

Basic Nonparametric Statistical Methods

H. Wakaki

January 25, 2016

Contents

1	Introduction	1
1.1	Binomial test	1
1.2	Point estimation	2
1.3	Interval estimation	2
2	The One-Sample Location Problem	3
2.1	Paired replicates analyses by way of signed ranks	3
2.1.1	A distribution-free signed rank test (Wilcoxon)	4
2.1.2	An estimator associated with Wilcoxon's signed rank statistic (Hodges-Lehmann)	7
2.1.3	A distribution-free confidence interval based on Wilcoxon's signed rank test (Tukey)	7
2.2	Paired replicates analyses by way of signs	8
2.2.1	A distribution-free signed test (Fisher)	8
2.2.2	An estimator associated with the sign statistic(Hodges-Lehmann) .	9
2.2.3	A distribution-free confidential interval based on the sign test (Thomp- son, Savur)	10
2.3	One-sample data	10
2.3.1	A distribution-free signed rank test	10
2.3.2	An estimator based on the signed rank test	11
2.3.3	A distribution free confidential interval based on the signed rank test	11
2.4	Procedures based on the sign statistic	11
3	The Two-Sample Location Problem	11
3.1	A distribution-free rank sum test (Wilcoxon, Mann and Whitney)	12
3.2	An estimator associated with Wilcoxon's rank sum statistic (Hodges-Lehmann)	16
3.3	A distribution-free confidence interval based on Wilcoxon's rank sum test (Moses)	16
4	The Two-Sample Dispersion Problem and Other Two-Sample Prob- lems	17
4.1	A distribution-free rank test for dispersion-case of equal medians (Ansari- Bradley)	17

4.2	An asymptotically distribution-free test for dispersion based on the Jackknife-medians not necessarily equal (Miller)	24
4.3	A distribution-free rank test for either location or dispersion (Lepage) . . .	28
4.4	A distribution-free test for general differences in two populations (Kolmogorov-Smirnov)	30
5	The One-Way Layout	31
5.1	A distribution-free test for general alternatives (Kruskal-Wallis)	32
5.2	A distribution-free test for ordered alternatives (Jonckheere-Terpstra)	36

Reference

“NONPARAMETRIC STATISTICAL METHODS second edition”
by Myles Hollander and Douglas A. Wolfe

1 Introduction

Thema Inference on Bernoulli trial

Data X_1, \dots, X_n

Assumptions

A1. X_1, \dots, X_n : independent

A2. $\Pr\{X_i = 1\} = 1 - \Pr\{X_i = 0\} = p$ ($0 < p < 1$)

1.1 Binomial test

Null hypothesis

$$H_0 : p = p_0 \quad (\text{a specified number})$$

Test statistic

$$T = \sum_{i=1}^n X_i$$

Null distribution Under H_0

$$T \sim B(n, p_0),$$

$$\Pr\{T = k\} = q_0(k) := \frac{n!}{k!(n-k)!} p_0^k (1-p_0)^{n-k} \quad (k = 0, 1, \dots, n)$$

a. One-sided upper-tail test

$$H_0 : p = p_0 \text{ versus } H_1 : p > p_0$$

$$T \geq b_\alpha \Rightarrow \text{Reject } H_0, \quad b_\alpha = \min \left\{ k : \sum_{j=k}^n q_0(j) \leq \alpha \right\}$$

Randomized test

In the case of $\alpha_0 = \Pr\{T \geq b_\alpha\} < \alpha$.

Let U : uniform random variable on $(0, 1)$.

$$T \geq b_\alpha \Rightarrow \text{Reject}$$

or

$$T = b_\alpha - 1 \text{ and } U < \frac{\alpha - \alpha_0}{q_0(b_\alpha - 1)} \Rightarrow \text{Reject}$$

The P -value

Let t : observed value

$\Pr\{B \geq t\}$ is called “ P -value”, $B \sim B(n, p_0)$

b. One-sided lower-tail test

$H_0 : p = p_0$ versus $H_2 : p < p_0$

$$T \leq c_\alpha \Rightarrow \text{Reject } H_0, \quad c_\alpha = \max\left\{k : \sum_{j=0}^k q_0(k) \leq \alpha\right\}$$

c. Two-sided test

$H_0 : p = p_0$ versus $H_3 : p \neq p_0$

$T \geq b_{\alpha_1}$ or $T \leq c_{\alpha_2} \Rightarrow \text{Reject } H_0, \quad \alpha_1 + \alpha_2 = \alpha$

Large-sample approximation

$$T^* = \frac{T - np_0}{\sqrt{np_0(1-p_0)}} \xrightarrow{d} N(0, 1)$$

$T^* \geq z_\alpha \Rightarrow \text{Reject (one-sided upper-tail test)}$

$T^* \leq -z_\alpha \Rightarrow \text{Reject (one-sided lower-tail test)}$

$|T^*| \geq z_{\alpha/2} \Rightarrow \text{Reject (Two-sided test)}$

z_α : upper 100α percent point of $N(0, 1)$

1.2 Point estimation

Estimator

$$\hat{p} = \frac{1}{n}T = \frac{1}{n}\#\{i : X_i = 1, 1 \leq i \leq n\}$$

Properties

$E[\hat{p}] = p$ (unbiased estimator)

$$\text{Var}[\hat{p}] = \frac{1}{n}p(1-p)$$

$$\frac{\sqrt{n}(\hat{p} - p)}{\sqrt{p(1-p)}} \xrightarrow{d} N(0, 1) \quad (n \rightarrow \infty)$$

1.3 Interval estimation

Define

$$\{p_L(k; \alpha)\}_{k=0,1,\dots,n} \quad \text{and} \quad \{p_U(k; \alpha)\}_{k=0,1,\dots,n}$$

such that

$$\begin{aligned}
 p_L(k; \alpha) &< p_U(k; \alpha) \quad (k = 0, 1, \dots, n) \\
 p_L(k; \alpha) &< p_L(k+1; \alpha), \quad p_U(k; \alpha) < p_U(k+1; \alpha) \quad (k = 0, 1, \dots, n-1) \\
 \sum_{k: p_L(k; \alpha) \leq p \leq p_U(k; \alpha)} \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k} &\geq 1 - \alpha \text{ for all } p \in (0, 1).
 \end{aligned}$$

\Rightarrow

$$\Pr\{p_L(T; \alpha) \leq p \leq p_U(T; \alpha)\} \geq 1 - \alpha \quad \text{for all } p$$

Confidence interval

$$[p_L(T; \alpha), p_U(T; \alpha)] : 1 - \alpha \text{ confidence interval for } p$$

Large-sample approximation

$$p_L(T; \alpha) = \hat{p} - z_{\alpha/2} \left[\frac{\hat{p}(1-\hat{p})}{n} \right]^{1/2}, \quad p_U(T; \alpha) = \hat{p} + z_{\alpha/2} \left[\frac{\hat{p}(1-\hat{p})}{n} \right]^{1/2}$$

\Rightarrow

$$\lim_{n \rightarrow \infty} \Pr\{p_L(T; \alpha) \leq p \leq p_U(T; \alpha)\} = 1 - \alpha$$

2 The One-Sample Location Problem

Thema Inferences on the location

Two types of data

- paired replicates data — “pretreatment” and “posttreatment”
- one-sample data

2.1 Paired replicates analyses by way of signed ranks

Data

Subject i	X_i	Y_i
1	X_1	Y_1
2	X_2	Y_2
\vdots	\vdots	\vdots
n	X_n	Y_n

Assumptions

A1. Let $Z_i = Y_i - X_i$ ($i = 1, \dots, n$).

