1}^J (2j+1) E_p L_i(a, Z) PL_i(a, Z) \right| \rightarrow 0, \text{ a.s.P.} \end{aligned}$$

令

$$S(a) = \sum_{j=1}^J (2j+1) E_p L_i(a, Z) \cdot PL_i(a, Z),$$

由于 $\|a - b\| \rightarrow 0$ 时, $P|L_i(a, Z) - L_i(b, Z)|^2 \rightarrow 0$, 所以, $S(a)$ 是 $\{a \in \mathbb{R}^r: \|a\| = 1\}$ 上的一致连续函数。

综上定理条件满足, 从而

$$\sqrt{n} (\sup_{\|a\|=1} I_s(a) - \sup_{\|a\|=1} I_p(a)) \rightsquigarrow \sup_{a \in B_1} \sum_{j=1}^J (2j+1) (E_p L_i(a, Z)) PL_i(a, Z).$$

因为

$$\begin{aligned} & \sqrt{n} (I_s(a, b) - I(a, b)) \\ & = \frac{1}{2} \sum_{j=1}^J (2j+1) [\sqrt{n} (P_s - P) L_i(a, Z)] PL_i(a, Z) \\ & + \frac{1}{2} \sum_{k=1}^J (2k+1) [\sqrt{n} (P_s - P) L_k(b, Z)] PL_k(b, Z) \\ & + \frac{1}{2} \sum_{j=1}^J \sum_{k=1}^{j-1} (2j+1)(2k+1) \sqrt{n} (P_s - P) (L_i(a, Z) L_k(b, Z)) \\ & \quad \cdot P(L_i(a, Z) L_k(b, Z)) + o_p(1). \end{aligned}$$

用类似的强逼近办法得到

$$\begin{aligned} & \sqrt{n} (\sup\{I_s(a, b); a^r b = 0, \|a\| = \|b\| = 1\} - \sup\{I(a, b); a^r b = 0, \|a\| = \|b\| = 1\}) \\ & \rightsquigarrow \sup_{a \in B_{11}} \left\{ \frac{1}{2} \sum_{j=1}^J (2j+1) (E_p L_i(a, Z)) PL_i(a, Z) \right. \\ & + \frac{1}{2} \sum_{k=1}^J (2k+1) (E_p L_k(a, Z)) PL_k(b, Z) \\ & \left. + \frac{1}{2} \sum_{j=1}^J \sum_{k=1}^{j-1} (2j+1)(2k+1) E_p (L_i(a, Z) L_k(b, Z)) P(L_i(a, Z) L_k(b, Z)) \right\}. \end{aligned}$$

推论 3. PP-Neyman 检验的功效

设 $H_0: G = P \leftrightarrow H_1: G \neq P$, P 的 a 方向分布函数记为 F^a . G, P 连续分布.

$$\begin{aligned} K_n(a) &= \sum_{i=1}^J [P_n \pi_i(F^a(a^r X))]^2 \\ K_G(a) &= \sum_{i=1}^J [G \pi_i(F^a(a^r X))]^2, \\ \sqrt{n} (\sup_{\|a\|=1} K_n(a) - \sup_{\|a\|=1} K_G(a)) \\ &\sim \sup_{a \in B} 2 \left\{ \sum_{i=1}^J E_G \pi_i(F^a(a^r X)) \cdot G \pi_i(F^a(a^r X)) \right\}. \end{aligned}$$

其中 $B = \{a \in \mathbb{R}^p: \|a\| = 1, K_G(a) = \sup_{\|a\|=1} K_G(a)\}$, E_G 是 $\{\pi_i(F^a(a^r X)): 1 \leq i \leq J, \|a\| = 1\}$ 上的 P -桥, 而 $\{\pi_i: 1 \leq i \leq J, \pi_i \triangleq 1\}$ 是 $[0, 1]$ 上的正交多项式.

证. 完全类似推论 2, 故略.

推论 4. PP-型 χ^2 拟合优度检验的功效.

