

# Two-Samples Problem について

田 村 亮 二

## § 1. 序

Distribution-free test といわれる統計的推測理論の一つの型である two-samples problem について exceedence number を利用した検定の方法を論ずる。此の問題に対して rank を使つた Wilcoxon [1], Mann and Whitney [2] の test 及び exceedence number による Rosenbaum [3] [4], Epstein [5] がある。又 Mood [6] は large sample で median test を提唱している。

今二つの分布函数  $F(x)$ ,  $G(x)$  からの random sample を夫々  $x_1, x_2, \dots, x_n; y_1, y_2, \dots, y_m$  とし帰無仮説  $H_0: F(x) = G(x)$  を適当な対立仮説に対して上記 Sample に基いて検定するといふ問題で  $F(x)$ ,  $G(x)$  の型は未知である。さて  $x_i$  の小さい方から大きさの順に並びかえ夫々  $r$  番目から  $s$  番目を改めて  $x_r, x_s$  で表す ( $r < s$  とす)。そして  $(x_r, x_s)$  に含まれる  $y$  の個数を  $U$  とするとき、この random variable  $U$  が検定基準として採用されるものである。分布函数を連続型と仮定すれば  $P_r(y_i = x_r) = 0$ ,  $P_r(y_i = x_s) = 0$  であるから確率 1 で  $x_r, x_s$  に等しくなる  $y$  はないと考えてよい。上の  $U$  による検定を便宜上  $U$ -test と名付ける。正整数  $r, s$  ( $< n$ ) は一般に任意であるが  $n$  が余り小さくない時 (例えば  $n > 5$  又は 6) は  $r=1$ ,  $s=n$  の如き両端をとらないのが妥当と考えられる。殊に寿命試験等の如く Sample 全部の検査が終了するまで待たず中途打切をする場合は時間に応じて適宜  $r$  及び  $s$  を定めればよい。普通両端の一つ又は二つの outlying observations を除いたもの即  $r=2$ ,  $s=n-1$  又は  $r=3$ ,  $s=n-2$  等を用ふるのが安全と考えられる。

§ 2. で  $U$ -statistics の確率分布, moment 等  $U$  の性質をしらべ、§ 3 で  $U$ -test の consistency を証明し、§ 4 で検定方法及び表を与える。

## § 2. $U$ の確率分布

仮説  $H_0$  の下で  $x_r, x_s$  の確率密度は

$$\frac{n!}{(r-1)!(s-r-1)!(n-s)!} \{F(x_r)\}^{r-1} \{F(x_s) - F(x_r)\}^{s-r-1} \{1 - F(x_s)\}^{n-s} dF(x_s) dF(x_r)$$

で表わされる。故に

$$P(x) = P_r(U = \lambda) = C \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} F_1^{r-1} (F_2 - F_1)^{s-r-1} (1 - F_2)^{n-s} (F_2 - F_1)^\lambda (1 - F_2 + F_1)^{m-\lambda} dF_1 dF_2$$

但  $F_1 = F(x_r)$ ,  $F_2 = F(x_s)$

$$C = \binom{m}{\lambda} \frac{n!}{(r-1)!(s-r-1)!(n-s)!}$$

$F_1 = t_1, F_2 = t_2$  から

$$P(\lambda) = C \int_0^1 \int_0^{1-t_1} t_2^{r-1} (t_2-t_1)^{\lambda+s-r-1} (1-t_2)^{n-s} (1-t_2+t_1)^{m-\lambda} dt_1 dt_2$$

$$0 < t_1 < t_2 < 1$$

之は Dirichlet integral の一種で  $t_2 - t_1 = u, t_1 = v$  で

$$= c \int_0^{1-u} \int_0^v v^{r-1} u^{\lambda+s-r-1} (1-u)^{m-\lambda} (1-u-v)^{n-s} dudv$$

$$\int_0^{1-u} v^{r-1} (1-u-v)^{n-s} dv = \frac{(r-1)!}{(n-s+1)(n-s+2)\cdots(n-s+r)} (1-u)^{n-s+r}$$

を利用すれば

$$P(\lambda) = c \cdot \frac{(r-1)!(n-s)!}{(n-s+r)!} B(\lambda+s-r, m+n-\lambda-s+r+1) = \frac{\binom{m}{\lambda} \binom{n}{s-r} \frac{s-r}{\lambda+s-r}}{\binom{m+n}{\lambda+s-r}}$$

こゝに注意すべきは  $P(\lambda)$  は  $r, s$  の各々に依存せず  $s-r$  ( $=t$  とおく) にのみ関係してゐることである。

次に  $U$  の moment を求めるため先ず  $\sum_{\lambda=0}^m P(\lambda) = 1$  を証明する。

$$\sum_{\lambda=0}^m P(\lambda) = t \binom{n}{t} \sum \frac{\binom{m}{\lambda}}{(t+\lambda) \binom{m+n}{\lambda+t}} = \frac{1}{\binom{m+n}{m}} \sum \binom{\lambda+t-1}{t-1} \binom{m+n-\lambda-t}{n-t}$$