Z_1, \dots, Z_n : independent

A2. Let $F_i(t) = \Pr\{Z_i \leq t\}$ ($i = 1, \dots, n$).

F_i is continuous and

$$F_i(\theta + t) + F_i(\theta - t) = 1, \text{ for every } t \quad (i = 1, \dots, n)$$

Note θ is the median which is referred as the *treatment effect*

2.1.1 A distribution-free signed rank test (Wilcoxon)

Null hypothesis

$$H_0 : \theta = 0$$

Test statistic

$$\psi_i = \begin{cases} 1, & (Z_i > 0) \\ 0, & (Z_i < 0) \end{cases}$$

R_i : the rank of $|Z_1|, \dots, |Z_n|$ ordered from least to greatest

$$\Rightarrow T^+ = \sum_{i=1}^n R_i \psi_i : \text{Wilcoxon signed rank statistic}$$

Null distribution Let

V_1, \dots, V_n : independent random variables

$$\Pr\{V_i = i\} = \Pr\{V_i = 0\} = \frac{1}{2} \quad (i = 1, \dots, n)$$

$$\Rightarrow T^+ \sim \sum_{i=1}^n V_i \quad \text{same distribution} \Rightarrow \text{distribution free (not depend on } F_i)$$

\therefore

$\theta = 0 \Rightarrow |Z_1|, \dots, |Z_n|, \psi_1, \dots, \psi_n$: independent

Define I_1, \dots, I_n by

$$I_j = i \Leftrightarrow R_i = j \quad (j = 1, \dots, n).$$

$\Rightarrow (I_1, \dots, I_n), \psi_1, \dots, \psi_n$: independent

$b_j = 0$ or 1 ($j = 1, \dots, n$).

S_n : set of all permutations of $(1, 2, \dots, n)$.

$$\begin{aligned}
& \Pr\{\psi_{I_1} = b_1, \dots, \psi_{I_n} = b_n\} \\
&= \sum_{(s_1, \dots, s_n) \in S_n} \Pr\{\psi_{I_1} = b_1, \dots, \psi_{I_n} = b_n, (I_1, \dots, I_n) = (s_1, \dots, s_n)\} \\
&= \sum_{(s_1, \dots, s_n) \in S_n} \Pr\{\psi_{s_1} = b_1, \dots, \psi_{s_n} = b_n, (I_1, \dots, I_n) = (s_1, \dots, s_n)\} \\
&= \sum_{(s_1, \dots, s_n) \in S_n} \Pr\{(I_1, \dots, I_n) = (s_1, \dots, s_n)\} \left(\prod_{i=1}^n \Pr\{\psi_{s_i} = b_i\} \right) \\
&= \sum_{(s_1, \dots, s_n) \in S_n} \Pr\{(I_1, \dots, I_n) = (s_1, \dots, s_n)\} \left(\frac{1}{2} \right)^n = \left(\frac{1}{2} \right)^n \\
&\Rightarrow \\
&\psi_{I_1}, \dots, \psi_{I_n} : \text{independent,} \\
&\Pr\{\psi_{I_j} = 0\} = \Pr\{\psi_{I_j} = 1\} = \frac{1}{2} \quad (j = 1, \dots, n)
\end{aligned}$$

Define

$$V_j = \psi_{I_j} j \quad (j = 1, \dots, n).$$

Then V_1, \dots, V_n : independent, $\Pr\{V_j = 0\} = \Pr\{V_j = j\} = \frac{1}{2}$,

$$T^+ = \sum_{j=1}^n V_j,$$

Three types of test

a. One-sided upper-tail test

$$H_0 : \theta = 0 \text{ versus } H_1 : \theta > 0$$

$$T^+ \geq t_\alpha \Rightarrow \text{Reject } H_0, \quad t_\alpha = \min \left\{ k : \Pr \left\{ \sum_{i=1}^n V_i \geq k \right\} \leq \alpha \right\}$$

b. One-sided lower-tail test

$$H_0 : \theta = 0 \text{ versus } H_2 : \theta < 0$$

$$T^+ \leq \frac{n(n+1)}{2} - t_\alpha \Rightarrow \text{Reject } H_0,$$

c. Two-sided test

$$H_0 : \theta = 0 \text{ versus } H_3 : \theta \neq 0$$

$$T^+ \geq t_{\alpha/2} \text{ or } T^+ \leq \frac{n(n+1)}{2} - t_{\alpha/2} \Rightarrow \text{Reject } H_0,$$

Large-sample approximation

Under H_0

$$E[T^+] = \frac{n(n+1)}{4}, \quad \text{Var}[T^+] = \frac{n(n+1)(2n+1)}{24}$$

$$T^* = \frac{T^+ - E[T^+]}{\sqrt{\text{Var}[T^+]}} \xrightarrow{d} N(0, 1) \quad (n \rightarrow \infty)$$

$T^* \geq z_\alpha \Rightarrow$ Reject (one-sided upper-tail test)

$T^* \leq -z_\alpha \Rightarrow$ Reject (one-sided lower-tail test)

$|T^*| \geq z_{\alpha/2} \Rightarrow$ Reject (Two-sided test)

z_α : upper 100α percent point of $N(0, 1)$

∴

Lemma (Lindeberg)

$\{X_n\}_{n=1,2,\dots}$: sequence of independent random variables

$E[X_n] = 0, E[X_n^2] = \sigma_n^2 < \infty$

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^2} \sum_{j=1}^n \int_{|x| > \varepsilon s_n} x^2 dF_j(x) = 0 \quad \text{for all } \varepsilon > 0$$

$$\frac{\sum_{j=1}^n X_n}{s_n} \xrightarrow{d} N(0, 1),$$

where

$$s_n^2 = \sum_{j=1}^n \sigma_n^2$$

Let $X_j = V_j - E[V_j]$. Then $\Pr(X_j = -\frac{j}{2}) = \Pr(X_j = \frac{j}{2}) = \frac{1}{2}$.

$$E[V_j] = \frac{j}{2}, \quad \sigma_j^2 = \text{Var}[V_j] = \frac{j^2}{2} - \left(\frac{j}{2}\right)^2 = \frac{j^2}{4}$$

$$s_n^2 = \text{Var}[T^+] = \frac{n(n+1)(2n+1)}{24}$$

$$\left| \frac{\frac{j}{2}}{s_n} \right| \leq \frac{n}{2s_n} = \sqrt{\frac{6n}{(n+1)(2n+1)}} \rightarrow 0 \quad (n \rightarrow \infty)$$

Hence

$$\exists n_\varepsilon \text{ s.t. } n > n_\varepsilon \Rightarrow \int_{|x| > \varepsilon s_n} x^2 dF_j(x) = 0,$$

which leads that the Lindeberg condition holds.