假定 G, P 都是连续分布,

$$H_0: G = P \leftrightarrow H_1: G \neq P, X_i \sim G, 1 \leq i \leq n.$$

取 $S_i \triangleq (t_i, t_{i+1}]$, $t_i \uparrow$, $t_{m+1} = +\infty$, $t_1 = -\infty$,

$$Z_n^a = \sum_{i=1}^m \frac{(P_n[a^r X \in S_i] - P[a^r X \in S_i])^2}{P[a^r X \in S_i]},$$

若 $\min_{1 \leq i \leq m} \inf_{\|a\|=1} P[a^r X \in S_i] > 0$, 则

$$\sqrt{n} (\sup_{\|a\|=1} Z_n^a - \sup_{\|a\|=1} Z^a) \sim 2 \sup_{a \in B} \sum_{i=1}^m \frac{E_G I_{[a^r X \in S_i]}}{P[a^r X \in S_i]} (G - P) I_{[a^r X \in S_i]},$$

其中 $B = \{a: Z^a = \sup_{\|a\|=1} Z^a, \|a\| = 1\}$, E_G 是 $\{I_{[a^r X \in S_i]}: 1 \leq i \leq m, \|a\| = 1\}$ 上的 G -桥.

证. 因为

$$\begin{aligned} S_n(a) &\triangleq \sqrt{n} (Z_n^a - Z^a) = 2 \sum_{i=1}^m \frac{\sqrt{n} (P_n - G) I_{[a^r X \in S_i]}}{P[a^r X \in S_i]} \\ &\quad \cdot (G - P)[a^r X \in S_i] + o_p(1). \end{aligned}$$

令

$$S(a) = 2 \sum_{i=1}^m \frac{E_G(I_{[a^r X \in S_i]})[(G - P)I_{[a^r X \in S_i]}]}{P[a^r X \in S_i]},$$

$E_G I_{[a^r X \in S_i]}$ 关于 $L_2(G)$ 半模在 $\{I_{[a^r X \in S_i]}: \|a\| = 1\}$ 上一致连续, 又因为

当 $\|a - b\| \rightarrow 0$ 时,

$$G(I_{[a^r X \in S_i]} - I_{[b^r X \in S_i]})^2 \rightarrow 0, \inf_{\|a\|=1} P[a^r X \in S_i] > 0, 1 \leq i \leq m.$$

从而 $S(a)$ 的样本轨迹 关于 a 在 $\{a \in \mathbb{R}^p: \|a\| = 1\}$ 上一致连续, 一致有界.

用 Dudley-Philipp^[3] 及 Pollard^[4] 的表现定理知, 定理的条件都成立, 从而即证.

推论 5. PP Kolmogorov-Smirnov 统计量的功效.

假定 G, P 是连续分布,

$$H_0: G = P \leftrightarrow H_1: G \neq P, X_i \sim G, 1 \leq i \leq n.$$

$$\beta_n = \sup_{\|a\|=1} \sup_{|t|<+\infty} |(P_n - P) I_{\{a^T X \leq t\}}|$$

$$\beta = \sup_{\|a\|=1} \sup_{|t|<+\infty} |(G - P) I_{\{a^T X \leq t\}}|,$$

$$V_n(a, t) = |(P_n - P) I_{\{a^T X \leq t\}}|$$

$$V(a, t) = |(G - P) I_{\{a^T X \leq t\}}|,$$

显然 $V(a, t)$ 是 (a, t) 的连续函数。

假设 H_1 成立，则

$$\sqrt{n}(\beta_n - \beta) \rightsquigarrow \sup_{(a, t) \in B} E_G I_{\{a^T X \leq t\}} \operatorname{sgn}((G - P) I_{\{a^T X \leq t\}}),$$

其中

$$B = \{(a, t) : a \in \mathbb{R}^p, \|a\| = 1, t \in \mathbb{R}, |(G - P) I_{\{a^T X \leq t\}}|$$

$$= \sup_{\|a\|=1} \sup_{|t|<+\infty} |(G - P) I_{\{a^T X \leq t\}}|\},$$

E_G 是 $\{I_{\{a^T X \leq t\}} : \|a\| = 1, |t| < +\infty\}$ 上的 G -桥。

注记. PP Anderson-Darling 统计量:

$$\sup_{\|a\|=1} \sup_{|t|<+\infty} \frac{|(P_n - P) I_{\{a^T X \leq t\}}|}{[1 - F^*(t)] F^*(t)} \text{ 也有类似的结论,}$$