さて等式  $(1-x)^{-(n-t)-1} (1-x)^{-(t-1)-1} = (1-x)^{-n-1}$  で  $x$  の係数を比較して

$$\sum_{\lambda=0}^m \binom{-n+t-1}{m-\lambda} \binom{-t}{\lambda} = \binom{-n-1}{m}$$

即ち  $\sum \binom{m+n-\lambda-t}{m-\lambda} \binom{t+\lambda-1}{\lambda} = \binom{m+n}{m}$  これを代入することにより  $\sum P(\lambda) = 1$  をうる。

後で使ふため上式を次のやうにかきかえる。

$$\sum_{\lambda=0}^m \frac{\binom{m}{\lambda}}{\binom{m+n-1}{\lambda+t-1}} = \frac{m+n}{t} \binom{m}{t} \quad (*)$$

a 次の factorial moment  $E\{U(U-1)\cdots(U-a+1)\}$  を  $\alpha_{(a)}$  で表すと

$$\alpha_{(a)} = \sum_{\lambda=0}^m \lambda^{(a)} P(\lambda) = \frac{m}{m+n} \sum_{\lambda=0}^m \frac{\binom{m-a}{\lambda-a}}{\binom{m+n-1}{\lambda+t-1}}$$

(\*) を利用して

$$\alpha_{(a)} = \frac{\binom{m}{(a)} \binom{t+a-1}{(a)}}{\binom{n+a}{(a)}} \quad (**)$$

(\*) より  $U$  の平均及び分散は次式で求められる。

$$E(U) = \frac{m}{n+1} t, \quad V(U) = \sigma^2 = \frac{m(m-1)}{(n+1)(n+2)} t(t+1) + \frac{m}{n+1} t - \left(\frac{m}{n+1} t\right)^2$$

特に  $m, n$  が大きくて  $m \approx n$  と考えてもよいとき (又は  $m, n \rightarrow \infty$  の極限に於いて)

$$E(U) = t, \quad \sigma^2 = 2t$$

となり  $U$  の確率分布は自由度  $t$  の  $X^2$  - 分布と見做してもよい。

### § 3. U-test の一貫性

U-test の棄却域は対立仮説によつて変つてくるが例えば  $H_1: F(x) > G(x)$  (for every  $x$ ) に対する  $H_0$  の test では有意水準  $\varepsilon$  をきめた時  $U \leq U_\varepsilon$  である。  $U_\varepsilon$  は次節の表から求められる。或は棄却域は  $U \leq \frac{m}{n+1} t - \beta(n, m) \sigma$  ( $\beta(n, m)$  は  $n, m$  の値に応じて定まり  $\lim_{n, m \rightarrow \infty} \beta(n, m) = \beta(\text{const})$ ) と考えてもよし。 U-test の power は他の distribution - free test の場合と同様計算できないが一貫性については次のやうに証明できる。

$H_1$  の下での expectation を  $E_{H_1}$  で表すものとし確率変数  $Z_i$  を  $Z_i = 1$ , for  $x_r < y_i < x_s$

$Z_i = 0$ , for  $y_i > x$  又は  $y_i < x_r$ . で定義すれば  $Z_1, Z_2, \dots, Z_m$  は independent. 且  $\sum_{i=1}^m Z_i = U$

$$\text{又 } E_{H_1}(Z_i) = P_r(x_r < y_i < x_s) = \frac{n!}{(r-1)!(s-r-1)!(n-s)!} \iiint_{x_r < y_i < x_s} F^{r-1}(x_r) \{F(x_s) - F(x_r)\}^{s-r-1}$$

$$\times \{1 - F(x_s)\}^{n-s} dF(x_r) dF(x_s) dG(y_i)$$

$$= c_1 \iint_{x_r < x_s} F_1^{r-1} (F_2 - F_1)^{s-r-1} (1 - F_2)^{n-s} (G_2 - G_1) dF_1 dF_2$$

$H_1$  の仮定から

$$< c_1 \iint_{x_r < x_s} F_1^{r-1} (F_2 - F_1)^{s-r-1} (1 - F_2)^{n-s} F_2 dF_1 dF_2$$

§ 2 の momentf の計算と同様にして

$$E_{H_1}(Z_i) < \frac{t}{n+1} + \frac{r}{n+1} \quad \text{をうる。}$$

同様にして  $E_{H_1}(Z_i^2) = P_r(x < y_i < x_s) = E_{H_1}(Z_i)$

今  $E_{H_1}(Z_i) = \frac{t}{n+1} + \frac{r}{n+1} - \lambda$  ( $\lambda > 0$ ) とおくと

$$E_{H_1}(U) = \frac{m t}{n+1} + \frac{m r}{n+1} - m \lambda, \quad V_{H_1}(U) = m \left\{ \frac{t}{n+1} + \frac{r}{n+1} - \lambda - \left( \frac{t}{n+1} + \frac{r}{n+1} - t \right)^2 \right\}$$