Testing equality to the specified value

$$H_0 : \theta = \theta_0 \quad (\text{specified value})$$

Test statistic

$$\psi_i = \begin{cases} 1, & (Z_i > \theta_0) \\ 0, & (Z_i < \theta_0) \end{cases}$$

R_i : the rank of $|Z_1 - \theta_0|, \dots, |Z_n - \theta_0|$ ordered from least to greatest

$$\Rightarrow T^+ = \sum_{i=1}^n R_i \psi_i : \text{Wilcoxon signed rank statistic}$$

2.1.2 An estimator associated with Wilcoxon's signed rank statistic (Hodges–Lehmann)

Derivation Choose the value of θ_0 such that $H_0 : \theta = \theta_0$ is least significant, that is

$$\hat{\theta} = \underset{\theta_0}{\operatorname{argmin}} \left| T_0^+ - \frac{n(n+1)}{4} \right|,$$

T_0^+ : test statistic for testing “ $\theta = \theta_0$ ”

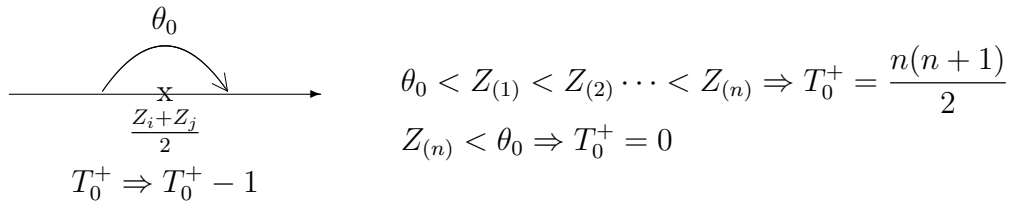
Estimator

$$\hat{\theta} = \operatorname{median} \left\{ \frac{Z_i + Z_j}{2}, 1 \leq i \leq j \leq n \right\}$$

$W^{(1)} < \dots < W^{(M)}$: ordered values of $\frac{Z_i + Z_j}{2}$

$$M = \frac{n(n+1)}{2} = 2k + 1 \Rightarrow \hat{\theta} = W^{(k+1)}$$

$$M = 2k \Rightarrow \hat{\theta} = \frac{W^{(k)} + W^{(k+1)}}{2}$$



2.1.3 A distribution-free confidence interval based on Wilcoxon's signed rank test (Tukey)

Derivation Collect θ_0 such that the two sided test for

$$H_0 : \theta = \theta_0 \text{ versus } H_3 : \theta \neq \theta_0$$

does not reject.

Confidence interval

$$C_\alpha = \frac{n(n+1)}{2} + 1 - t_{\alpha/2},$$

$$\theta_L = W^{(C_\alpha)}, \quad \theta_U = W^{M+1-C_\alpha} = W^{(t_{\alpha/2})}$$

\Rightarrow

$$\Pr_\theta \{ \theta_L < \theta < \theta_U \} = 1 - \alpha \text{ for all } \theta$$

Large-sample approximation

$$C_\alpha \approx \frac{n(n+1)}{4} - z_{\alpha/2} \left\{ \frac{n(n+1)(2n+1)}{24} \right\}^{1/2}$$

2.2 Paired replicates analyses by way of signs

Data

Subject i	X_i	Y_i
1	X_1	Y_1
2	X_2	Y_2
\vdots	\vdots	\vdots
n	X_n	Y_n

Assumptions

B1. Let $Z_i = Y_i - X_i$ ($i = 1, \dots, n$).

Z_1, \dots, Z_n : independent

B2.

$$\Pr\{Z_i < \theta\} = \Pr\{Z_i > \theta\} = \frac{1}{2} \quad (n = 1, \dots, n)$$

2.2.1 A distribution-free signed test (Fisher)

Null hypothesis

$$H_0 : \theta = 0$$

Test statistic

$$\psi_i = \begin{cases} 1, & (Z_i > 0) \\ 0, & (Z_i < 0) \end{cases}$$

$$\Rightarrow B = \sum_{i=1}^n \psi_i : \text{Fisher's signed statistic}$$

Null distribution

$$B \sim B\left(n, \frac{1}{2}\right)$$

Three types of test

a. One-sided upper-tail test

$$H_0 : \theta = 0 \text{ versus } H_1 : \theta > 0$$

$$B \geq b_\alpha \Rightarrow \text{Reject } H_0, \quad b_\alpha = \min \left\{ k : \sum_{i=k}^n \frac{n!}{i!(n-i)!} \left(\frac{1}{2}\right)^n \leq \alpha \right\}$$

b. One-sided lower-tail test

$$H_0 : \theta = 0 \text{ versus } H_2 : \theta < 0$$

$$B \leq n - b_\alpha \Rightarrow \text{Reject } H_0,$$

c. Two-sided test

$$H_0 : \theta = 0 \text{ versus } H_3 : \theta \neq 0$$

$$B \geq b_{\alpha/2} \text{ or } B \leq n - b_{\alpha/2} \Rightarrow \text{Reject } H_0,$$

Large-sample approximation

Under H_0

$$E[B] = \frac{n}{2}, \quad \text{Var}[B] = \frac{n}{4}$$

$$B^* = \frac{B - E[B]}{\sqrt{\text{Var}[B]}} \xrightarrow{d} N(0, 1) \quad (n \rightarrow \infty)$$

$$B^* \geq z_\alpha \Rightarrow \text{Reject (one-sided upper-tail test)}$$

$$B^* \leq -z_\alpha \Rightarrow \text{Reject (one-sided lower-tail test)}$$

$$|B^*| \geq z_{\alpha/2} \Rightarrow \text{Reject (Two-sided test)}$$

$$z_\alpha : \text{upper } 100\alpha \text{ percent point of } N(0, 1)$$

Testing equality to the specified value

$$H_0 : \theta = \theta_0$$

Test statistic

$$\psi_i = \begin{cases} 1, & (Z_i > \theta_0) \\ 0, & (Z_i < \theta_0) \end{cases}$$

$$\Rightarrow B = \sum_{i=1}^n \psi_i$$

2.2.2 An estimator associated with the sign statistic (Hodges-Lehmann)

Derivation Choose the value of θ_0 such that $H_0 : \theta = \theta_0$ is least significant, that is

$$\hat{\theta} = \underset{\theta_0}{\text{argmin}} \left| B_0 - \frac{n}{2} \right|,$$

B_0 : test statistic for testing “ $\theta = \theta_0$ ”

Estimator

$$\hat{\theta} = \text{median}\{Z_i : 1 \leq i \leq n\}$$

2.2.3 A distribution-free confidential interval based on the sign test (Thompson, Savur)

Derivation Collect θ_0 such that the two-sided test for

$$H_0 : \theta = \theta_0 \text{ versus } H_3 : \theta \neq \theta_0$$

does not reject.

Confidence interval

$$C_\alpha = n + 1 - b_{\alpha/2}$$

$Z_{(1)} < Z_{(2)} < \dots < Z_{(n)}$: ordered statistics of Z_1, \dots, Z_n

$$\theta_L = Z_{(C_\alpha)}, \quad \theta_U = Z_{(b_{\alpha/2})}$$

\Rightarrow

$$\Pr_\theta\{\theta_L < \theta < \theta_U\} = 1 - \alpha \text{ for all } \theta$$

Large-sample approximation

$$C_\alpha \approx \frac{n}{2} - z_{\alpha/2} \left\{ \frac{n}{4} \right\}^{1/2}$$

2.3 One-sample data

Data

$$X_1, \dots, X_n$$

Assumptions

C1. X_1, \dots, X_n : independent

C2. Let $F_i(t) = \Pr\{X_i \leq t\}$ ($i = 1, \dots, n$).

F_i is continuous and

$$F_i(\theta + t) + F_i(\theta - t) = 1, \text{ for every } t \quad (n = 1, \dots, n)$$

2.3.1 A distribution-free signed rank test

Null hypothesis

$$H_0 : \theta = \theta_0$$

Define

$$Z_i = X_i - \theta_0 \quad (i = 1, \dots, n)$$

The same procedure as §2.1.1 can be used.