其中 F^* 是 P 在 a 方向上的分布函数,

证. 当 $(G - P) I_{\{a^T X \leq t\}} \neq 0$ 时,

$$\begin{aligned} S_n(a, t) &\triangleq \sqrt{n}(V_n(a, t) - V(a, t)) \\ &= \sqrt{n} (P_n - G) I_{\{a^T X \leq t\}} \operatorname{sgn}((G - P) I_{\{a^T X \leq t\}}) + o_p(1), \\ &\quad o_p(1) \text{ 与 } a \text{ 无关;} \end{aligned}$$

当 $(G - P) I_{\{a^T X \leq t\}} = 0$ 时,

$$\begin{aligned} S_n(a, t) &\triangleq \sqrt{n}(V_n(a, t) - V(a, t)) \\ &= \sqrt{n} |(P_n - G) I_{\{a^T X \leq t\}}| \\ S(a, t) &= \begin{cases} E_G(I_{\{a^T X \leq t\}}) \operatorname{sgn}((G - P) I_{\{a^T X \leq t\}}), & (G - P) I_{\{a^T X \leq t\}} \neq 0, \\ |E_G(I_{\{a^T X \leq t\}})| \operatorname{sgn}((G - P) I_{\{a^T X \leq t\}}), & \text{否则.} \end{cases} \end{aligned}$$

因为 H_1 成立, 所以,

存在 $\delta > 0$, 使在 B 的 δ 闭邻域内 $(G - P) I_{\{a^T X \leq t\}} \neq 0$. 由于 $S(a, t)$ 在 $(G - P) I_{\{a^T X \leq t\}} \neq 0$ 上是 $G(a^T X \leq t)$ 的一致连续函数。

又因为 G, P 是连续分布, 即 $G[a^T X \leq t], P[a^T X \leq t]$ 是 (a, t) 的一致连续函数, 从而 $S(a, t)$ 在 B 的 δ 闭邻域内关于 (a, t) 一致连续。

从而定理注 2) 的条件满足, 即证。

推论 6. PP-VonMise-Smirnov 统计量.

设 $H_0: G = P \leftrightarrow H_1: G \neq P; G, P$ 连续;

$$\begin{aligned} w_* &= \sup_{\|a\|=1} \int ((P_a - P) I_{[a^T X \leq t]})^2 dF^*(t) \\ &= \sup_{\|a\|=1} \int ((P_a - P) I_{[a^T X \leq a^T Y]})^2 dF(Y); \end{aligned}$$

F^* 是 P 在 a 方向上的边缘分布;

$$w = \sup_{\|a\|=1} \int ((G - P) I_{[a^T X \leq t]})^2 dF^*(t);$$

$$V_*(a) = \int ((P_a - P) I_{[a^T X \leq t]})^2 dF^*(t);$$

$$V(a) = \int ((G - P) I_{[a^T X \leq t]})^2 dF^*(t); \text{ 则}$$

$$\sqrt{n}(w_* - w) \rightsquigarrow 2 \sup_{a \in B} \int E_G(I_{[a^T X \leq t]}) \cdot (G - P) I_{[a^T X \leq t]} dF^*(t),$$

其中 $B = \{a : \|a\| = 1, V(a) = \sup_{\|a\|=1} V(b)\}$, E_G 是 $\{I_{[a^T X \leq t]} : \|a\| = 1, |t| < +\infty\}$ 上的 G -桥.

$$\begin{aligned} \text{证. } S_*(a) &\triangleq \sqrt{n}(V_*(a) - V(a)) \\ &= \int \sqrt{n} (P_a - G) I_{[a^T X \leq t]} (P_a + G - 2P) I_{[a^T X \leq t]} dF^*(t) \\ &= 2 \int \sqrt{n} (P_a - G) I_{[a^T X \leq t]} (G - P) I_{[a^T X \leq t]} dF^*(t) + o_p(1). \end{aligned}$$

取强逼近过程。

$$S(a) = 2 \int E_G I_{[a^T X \leq t]} \cdot (G - P) I_{[a^T X \leq t]} dF^*(t).$$

由 P, G 连续及 E_G 在 $\{I_{[a^T X \leq t]} : \|a\| = 1, |t| < +\infty\}$ 上关于 $L_1(G)$ 一致连续知,

$S(a)$ 关于 a 连续, 从而 $S(a)$ 在 $\{a \in \mathbb{R}^n : \|a\| = 1\}$ 上一致连续且一致有界. 由 Pollard^[6] 表现定理知所需要的条件成立, 即证.