一方 power function  $P(H_1) = P_r(U < \frac{m t}{n+1} - \beta(n, m) \sigma | H_1)$  は

$$\begin{aligned} P(H_1) &= P_r \left\{ U - \left( \frac{m t}{n+1} + \frac{m r}{n+1} - m \lambda \right) < m \lambda - \frac{m r}{n+1} - \beta(n, m) \sigma | H_1 \right\} \\ &= P_r(U - E_{H_1}(U) < k \sigma_{H_1}) \quad \text{但 } k = \frac{m \lambda - \frac{m r}{n+1} - \beta(n, m) \sigma}{\sigma_{H_1}} \end{aligned}$$

と変形され、更に Tchebycheff の定理を利用すれば

$$P(H_1) \geq 1 - \frac{\left\{ m \frac{t}{n+1} - \lambda - \left( \frac{t}{n+1} - \lambda \right)^2 \right\}}{\left\{ n \lambda - \frac{m r}{n+1} - \beta(n, m) \sigma \right\}^2} \quad (\sigma \text{ は } \S 2. \text{ 参照})$$

こゝで  $n, m \rightarrow \infty$  とすれば右辺の第二項は 0 に収斂し  $\lim_{n, m \rightarrow \infty} P(H_1) = 1$  をうる。

#### § 4. 検定方法及び表

対立仮説として考えられるものは location に関しては (i)  $F(x) > G(x)$  (ii)  $F(x) < G(x)$  であり dispersion については (iii)  $d(F(x)) > d(G(x))$  (iv)  $d(F(x)) < d(G(x))$  である。 $d(F(x))$  は  $F(x)$  の dispersion measure とする。

上の型の検定は實際上屢々出会ふもので例えば在来の結果に対してある処置をした時の効果を問題にすると dispersion に変化は余らないといふ apriori information があるとき location の増減についての検査は (i), (ii), であり逆の場合は (iii), (iv), である。(i) 及 (ii) の型に対しては適当な  $r, s$  を定め  $U$  を勘定して表から得られる  $U_\varepsilon$  より小さいならば  $H_0$  を棄却する。(iii), (iv) の型の対立仮説に関しては sample から勘定された  $U$  が表より得られた  $U_{\varepsilon_1}$  より小さいか  $U_{\varepsilon_2}$  より大きいならば  $H_0$  を棄却するといふ立場に立てばよい。



$t$	$\frac{m}{\lambda}$	5	6	7	8	9	10
4	0		0.030	0.034	0.036	0.041	0.043
	1		0.090	0.097	0.102	0.105	0.108
	2		0.102	0.163	0.163	0.162	0.162
	3		0.211	0.203	0.195	0.190	0.185
	4		0.227	0.203	0.190	0.181	0.175
	5		0.181	0.163	0.152	0.145	0.140
	6		0.090	0.097	0.097	0.120	0.095
	7			0.034	0.046	0.081	0.054
	8				0.012	0.040	0.025
	9					0.011	0.008
10						0.001	
5	0			0.010	0.012	0.014	0.015
	1			0.040	0.046	0.050	0.054
	2			0.091	0.097	0.101	0.104
	3			0.151	0.152	0.151	0.150
	4			0.203	0.190	0.181	0.175
	5			0.220	0.195	0.181	0.171
	6			0.183	0.162	0.151	0.143
	7			0.096	0.102	0.101	0.100
	8				0.038	0.050	0.056
	9					0.014	0.023
10						0.005	
6	0				0.003	0.004	0.005
	1				0.016	0.020	0.023
	2				0.045	0.051	0.056
	3				0.090	0.096	0.100
	4				0.146	0.145	0.143
	5				0.195	0.181	0.171
	6				0.215	0.190	0.175
	7				0.184	0.162	0.150
	8				0.100	0.105	0.104
	9					0.041	0.054
10						0.015	

t	$\frac{m}{\lambda}$	8	9	10	t	9	10
7	0		0.001	0.001	8		0.000
	1		0.006	0.008			0.002
	2		0.020	0.025			0.008
	3		0.048	0.054			0.023
	4		0.090	0.095			0.050
	5		0.142	0.140			0.090
	6		0.190	0.175			0.139
	7		0.211	0.185			0.185
	8		0.185	0.162			0.208
	9		0.102	0.108			0.185
	10			0.043			0.105

## 文 献

- [1] F. Wilcoxon : Individual comparisons by ranking methods.  
Biometrics Bull. Vol. 1 (1945)
- [2] H. B. Mann and D. R. whitney : On a test of whether one of two random variables is stochastically larger than the other.  
Ann. Math. Statis. Vol 18 (1947)
- [3] S. Rosenbaum : Tables for a nonparsmetric test of dispersion  
Ann. Math. Statis. Vol 24 (1953)
- [4] S. Rosenbeum : Tables for a nonparametric test of location.  
Ann. Math. statis. Vol 25 (1954)
- [5] B. Epstein : Tables for the distribution of the number of exceedances.  
Ann. Matq. Statis. Vol 25 (1954)
- [6] A. M. Mood : Introduction to the theory of statistics. (1950)