2.3.2 An estimator based on the signed rank test

The same procedure as §2.1.2 can be used

2.3.3 A distribution free confidential interval based on the signed rank test

The same procedure as §2.1.3 can be used.

2.4 Procedures based on the sign statistic

Data

$$X_1, \dots, X_n$$

Assumptions

D1.

$$X_1, \dots, X_n : \text{independent}$$

D2.

$$\Pr\{X_i < \theta\} = \Pr\{X_i > \theta\} = \frac{1}{2} \quad (n = 1, \dots, n)$$

Define

$$Z_i = X_i - \theta_0 \quad (i = 1, \dots, n)$$

Procedures The same procedures in §2.2 can be used for testing problem, point estimation and interval estimation.

3 The Two-Sample Location Problem

Data

$$X_1, \dots, X_m, Y_1, \dots, Y_n$$

Assumptions

A1. X_1, \dots, X_m : a random sample from population 1

Y_1, \dots, Y_n : a random sample from population 2

that is

X_1, \dots, X_m : i.i.d. (independent and identically distributed)

Y_1, \dots, Y_n : i.i.d. (independent and identically distributed)

A2. $(X_1, \dots, X_m), (Y_1, \dots, Y_n)$: independent

A3. Population distribution functions of 1 ($F(t)$) and 2 ($G(t)$) are continuous.

Null hypothesis

$$H_0 : F(t) = G(t) \text{ for all } t$$
$$F(t) = \Pr\{X_i \leq t\}, \quad G(t) = \Pr\{Y_i \leq t\}$$

Alternative hypothesis(Assumption)

$$H : G(t) = F(t - \Delta)$$

Under H

$$\Pr\{X_i + \Delta \leq t\} = \Pr\{X_i \leq t - \Delta\} = F(t - \Delta) = G(t)$$
$$\Rightarrow X_1 + \Delta, \dots, X_m + \Delta, Y_1, \dots, Y_n : i.i.d.$$

F : control (对照群)
 G : case (処置群)
 Δ : treatment effect

3.1 A distribution-free rank sum test (Wilcoxon, Mann and Whitney)

Testing problems

$$H_0 : \Delta = 0 \text{ versus } \begin{cases} H_1 : \Delta > 0 \\ H_2 : \Delta < 0 \\ H_3 : \Delta \neq 0 \end{cases}$$

Derivation

Under H_0

$$X_1, \dots, X_m, Y_1, \dots, Y_n : i.i.d.$$

Test statistic

$$W = \sum_{j=1}^n S_j$$
$$S_j : \text{rank of } Y_j \text{ in } \{X_1, \dots, X_m, Y_1, \dots, Y_n\}$$

Null distribution

$\mathcal{S}_{n,N}$: set of subsets of $\{1, 2, \dots, N\}$ ($N = m + n$) of size n

\Rightarrow

$\{S_1, \dots, S_n\}$ is uniformly distributed random set on $\mathcal{S}_{n,N}$

Example $m = 4, n = 3$

$$\mathcal{S}_{4,7} = \left\{ \begin{array}{l} \{1, 2, 3\}, \{1, 2, 4\}, \{1, 2, 5\}, \{1, 2, 6\}, \{1, 2, 7\}, \{1, 3, 4\}, \{1, 3, 5\}, \\ \{1, 3, 6\}, \{1, 3, 7\}, \{1, 4, 5\}, \{1, 4, 6\}, \{1, 4, 7\}, \{1, 5, 6\}, \{1, 5, 7\}, \\ \{1, 6, 7\}, \{2, 3, 4\}, \{2, 3, 5\}, \{2, 3, 6\}, \{2, 3, 7\}, \{2, 4, 5\}, \{2, 4, 6\}, \\ \{2, 4, 7\}, \{2, 5, 6\}, \{2, 5, 7\}, \{2, 6, 7\}, \{3, 4, 5\}, \{3, 4, 6\}, \{3, 4, 7\}, \\ \{3, 5, 6\}, \{3, 5, 7\}, \{3, 6, 7\}, \{4, 5, 6\}, \{4, 5, 7\}, \{4, 6, 7\}, \{5, 6, 7\} \end{array} \right\}$$

w	6	7	8	9	10	11	12	13	14	15	16	17	18
$\Pr\{W = w\}$	$\frac{1}{35}$	$\frac{1}{35}$	$\frac{2}{35}$	$\frac{3}{35}$	$\frac{4}{35}$	$\frac{4}{35}$	$\frac{5}{35}$	$\frac{4}{35}$	$\frac{4}{35}$	$\frac{3}{35}$	$\frac{2}{35}$	$\frac{1}{35}$	$\frac{1}{35}$

Three types of test

a. One-sided upper-tail test

$$H_0 : \Delta = 0 \text{ versus } H_1 : \Delta > 0$$

$$W \geq w_\alpha \Rightarrow \text{Reject } H_0,$$

b. One-sided lower-tail test

$$H_0 : \Delta = 0 \text{ versus } H_2 : \Delta < 0$$

$$W \leq n(m + n + 1) - w_\alpha \Rightarrow \text{Reject } H_0,$$

c. Two-sided test

$$H_0 : \Delta = 0 \text{ versus } H_3 : \Delta \neq 0$$

$$W \geq w_{\alpha/2} \text{ or } W \leq n(m + n + 1) - w_{\alpha/2} \Rightarrow \text{Reject } H_0,$$

Large-sample approximation

Under H_0

$$E[W] = \frac{n(m + n + 1)}{2}, \quad \text{Var}_0[W] = \frac{mn(m + n + 1)}{12}$$

$$W^* = \frac{W - E[W]}{\sqrt{\text{Var}_0[W]}} \xrightarrow{d} N(0, 1) \quad (n \rightarrow \infty)$$

$$W^* \geq z_\alpha \Rightarrow \text{Reject (one-sided upper-tail test)}$$

$$W^* \leq -z_\alpha \Rightarrow \text{Reject (one-sided lower-tail test)}$$

$$|W^*| \geq z_{\alpha/2} \Rightarrow \text{Reject (Two-sided test)}$$

$$z_\alpha : \text{upper } 100\alpha \text{ percent point of } N(0, 1)$$

Proof

$$\phi(x, y) = \begin{cases} 1 & \text{if } x < y \\ 0 & \text{otherwise} \end{cases}$$

$$U = \sum_{i=1}^m \sum_{j=1}^n \phi(X_i, Y_j)$$

$$\Rightarrow W = U + \frac{n(n + 1)}{2}$$

∴

$R(Y_j)$: rank of Y_j

$$\Rightarrow R(Y_j) = \sum_{i=1}^m \phi(X_i, Y_j) + \sum_{j'=1}^n \phi(Y_{j'}, Y_j) + 1,$$

$$W = \sum_{j=1}^n \sum_{i=1}^m \phi(X_i, Y_j) + \sum_{j=1}^n \sum_{j'=1}^n \phi(Y_{j'}, Y_j) + n$$

$$\sum_{j=1}^n \sum_{j'=1}^n \phi(Y_{j'}, Y_j) = 0 + 1 + \dots + n - 1 = \frac{n(n-1)}{2}$$

$$\mathbb{E}[U] = \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}[\phi(X, Y)], \quad X, Y \stackrel{i.i.d.}{\sim} F,$$

$$\mathbb{E}[\phi(X, Y)] = \mathbb{E}_Y[\phi(X, Y)|Y] = \mathbb{E}_Y[F(Y)] = \frac{1}{2}$$