推论 7. PP-L 统计量(作为位置和刻度的估计).

设 $X_i \sim G$, $1 \leq i \leq n$, G 连续分布函数, G_n 是 G 的经验分布函数, G^* 是 G 在 a 方向上的边缘分布函数, 同理 G_n^* 是 G_n 在 a 方向上的边缘分布函数. h 是偏对称(Skew-Symmetric), 有界连续(或对称有界连续).

$$T(G^*) = \int x h(G^*(x)) dG^*(x),$$

$$T(G_n^*) = \int x h(G_n^*(x)) dG_n^*(x).$$

设 G 满足[10]定理 3.2 的条件, 则

$$T(G_n^*) - T(G^*) = - \int \frac{G_n^*(x) - G^*(x)}{g_n(x)} h(G^*(x)) dG^*(x) + O(n^{-\frac{1}{2}} \log n),$$

$g_n(x)$ 是 G^* 的密度(x 处的).

令 $V(a) = T(G^*)$, $V_*(a) = T(G_n^*)$, 有

$$\sqrt{n}(V_*(a) - V(a)) \rightsquigarrow \sup_{a \in B} \left[- \int \frac{E_G I_{[a^T X \leq t]}}{g_n(t)} h(G^*(t)) dG^*(t) \right],$$

其中 $B = \{a : \|a\| = 1, T(G^a) = \sup_{1 \leq i \leq n} T(G^a_i)\}$, E_G 是 $\{I_{\{a^r X_i < 0\}} : \|a\| = 1, |r| < +\infty\}$ 上的 G -桥。

证. 类同推论 6, 故略。

注. 作为检验统计量 $\sup_{1 \leq i \leq n} |T(G^a_i) - T(G^a)|$ 或 $\sup_{1 \leq i \leq n} \left| \frac{T(G^a_i)}{T(G^a)} - 1 \right|$, 其中 $H_0: P = G \leftrightarrow H_1: P \neq G$, $X_i \sim P$, $1 \leq i \leq n$, P , G 是位置或刻度参数族中的元素。此时, 用本文的方法, 可得到这些统计量的功效。

推论 8. PP-Mann-Whitney 检验。

设 $G(X) = P(X - r)$, r 为未知 p 维参数向量, P 是连续分布。

$$H_0: r = 0 \leftrightarrow H_1: r \neq 0.$$

取独立样本: X_1, \dots, X_n i. i. d. $\sim P$,

$$Y_1, \dots, Y_n$$
 i. i. d. $\sim G$.

$$\begin{aligned} U_n^* &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left[I_{\{a^r X_i < a^r Y_j\}} - \frac{1}{2} \right] \\ &= (P \times G)_* I_{\{a^r X_i < a^r Y_j\}} - \frac{1}{2}, \end{aligned}$$

其中 $(P \times G)_* = \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \delta_{(X_i, Y_j)}$, $\delta_{(X_i, Y_j)}$ 为单点测度。

$$V_n(a) = (P \times G)_* I_{\{a^r X_i < a^r Y_j\}} - \frac{1}{2},$$

$$V(a) = (P \times G) I_{\{a^r X_i < a^r Y_j\}} - \frac{1}{2}.$$

当 H_0 成立时,

$$\sqrt{6n} (\sup_{1 \leq i \leq n} V_n(a) - \sup_{1 \leq i \leq n} V(a)) \sim \sup_{a \in B} Z^a, Z^a$$