∵ $F(Y) \sim U(0, 1)$ Uniform distribution

$$\mathbb{E}[W] = \frac{mn}{2} + \frac{n(n+1)}{2} = \frac{n(m+n+1)}{2}$$

$$\text{Var}[U] = \sum_{i=1}^m \sum_{i'=1}^m \sum_{j=1}^n \sum_{j'=1}^n \text{Cov}[\phi(X_i, Y_j), \phi(X_{i'}, Y_{j'})]$$

If $i \neq i', j \neq j'$

$$\text{Cov}[\phi(X_i, Y_j), \phi(X_{i'}, Y_{j'})] = 0$$

If $i \neq i', j = j'$

$$\begin{aligned} \text{Cov}[\phi(X_i, Y_j), \phi(X_{i'}, Y_j)] &= \mathbb{E}[\phi(X_i, Y_j)\phi(X_{i'}, Y_j)] - \mathbb{E}[\phi(X_i, Y_j)] \cdot \mathbb{E}[\phi(X_{i'}, Y_j)] \\ &= \mathbb{E}_{X_i, X_{i'}} \mathbb{E}[\phi(X_i, Y_j)\phi(X_{i'}, Y_j)|X_i, X_{i'}] - \frac{1}{4} \\ &= \mathbb{E}_{X_i, X_{i'}} [\text{Pr}\{\max\{X_i, X_{i'}\} < Y | X_i, X_{i'}\}] - \frac{1}{4} \\ &= \mathbb{E}_{X_i, X_{i'}} [\min\{1 - F(X_i), 1 - F(X_{i'})\}] - \frac{1}{4} = \frac{1}{3} - \frac{1}{4} = \frac{1}{12} \end{aligned}$$

If $i = i', j \neq j'$

$$\text{Cov}[\phi(X_i, Y_j), \phi(X_i, Y_{j'})] = \frac{1}{12} \quad \text{similarly as above}$$

If $i = i', j = j'$

$$\text{Cov}[\phi(X_i, Y_j), \phi(X_i, Y_j)] = \mathbb{E}[\phi(X_i, Y_j)^2] - \frac{1}{4} = \mathbb{E}[\phi(X_i, Y_j)] - \frac{1}{4} = \frac{1}{4}$$

$$\begin{aligned} \text{Var}[U] &= \sum_{i=1}^m \sum_{i'=1}^m \sum_{j=1}^n \sum_{j'=1}^n \left\{ \frac{(1 - \delta_{ii'})\delta_{jj'}}{12} + \frac{\delta_{ii'}(1 - \delta_{jj'})}{12} + \frac{\delta_{ii'}\delta_{jj'}}{4} \right\} \\ &= \frac{1}{12} \sum_{i=1}^m \sum_{i'=1}^m \sum_{j=1}^n \sum_{j'=1}^n \{\delta_{jj'} + \delta_{ii'} + \delta_{ii'}\delta_{jj'}\} = \frac{1}{12}(m^2n + mn^2 + mn) = \frac{mn(m+n+1)}{12} \end{aligned}$$

Theorem (CLT for a general Mann–Whittney statistic)

Assume

$$\begin{aligned} X_1, \dots, X_m &: i.i.d., \quad Y_1, \dots, Y_n : i.i.d. \\ (X_1, \dots, X_m), (Y_1, \dots, Y_n) &: \text{independent} \\ \mathbb{E}[\varphi(X_i, Y_j)] &= 0, \quad \text{Var}[\varphi(X_i, Y_j)] = \sigma^2 < \infty. \end{aligned}$$

Let

$$U = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \varphi(X_i, Y_j).$$

If

$$\lim_{N \rightarrow \infty} \frac{N}{m} = \exists \rho_{10}, \quad \lim_{N \rightarrow \infty} \frac{N}{n} = \exists \rho_{01}$$

then

$$\sqrt{N}U \xrightarrow{d} N(0, \sigma_*^2), \quad \sigma_*^2 = \rho_{10}\sigma_{10}^2 + \rho_{01}\sigma_{01}^2$$

where $N = m + n$, and

$$\sigma_{10}^2 = \text{Var}_X[\mathbb{E}[\varphi(X, Y)|X]], \quad \sigma_{01}^2 = \text{Var}_Y[\mathbb{E}[\varphi(X, Y)|Y]]$$

Proof

Let

$$\begin{aligned} U^* &= \frac{1}{m} \sum_{i=1}^m \varphi_{10}(X_i) + \frac{1}{n} \sum_{j=1}^n \varphi_{01}(Y_j), \\ \varphi_{10}(x) &= \mathbb{E}[\varphi(x, Y)], \quad \varphi_{01}(y) = \mathbb{E}[\varphi(X, y)]. \end{aligned}$$

Then

$$\begin{aligned} \text{Var}[U^*] &= \frac{1}{m}\sigma_{10}^2 + \frac{1}{n}\sigma_{01}^2, \\ \sqrt{N}U^* &\xrightarrow{d} N(0, \sigma_*^2). \end{aligned}$$

$$\begin{aligned} \text{Var}[U] &= \frac{1}{m}\sigma_{10}^2 + \frac{1}{n}\sigma_{01}^2 + \frac{1}{mn}(\sigma^2 - \sigma_{10}^2 - \sigma_{01}^2), \\ \text{Cov}[U, U^*] &= \frac{1}{m}\sigma_{10}^2 + \frac{1}{n}\sigma_{01}^2. \end{aligned}$$

Hence

$$\begin{aligned} \mathbb{E}[N(U - U^*)^2] &= N\{\text{Var}[U] - 2\text{Cov}[U, U^*] + \text{Var}[U^*]\} \\ &= \frac{N}{mn}(\sigma^2 - \sigma_{10}^2 - \sigma_{01}^2) \rightarrow 0 \quad \Rightarrow \quad \sqrt{N}(U - U^*) \xrightarrow{P} 0. \\ \sqrt{N}U &= \sqrt{N}U^* + \{\sqrt{N}(U - U^*)\} \xrightarrow{d} N(0, \sigma_*^2) \end{aligned}$$

3.2 An estimator associated with Wilcoxon's rank sum statistic (Hodges–Lehmann)

Testing equality to the specified value

$$H_0 : \Delta = \Delta_0 \text{ (specified value)}$$

Let

$$S_j(\Delta_0) : \text{rank of } Y_j - \Delta_0 \text{ in } \{X_1, \dots, X_m, Y_1 - \Delta_0, \dots, Y_n - \Delta_0\}$$

Two-sided test

$$H_0 \text{ versus } H_3 : \Delta \neq \Delta_0$$

$$W(\Delta_0) = \sum_{j=1}^n S_j(\Delta_0) \geq w_{\alpha/2} \text{ or } W \leq n(m+n+1) - w_{\alpha/2} \Rightarrow \text{Reject } H_0$$

Derivation Choose the value of Δ_0 such that $H_0 : \Delta = \Delta_0$ is least significant, that is

$$\hat{\Delta} = \operatorname{argmin}_{\Delta_0} \left| W(\Delta_0) - \frac{n(m+n+1)}{2} \right|$$

Estimator

$$\hat{\Delta} = \operatorname{median}\{(Y_j - X_i); i = 1, \dots, m; j = 1, \dots, n\}$$

\therefore

$$W(\Delta_0) = \sum_{i=1}^m \sum_{j=1}^n \phi(X_i, Y_j - \Delta_0) + \frac{n(n+1)}{2}.$$

3.3 A distribution-free confidence interval based on Wilcoxon's rank sum test (Moses)

Derivation Collect Δ_0 such that the two-sided test for

$$H_0 : \Delta = \Delta_0 \text{ versus } H_3 : \Delta \neq \Delta_0$$

does not reject.