$\mathcal{F} = \{I_{\{a^r X_i < a^r Y_j\}} : \|a\| = 1\}$ 上的 $P \times G$ -桥,

$$B = \{a : \|a\| = 1, V(a) = \sup_{1 \leq i \leq n} V(a)\}.$$

注. 类似, 对 PP-U 统计量也有类似的结论。

证. 类似前面的证法, 故略。

推论 9. PP 两样本检验

设 F 为 k 维分布函数, $G(X) = F(X - r)$, r 为 p 维未知参数。

$$H_0: r = 0, \leftrightarrow H_1: r \neq 0.$$

设有样本 X_1, \dots, X_n i. i. d. F

$$Y_1, \dots, Y_m$$
 i. i. d. G , $\frac{m}{n} = \lambda$, $0 < \lambda < +\infty$,

$\text{med}(a^r X) \triangleq a^r X_1, \dots, a^r X_n$ 的中位数,

$\text{mad}(a^r X) = |a^r X_1 - \text{med } a^r X|, \dots, |a^r X_n - \text{med } a^r X|$ 的中位数,

$$T_{n,m}^* = \frac{\text{med } a^r X - \text{med } a^r Y}{\text{mad}(a^r (X \cup Y))},$$

$$T_{s,m} = \sup_{\|s\|=1} T_{s,m}^s.$$

设 F^* 为 F 在 a 方向上边缘分布函数, 密度为 f^* , $\bar{F}^* \triangleq \frac{1}{1+\lambda} F^* + \frac{\lambda}{1+\lambda} G^*$, \bar{d}_s

为 \bar{F}^* 的中位数。

$|a^*X - \bar{d}_s|$ 的分布记作 F_s , 密度 f_s ,

$|a^*Y - \bar{d}_s|$ 的分布记作 G_s , 密度 g_s ,

$\bar{F}_s \triangleq \frac{1}{1+\lambda} F_s + \frac{\lambda}{1+\lambda} G_s$, \bar{d}_s 为 \bar{F}_s 的中位数。则

$$1) \quad \sqrt{\frac{mn}{m+n}} T_{s,m}^s = \sqrt{\frac{mn}{m+n}} \left(\frac{G_s I_{\{a^*Y < 0\}}}{f(-a^*r)} - \frac{F_s I_{\{a^*X < 0\}}}{f(0)} \right) \frac{1}{\bar{d}_s} + o_p(1).$$

其中, G_s , F_s 分别是 G 和 F 的经验分布, \bar{F}^* , \bar{F}_s 满足 [6] 定理 3.2 条件。

$$2) \quad \text{令 } V(a) = \left(\frac{F^*(-a^*r)}{f(-a^*r)} - \frac{F^*(0)}{f^*(0)} \right) \frac{1}{\bar{d}_s}, \quad \|a\| = 1; \text{ 那么}$$

$$\begin{aligned} & \sqrt{\frac{mn}{m+n}} \left(\sup_{\|s\|=1} T_{s,m}^s - \sup_{\|s\|=1} V(a) \right) \\ & \sim \sup_{s \in B} \left(\sqrt{\frac{1}{1+\lambda}} E_G \frac{I_{\{a^*Y < 0\}}}{f(-a^*r)} - \sqrt{\frac{\lambda}{1+\lambda}} E_F \frac{I_{\{a^*X < 0\}}}{f^*(0)} \right) \frac{1}{\bar{d}_s}. \end{aligned}$$

其中, $B = \{a: \|a\| = 1, V(a) = \sup_{\|b\|=1} V(b)\}$, E_G , E_F 分别是 $\{I_{\{a^*Y < 0\}}: \|a\| = 1\}$, $\{I_{\{a^*X < 0\}}: \|a\| = 1\}$ 上的 G -桥与 F -桥。

证。由 [6] 定理 3.2 知, 欲证(1), 剩下只须证:

$$i) \quad \sup_{\|s\|=1} |\text{med}(a^*(X \cup Y)) - \bar{d}_s| \xrightarrow{n \rightarrow +\infty} 0,$$

ii) $\text{med}(|a^*X - \bar{d}_s| \cup |a^*Y - \bar{d}_s|) \rightarrow \bar{d}_s$, 对 $\|a\| = 1$ 一致, 只须证 ii), (i) 类似)首先令

$$\mathcal{F}_s = \{I_{\{|a^*X| < s\}} - I_{\{|a^*X| \geq s\}}: \|a\| = 1, s, t \in \mathbb{R}^1\}$$

$$\mathcal{F}_s = \{I_{\{|a^*X - \bar{d}_s| < s\}} - I_{\{|a^*X - \bar{d}_s| \geq s\}}: \|a\| = 1, s \in I_s^*\} \leq \mathcal{F}_s,$$

$$I_s^* = \left(\xi_s - \frac{\ln n}{n^{\frac{1}{2}}}, \xi_s + \frac{\ln n}{n^{\frac{1}{2}}} \right), \text{ 则有}$$

$$\sup_Q N_1(s, Q, \mathcal{F}_s) \leq A_0 s^{-\alpha_0}, \quad \forall n \geq 1.$$

$$\begin{aligned} \sup_{f \in \mathcal{F}_s} Pf^2 &= \sup_{\|s\|=1} \sup_{x \in I_s^*} |F_s(x) - F_s(\xi_s)| \\ &= \sup_{\|s\|=1} \sup_{x \in I_s^*} |F^*(x) - F^*(\xi_s) + F^*(-\xi_s) - F^*(-x)| \\ &\leq 2 \sup_{\|s\|=1} \sup_{x \in I_s^*} [f^*(x)Vf^*(-x)] \frac{\log n}{n^{\frac{1}{2}}} \\ &\leq \sup_{\|s\|=1} M(\xi_s) \cdot \frac{\log n}{n^{\frac{1}{2}}}, \quad f^*(x)Vf^*(-x) \triangleq \max\{f^*(x), f^*(-x)\}. \end{aligned}$$

f 在 ξ_s , $-\xi_s$ 附近有界时, $M(\xi_s)$ 有界。

令 $\delta_n^2 = \sup_{\|\alpha\|=1} M(\xi_\alpha) \cdot \frac{\log n}{n}$, $a_n^2 = \frac{1}{\sup_{\|\alpha\|=1} M(\xi_\alpha) n^{\frac{1}{2}}}$, 则由 [4] 中 p34, 定理 37 知, 存在

$L > 0$ 使

$$\sup_{\alpha \in I_n^d} |F_{\alpha, n} f - F_\alpha f| \leq L \delta_n^2 a_n = L (\sup_{\|\alpha\|=1} M(\xi_\alpha))^{\frac{1}{2}} \log n / n^{\frac{1}{2}}, \text{ a.s.}$$

也即 $\sup \{|[F_{\alpha, n}(x) - F_{\alpha, n}(\xi_\alpha)] - [F_\alpha(x) - F_\alpha(\xi_\alpha)]| : \|\alpha\|=1, x \in I_n^d\}$
 $= O(n^{-\frac{1}{2}} \log n)$ a.s.

同理 $\sup \{|[G_{\alpha, m}(x) - G_{\alpha, m}(\xi_\alpha)] - [G_\alpha(x) - G_\alpha(\xi_\alpha)]| : \|\alpha\|=1, x \in Z_n^d\}$
 $= O(m^{-\frac{1}{2}} \log m)$ a.s.

现设 $\bar{F}_{\alpha, n, m}$ 为 $|\alpha^T X_1 - \bar{d}_\alpha|, \dots, |\alpha^T X_n - \bar{d}_\alpha|, |\alpha^T Y_1 - \bar{d}_\alpha|, \dots, |\alpha^T Y_m - \bar{d}_\alpha|$ 的经验分布, 则

$$\begin{aligned} \bar{F}_{\alpha, n, m} &= \frac{1}{n+m} \left(\sum_{i=1}^n \delta_{(|\alpha^T X_i - \bar{d}_\alpha|)} + \sum_{j=1}^m \delta_{(|\alpha^T Y_j - \bar{d}_\alpha|)} \right) \\ &= \frac{1}{1+\lambda} F_{\alpha, n} + \frac{\lambda}{1+\lambda} G_{\alpha, m} \rightarrow \bar{F}_\alpha \text{ a.s. 对 } \|\alpha\|=1 \text{ 一致.} \end{aligned}$$