Confidence interval

$$C_\alpha = \frac{n(2m+n+1)}{2} + 1 - w_{\alpha/2}$$

$U_{(1)} < U_{(2)} < \dots < U_{mn}$: ordered statistic of

$$\{(Y_j - X_i); i = 1, \dots, m; j = 1, \dots, n\}$$

$$\Delta_L = U_{(C_\alpha)}, \quad \Delta_U = U_{mn+1-(C_\alpha)}$$

\Rightarrow

$$\Pr_{\Delta} \{\Delta_L < \Delta < \Delta_U\} = 1 - \alpha \text{ for all } \Delta.$$

Large-sample approximation

$$C_\alpha \approx \frac{mn}{2} - z_{\alpha/2} \left\{ \frac{mn(m+n+1)}{12} \right\}^{1/2}.$$

4 The Two-Sample Dispersion Problem and Other Two-Sample Problems

Data

$$X_1, \dots, X_m, Y_1, \dots, Y_n$$

Assumptions

- A1. X_1, \dots, X_m : a random sample from continuous population 1
 Y_1, \dots, Y_n : a random sample from continuous population 2
that is
 X_1, \dots, X_m : i.i.d. (independent and identically distributed)
 Y_1, \dots, Y_n : i.i.d. (independent and identically distributed)
- A2. $(X_1, \dots, X_m), (Y_1, \dots, Y_n)$: independent

4.1 A distribution-free rank test for dispersion—case of equal medians (Ansari–Bradley)

Null hypothesis

F, G : the distribution functions corresponding to population 1 and 2.

$$H_0 : F(t) = G(t), \text{ for every } t$$

Alternative hypothesis (Assumption)

$$H : F(t) = H\left(\frac{t - \theta_1}{\eta_1}\right) \text{ and } G(t) = H\left(\frac{t - \theta_2}{\eta_2}\right)$$

Additional assumption

A3. $\theta_1 = \theta_2$

Testing problems

$$H_0 : \gamma = 1 \text{ versus } \begin{cases} H_1 : \gamma^2 > 1 \\ H_2 : \gamma^2 < 1, \\ H_3 : \gamma^2 \neq 1 \end{cases}$$

where

$$\gamma = \frac{\eta_1}{\eta_2}$$

Test statistic

$$C = \sum_{j=1}^n S_j,$$
$$S_j = \frac{N+1}{2} - \left| \frac{N+1}{2} - R_j \right|$$
$$R_j : \text{rank of } Y_j \text{ in } \{X_1, \dots, X_m, Y_1, \dots, Y_n\},$$
$$N = m + n$$

Note

$$R_j \leq \frac{N+1}{2} \Rightarrow S_j = R_j$$
$$R_j > \frac{N+1}{2} \Rightarrow S_j = N+1 - R_j$$

Three types of test

a. One-sided upper-tail test

$$H_0 : \gamma^2 = 1 \text{ versus } H_1 : \gamma^2 > 1$$
$$C \geq c_\alpha \Rightarrow \text{Reject } H_0,$$

b. One-sided lower-tail test

$$H_0 : \gamma^2 = 1 \text{ versus } H_2 : \gamma^2 < 1$$
$$C \leq [c_{1-\alpha} - 1] \Rightarrow \text{Reject } H_0,$$

c. Two-sided test

$$H_0 : \gamma^2 = 1 \text{ versus } H_3 : \gamma^2 \neq 1$$
$$C \geq c_{\alpha/2} \text{ or } C \leq [c_{1-\alpha/2} - 1] \Rightarrow \text{Reject } H_0,$$

Null distribution

(S_1, \dots, S_n) : random sample of n integers from

$$\begin{cases} \{1, 2, \dots, \frac{N}{2}, \frac{N}{2}, \frac{N}{2} - 1, \dots, 1\} & (N : \text{even}) \\ \{1, 2, \dots, \frac{N+1}{2}, \frac{N+1}{2} - 1, \dots, 1\} & (N : \text{odd}) \end{cases}$$

without replacement

Lemma

(X_1, \dots, X_n) : random sample of size n from finite population

$$\Pi = (a_1, a_2, \dots, a_N)$$

without replacement.

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n X_j$$

\Rightarrow

$$\begin{aligned} \mathbb{E}[\bar{X}] &= \mu = \frac{1}{N} \sum_{j=1}^N a_j, \\ \text{Var}[\bar{X}] &= \frac{N-n}{n(N-1)} \sigma^2, \quad \sigma^2 = \frac{1}{N} \sum_{j=1}^N (a_j - \mu)^2 \end{aligned}$$

Proof.

$$\Pr\{X_j = a_k\} = \frac{(N-1)(N-2)\cdots(N-n+1)}{N(N-1)\cdots(N-n+1)} = \frac{1}{N}$$

$$\Rightarrow \mathbb{E}[X_j] = \mu, \quad \text{Var}[X_j] = \mathbb{E}[X_j^2] - \mu^2 = \frac{1}{N} \sum_{j=1}^N a_j^2 - \mu^2 = \sigma^2$$

$$\Pr\{X_i = a_k, X_j = a_m\} = \frac{(N-2)\cdots(N-n+1)}{N(N-1)\cdots(N-n+1)} = \frac{1}{N(N-1)}$$

$$\frac{1}{N} \left(\sum_{k=1}^N a_k \right)^2 = \frac{1}{N} \sum_{k=1}^N a_k^2 + \frac{1}{N} \sum_{k \neq m} a_k a_m$$

$(i \neq j) \Rightarrow$

$$\mathbb{E}[X_i X_j] = \frac{1}{N(N-1)} \sum_{k \neq m} a_k a_m = \frac{1}{N-1} \left\{ N\mu^2 - (\sigma^2 + \mu^2) \right\} = \mu^2 - \frac{1}{N-1} \sigma^2$$

$$\mathbb{E}[\bar{X}] = \frac{1}{n} \sum_{j=1}^n \mathbb{E}[X_j] = \mu$$

$$\text{Var}[\bar{X}] = \frac{1}{n^2} \sum_{j=1}^n \text{Var}[X_j] + \frac{1}{n^2} \sum_{i \neq j} \text{Cov}[X_i, X_j] = \frac{1}{n} \sigma^2 - \frac{n(n-1)}{n^2(N-1)} \sigma^2$$

Case of N : even

$$\mathbb{E}_0[C] = \frac{2}{N} \sum_{k=1}^{N/2} k = \frac{n(N+2)}{4}$$

$$\text{Var}_0[C] = \frac{n(N-n)}{N-1} \left\{ \frac{2}{N} \sum_{k=1}^{N/2} k^2 - \left(\frac{N+2}{4} \right)^2 \right\} = \frac{mn(N-2)(N+2)}{48(N-1)}$$

Case of N : odd

$$\mathbb{E}_0[C] = \frac{n(N+1)^2}{4N}$$

$$\text{Var}_0[C] = \frac{mn(N+1)(3+N^2)}{48N^2}$$

Large-sample approximation

$$C^* = \frac{C - E_0[C]}{\sqrt{\text{Var}_0[C]}} \rightarrow N(0, 1) \quad (m, n \rightarrow \infty)$$

\therefore Under H_0

$$F^{-1}(Y_1), \dots, F^{-1}(Y_n), F^{-1}(X_1), \dots, F^{-1}(X_m) \stackrel{i.i.d.}{\sim} \mathcal{U}(0, 1)$$

(uniform distribution on $(0, 1)$)

The following theorem can be applied since

$$C = \frac{n(N+1)}{2} - \frac{N+2}{2} \sum_{i=1}^n \left| 1 - \frac{2R_i}{N+1} \right|.$$

Theorem (CLT for a general rank statistic) Let

$$U_1, \dots, U_n, U_{n+1}, \dots, U_{m+n} \stackrel{i.i.d.}{\sim} \mathcal{U}(0, 1)$$

R_i : rank of U_i

$$S_N = \sum_{i=1}^N c_i \phi\left(\frac{R_i}{N+1}\right),$$

$$N = m + n,$$

ϕ : Lipschitz continuous function on $[0, 1]$ which is not constant

If

$$\lim_{N \rightarrow \infty} \frac{\max_i (c_i^{(N)} - \bar{c}^{(N)})^2}{\sum_{i=1}^N (c_i^{(N)} - \bar{c}^{(N)})^2} = 0 \quad \left(\bar{c}^{(N)} = \frac{1}{N} \sum_{i=1}^N c_i^{(N)} \right),$$

Then

$$\frac{S_N - E[S_N]}{\sqrt{\text{Var}[S_N]}} \rightarrow N(0, 1) \quad (N \rightarrow \infty)$$

Proof

Define

$$T_N = \sum_{i=1}^N d_i \phi(U_i), \quad d_i = c_i^{(N)} - \bar{c}^{(N)}$$

Then

$$(1) \quad \frac{T_N - E[T_N]}{\sqrt{\text{Var}[T_N]}} \rightarrow N(0, 1)$$

$$(2) \quad \frac{\text{Var}[S_N - T_N]}{\text{Var}[T_N]} \rightarrow 0$$

$$(3) \quad \frac{S_N - \mathbb{E}[S_N]}{\sqrt{\text{Var}[S_N]}} - \frac{T_N - \mathbb{E}[T_N]}{\sqrt{\text{Var}[T_N]}} \xrightarrow{P} 0$$

when $(N \rightarrow \infty)$.

Proof of (1)

$$\begin{aligned} g(t) &:= \log \mathbb{E}[\exp(it\phi(U_i))] = i\mu t - \frac{t^2}{2}\sigma^2 + \varepsilon(t)t^2 \\ \varepsilon(t) &= \frac{1}{2}\{g''(\theta t) - g''(0)\} \quad (0 < \theta < 1) \\ \mu &= \mathbb{E}[\phi(U_i)] = \frac{1}{i}g'(0), \quad \sigma^2 = \text{Var}[\phi(U_i)] = -g''(0) \\ \left| \log \mathbb{E} \left[\exp \left(\frac{T_N - \mathbb{E}[T_N]}{\sqrt{\text{Var}[T_N]}} \right) \right] + \frac{t^2}{2} \right| &= \left| \sum_{j=1}^N \varepsilon \left(\frac{d_j t}{\sqrt{\sum_{j=1}^N d_j^2}} \right) \frac{d_j^2 t^2}{\sum_{j=1}^N d_j^2} \right| \\ &\leq \max_j \left| \varepsilon \left(\frac{d_j t}{\sqrt{\sum_{j=1}^N d_j^2}} \right) \right| \rightarrow 0 \end{aligned}$$

Proof of (3)

$\mathcal{L} = \{X : \text{random variable such that } \mathbb{E}[X] = 0, \text{Var}[X] < \infty\}$
 $\implies \text{Cov}[X, Y] : \text{an inner product on the linear space } \mathcal{L}$

$$\begin{aligned} &\sqrt{\mathbb{E} \left[\left(\frac{S_N - \mathbb{E}[S_N]}{\sqrt{\text{Var}[S_N]}} - \frac{T_N - \mathbb{E}[T_N]}{\sqrt{\text{Var}[T_N]}} \right)^2 \right]} \\ &= \left\{ \mathbb{E} \left[\left(\frac{S_N - \mathbb{E}[S_N]}{\sqrt{\text{Var}[S_N]}} - \frac{S_N - \mathbb{E}[S_N]}{\sqrt{\text{Var}[T_N]}} + \frac{S_N - \mathbb{E}[S_N]}{\sqrt{\text{Var}[T_N]}} - \frac{T_N - \mathbb{E}[T_N]}{\sqrt{\text{Var}[T_N]}} \right)^2 \right] \right\}^{1/2} \\ &\leq \frac{|\sqrt{\text{Var}[T_N]} - \sqrt{\text{Var}[S_N]}|}{\sqrt{\text{Var}[T_N]}} + \sqrt{\frac{\text{Var}[S_N - T_N]}{\text{Var}[T_N]}} \leq 2\sqrt{\frac{\text{Var}[S_N - T_N]}{\text{Var}[T_N]}} \end{aligned}$$

Proof of (2)

$$(2-1) \quad U_1|_{R_1=i, R_2=j} \sim U_1|_{R_1=i} \quad (i \neq j)$$

$$(2-2) \quad \text{Cov} \left[\phi \left(\frac{R_1}{N+1} \right), \phi(U_2) \right] = -\frac{1}{N-1} \text{Cov} \left[\phi \left(\frac{R_1}{N+1} \right), \phi(U_1) \right]$$

$$(2-3) \quad \text{Cov} \left[\phi \left(\frac{R_1}{N+1} \right), \phi \left(\frac{R_2}{N+1} \right) \right] = -\frac{1}{N-1} \text{Var} \left[\left(\frac{R_1}{N+1} \right) \right]$$

$$(2-4) \quad \frac{\text{Var}[S_N - T_N]}{\text{Var}[T_N]} = \left\{ \frac{N}{N-1} \frac{\text{Var} \left[\phi \left(\frac{R_1}{N+1} \right) - \phi(U_1) \right]}{\text{Var}[\phi(U_1)]} - \frac{1}{N-1} \right\}$$

$$(2-5) \quad \text{Var} \left[\phi \left(\frac{R_1}{N+1} \right) - \phi(U_1) \right] \rightarrow 0$$

Proof of (2-1)

$$\text{Formula : } \int_a^b (u-a)^j (b-u)^k = (b-a)^{j+k+1} \frac{j!k!}{(j+k+1)!}$$

If $i < j$

$$\Pr\{R_1 = i, R_2 = j\} = \frac{1}{N(N-1)}$$

$$\Pr\{U_1 \leq t \mid R_1 = i, R_2 = j\}$$

$$= \frac{1}{\Pr(R_1 = i, R_2 = j)} \frac{(N-2)!}{(i-1)!(j-i-1)!(N-j)!} \cdot \Pr\{U_3 < U_1, \dots, U_{i+1} < U_1, U_1 < U_{i+2}, \dots, U_{j-1} < U_2, U_2 < U_{j+1}, \dots, U_n, U_1 \leq t\}$$

$$= \frac{N!}{(i-1)!(j-i-1)!(N-j)!} \int_0^t \left\{ \int_{u_1}^1 u_1^{i-1} (u_2 - u_1)^{j-i-1} (1-u_2)^{N-j} du_2 \right\} du_1$$

$$= \frac{N!}{(i-1)!(N-i)!} \int_0^t u_1^{i-1} (1-u_1)^{N-i} du_1,$$

$$\Pr\{U_1 \leq t \mid R_1 = i\}$$

$$= N \cdot \frac{(N-1)!}{(i-1)!(N-i)!} \Pr\{U_2 < U_1, \dots, U_i < U_1, U_1 < U_{i+1}, \dots, U_1 < U_N, U_1 \leq t\}$$