由此, $\sup_{\|\alpha\|=1} \sup_{x \in I_n^d} \left\{ \left| \frac{1}{1+\lambda} (F_{\alpha, n}(x) - F_{\alpha, n}(\xi_\alpha) - F_\alpha(x) + F_\alpha(\xi_\alpha)) \right. \right. \\ \left. \left. + \frac{\lambda}{1+\lambda} (G_{\alpha, m}(x) - G_{\alpha, m}(\xi_\alpha) - G_\alpha(x) + G_\alpha(\xi_\alpha)) \right| \right\} \\ = \sup_{\|\alpha\|=1} \sup_{x \in I_n^d} \{[\bar{F}_{\alpha, n, m}(x) - \bar{F}_{\alpha, n, m}(\xi_\alpha)] - [\bar{F}_\alpha(x) - \bar{F}_\alpha(\xi_\alpha)]\} \\ = O(n^{-\frac{1}{2}} \log n). \quad (3)$

令 $k_n = (n+m)p + o(n^{\frac{1}{2}} \log n)$. 下证

若设 $V_{\alpha, m}^*$ 为 $|\alpha^T X_1 - \bar{d}_\alpha|, \dots, |\alpha^T X_n - \bar{d}_\alpha|$ 的第 k_n 秩次统计量. 则 n 充分大时,

$$\forall \|\alpha\|=1, V_{\alpha, m}^* \in I_n^d.$$

实际上, 令 $a_n = \frac{\ln n}{\sqrt{n}}$, ξ_α 是 \bar{F}_α p 分位点.

$$\begin{aligned} \{\inf_{\alpha} (V_{\alpha, m}^* - \xi_\alpha + a_n) \leq 0\} &= \left\{ \sup_{\alpha} \bar{F}_{\alpha, n, m}(\xi_\alpha - a_n) \geq \frac{k_n}{n+m} \right\} \\ &\leq \left\{ \sup_{\|\alpha\|=1} \{|\bar{F}_{\alpha, n, m}(\xi_\alpha - a_n) - \bar{F}_\alpha(\xi_\alpha - a_n)|\} \geq r_n \right\} \\ &\leq \left\{ \sup_{\|\alpha\|=1} \sup_{x \in I_n^d} \{|\bar{F}_{\alpha, n, m}(x) - \bar{F}_{\alpha, n, m}(\xi_\alpha) - \bar{F}_\alpha(x) + \bar{F}_\alpha(\xi_\alpha)|\} \geq r_n \right\} \\ r_n &= \inf_{\|\alpha\|=1} \left(\frac{k_n}{n+m} - \bar{F}_\alpha(\xi_\alpha - a_n) \right) \geq \frac{1}{2} M \log n / n^{\frac{1}{2}}, M \text{ 是常数.} \end{aligned}$$

由(3)式知, $\inf_{\|\alpha\|=1} (V_{\alpha, m}^* - \xi_\alpha + a_n) > 0$ a.s., n 充分大.

同理 $\sup_{\|\alpha\|=1} (V_{\alpha, m}^* - \xi_\alpha + a_n) < 0$, a.s., n 充分大.

最后再证

$$V_{n,m}^* = \xi_n + [k_n - (n+m)\bar{F}_{n,m}(\xi_n)]/(n+m)\bar{f}_n(\xi_n) + O(n^{-\frac{1}{2}} \log n), \quad (4)$$

其中,

$$\bar{f}_n(\xi_n) = \frac{1}{1+\lambda} f_n(\xi_n) + \frac{\lambda}{1+\lambda} g_n(\xi_n).$$

实际上, 由 $\bar{F}_n(V_{n,m}^*) - \bar{F}_n(\xi_n) = \bar{f}_n(\xi_n)(V_{n,m}^* - \xi_n) + \bar{F}'_n(\theta_n^*)(V_{n,m}^* - \xi_n)^2$ 及 (3) 式和 $V_{n,m}^* \in I_n^*$, $\forall \|a\| = 1$, n 充分大, 马上得到(4)式. 特别有 ii) 成立.

推论 9 的第二部分的证明. 与前面方法类同, 故略.

推论 10. 两样本检验.

$H_0: F = G$, $H_1: F \neq G$, F , G 连续, $\frac{n_1}{n_1 + n_2} = \lambda$, $0 < \lambda < 1$.

$$V_{n_1, n_2}(a) = \int (F_{n_1}^*(t) - G_{n_1}^*(t))^2 dF^*(t),$$

$$V(a) = \int (F^*(t) - G^*(t))^2 dF^*(t), \text{ 则}$$

当 H_1 成立时,

$$\sqrt{\frac{n_1 n_2}{n_1 + n_2}} [\sup_{\|a\|=1} V_{n_1, n_2}(a) - \sup_{\|a\|=1} V(a)]$$

$$\rightsquigarrow 2 \sup_{a \in B} [\sqrt{1-\lambda} (E_F I_{\{a^T X \leq 0\}}) - \sqrt{\lambda} (E_G I_{\{a^T Y \leq 0\}})]$$

$$\cdot (F^*(t) - G^*(t)) dF^*(t)$$

其中 E_F , E_G 分别是 $\{I_{\{a^T X \leq 0\}}: \|a\| = 1\}$, $\{I_{\{a^T Y \leq 0\}}: \|a\| = 1\}$ 上的 F -桥和 G -桥.

注记: 1) 正态性检验

$$\delta_n = \sup_{a \in B} \left| F_n^*(t) - \Phi \left(\frac{t - a^T \mu_n}{\sqrt{a^T V_n a}} \right) \right|,$$

其中 μ_n 是均值 μ 的样本均值估计, V_n 是协方差 V 的样本协方差估计, Φ 是标准正态分布函数, $F_n^*(t) = \frac{1}{n} \sum_{i=1}^n I_{\{a^T X_i \leq 0\}}$. δ_n 的功效也有与前面类同的结论, 这里不再赘述.

2) M -型 PP 位置检验.

T_n^* 是 $\int \phi(a^T X, t) dP_n = 0$ 的最小解,

$T^* \triangleq T^*(P)$ 是 $\int \phi(a^T X, t) dP = 0$ 的最小解;

如果 P_n 是 G 的经验分布, $G \neq P$, $\sup_{\|a\|=1} |T^*(G) - T^*(P)| > 0$. 则有:

$$\sqrt{n} \left\{ \sup_{\|a\|=1} |T_n^* - T^*(P)| - \sup_{\|a\|=1} |T^*(G) - T^*(P)| \right.$$

$$\left. \rightsquigarrow \sup_{a \in B} E_G h(a) \cdot \text{sgn}(T^*(G) - T^*(P)), \right.$$

其中 $B = \{a: \|a\| = 1, \sup_{\|b\|=1} |T^b(G) - T^b(P)| = |T^*(G) - T^*(P)|\}$,

$$\text{而 } h(a) \triangleq \frac{\phi(a^T X, T^*(G))}{-G \frac{\partial \phi(a^T X, T^*(G))}{\partial t_2}}, E_G \text{ 是 } \{h(a): \|a\| = 1\} \text{ 上的 } G\text{-桥}$$

证。本推论的证明也类同前面的方法，故略。

参 考 文 献

- [1] Dudley, R. M., Philipp, W., Invariance principles for sums of Banach space valued random elements and Empirical processes, *Z. W. verw. Geb.*, 62 (1983), 509—552.
- [2] Gideon, R. A., Bentice, M. J., Pyke, R. (1989), The Limiting distribution of the rank correlation coefficient R_g . (to appear in a Festschrift for Ingram Olkin, Springer).
- [3] Li, G., Chen, Z., Projection-Pursuit approach to robust dispersion matrices and principal components: primary theory and Monte Carlo, *JASA*, 80(1985), 759—766.
- [4] Pollard, D., *Convergence of Stochastic Processes*, Springer-Verlag, New York, 1984.
- [5] 张 航, PP 型拟合优度检验, 系统科学与数学, 8:3(1988), 234—242.
- [6] 张 航, 几类 PP 型检验统计量的性质, 应用数学学报, 12:1(1989), 82—95.
- [7] 张 健, 朱力行, 成平 M-型 PP 指标其特征值和特征向量的渐近理论及其应用(I)(II), 投《中国科学》, 1988.

ASYMPTOTIC POWERS OF SOME PP TESTS

ZHANG JIAN CHENG PING

(Institute of Systems Science, Academia Sinica)

ABSTRACT

A general method of dealing with the asymptotic powers of the PP tests is given and by this method asymptotic powers of a number of useful PP tests and estimators are obtained. In particular, the asymptotic distribution of the largest eigenvalue of M-type PP index is derived.