$$= \frac{N!}{(i-1)!(N-i)!} \int_0^t u_1^{i-1} (1-u_1)^{N-i} du_1,$$

Similarly,

$$\Pr\{U_2 \leq t \mid R_1 = i, R_2 = j\}$$

$$= \frac{1}{\Pr\{R_1 = i, R_2 = j\}} \frac{(N-2)!}{(i-1)!(j-i-1)!(N-j)!} \cdot \Pr\{U_3 < U_1, \dots, U_{i+1} < U_1, U_1 < U_{i+2}, \dots, U_{j-1} < U_2, U_2 < U_j, \dots, U_n, U_2 \leq t\}$$

$$= \frac{N!}{(i-1)!(j-i-1)!(N-j)!} \mathbb{E}[1_{U_2 \leq t, U_1 < U_2} U_1^{i-1} (U_2 - U_1)^{j-i-1} (1-U_2)^{N-j}]$$

$$= \frac{N!}{(i-1)!(j-i-1)!(N-j)!} \int_0^t \left\{ \int_0^{u_2} u_1^{i-1} (u_2 - u_1)^{j-i-1} (1-u_2)^{N-j} du_1 \right\} du_2$$

$$= \frac{N!}{(j-1)!(N-j)!} \int_0^t u_2^{j-1} (1-u_2)^{N-j} du_2 = \Pr\{U_2 \leq t \mid R_2 = j\}$$

Proof of (2-2)

$$\begin{aligned}
\psi(j) &= \phi\left(\frac{j}{N+1}\right) - \frac{1}{N} \sum_{i=1}^N \phi\left(\frac{i}{N+1}\right) \\
\text{Cov}\left[\phi\left(\frac{R_1}{N+1}\right), \phi(U_2)\right] &= \text{E}[\psi(R_1)\phi(U_2)] \\
&= \frac{1}{N(N-1)} \sum_{i \neq j} \psi(i) \text{E}[\phi(U_2) | R_1 = i, R_2 = j] \\
&= \frac{1}{N(N-1)} \sum_{i \neq j} \psi(i) \text{E}[\phi(U_2) | R_2 = j] \\
&= \frac{1}{N(N-1)} \left\{ \sum_{i=1}^N \psi(i) \right\} \left\{ \sum_{j=1}^N \text{E}[\phi(U_2) | R_2 = j] \right\} \\
&\quad - \frac{1}{N(N-1)} \sum_{i=1}^N \psi(i) \text{E}[\phi(U_2) | R_2 = i] \\
&= -\frac{1}{N-1} \text{E}[\psi(R_2)\phi(U_2)] = -\frac{1}{N-1} \text{Cov}[\psi(R_2), \phi(U_2)]
\end{aligned}$$

Proof of (2-3)

$$\begin{aligned}
\text{Cov}\left[\phi\left(\frac{R_1}{N+1}\right), \phi\left(\frac{R_2}{N+1}\right)\right] &= \text{E}[\psi(R_1), \psi(R_2)] \\
&= \frac{1}{N(N-1)} \sum_{i \neq j} \psi(i)\psi(j) = \frac{1}{N(N-1)} \left\{ \sum_{i=1}^N \psi(i) \right\}^2 - \frac{1}{N(N-1)} \sum_{i=1}^N \{\psi(i)\}^2 \\
&= -\frac{1}{N-1} \text{Var}\left[\phi\left(\frac{R_1}{N+1}\right)\right]
\end{aligned}$$

Proof of (2-4)

$$\begin{aligned}
\text{Cov}[\psi(R_1) - \phi(U_1), \psi(R_2) - \phi(U_2)] &= \text{Cov}[\psi(R_1), \psi(R_2)] - 2\text{Cov}[\psi(R_1), \phi(U_2)] \\
&= -\frac{1}{N-1} \left\{ \text{Var}[\psi(R_1)] - 2\text{Cov}[\psi(R_1), \phi(U_1)] \right\} \\
0 &= \left\{ \sum_{i=1}^N d_i \right\}^2 = \sum_{i=1}^N d_i^2 + \sum_{i \neq j} d_i d_j \\
\text{Var}[S_N - T_N] &= \sum_{i=1}^N d_i^2 \text{Var}[\psi(R_i) - \phi(U_i)] + \sum_{i \neq j} d_i d_j \text{Cov}[\psi(R_i) - \phi(U_i), \psi(R_j) - \phi(U_j)] \\
&= \left(\sum_{i=1}^N d_i^2 \right) \left\{ \frac{N}{N-1} \left\{ \text{Var}[\psi(R_1)] - 2\text{Cov}[\psi(R_1), \phi(U_1)] \right\} + \text{Var}[\phi(U_1)] \right\} \\
\text{Var}[T_N] &= \left(\sum_{i=1}^N d_i^2 \right) \text{Var}[\phi(U_1)]
\end{aligned}$$

Proof of (2-5)

$$\begin{aligned}
f_{U_1|R_1=j}(u) &= \frac{N!}{(j-1)!(N-j)!} u^{j-1}(1-u)^{N-j} \\
\text{Var}[U_1|R_1=j] &= \text{E}[U_1^2|R_1=j] - \text{E}[U_1|R_1=j]^2 \\
&= \frac{N!}{(j-1)!(N-j)!} \frac{(j+1)!(N-j)!}{(N+2)!} - \left(\frac{N!}{(j-1)!(N-j)!} \frac{j!(N-j)!}{(N+1)!} \right)^2 \\
&= \frac{j(N+1-j)}{(N+1)^2(N+2)} \quad (\text{Formula in page 22}) \\
\text{E}\left[\left\{\phi\left(\frac{R_1}{N+1}\right) - \phi(U_1)\right\}^2\right] &= \frac{1}{N} \sum_{j=1}^N \text{E}\left[\left\{\phi\left(\frac{j}{N+1}\right) - \phi(U_1)\right\}^2 \middle| R_1=j\right] \\
&= \frac{1}{N} \sum_{j=1}^N \left[\int_{|u-\frac{j}{N+1}|<\varepsilon} \left\{\phi\left(\frac{j}{N+1}\right) - \phi(u)\right\}^2 f_{U_1|R_1=j}(u) du \right. \\
&\quad \left. \int_{|u-\frac{j}{N+1}|>\varepsilon} \left\{\phi\left(\frac{j}{N+1}\right) - \phi(u)\right\}^2 f_{U_1|R_1=j}(u) du \right] \\
&< \frac{1}{N} \sum_{j=1}^N \left[D\varepsilon^2 \Pr\left\{\left|U_1 - \frac{j}{N+1}\right| < \varepsilon \mid R_1=j\right\} + M \Pr\left\{\left|U_1 - \frac{j}{N+1}\right| > \varepsilon \mid R_1=j\right\} \right] \\
&\quad \left(D = \max_{x<y} \left\{ \frac{\phi(x) - \phi(y)}{x-y} \right\}^2, M = \max_{x<y} \{\phi(x) - \phi(y)\}^2 \right) \\
&< \frac{1}{N} \sum_{j=1}^N \left\{ D\varepsilon^2 + \frac{M}{\varepsilon^2} \text{Var}[U_1|R_1=j] \right\} = \frac{1}{N} \sum_{j=1}^N \left\{ D\varepsilon^2 + \frac{M}{\varepsilon^2} \frac{j(N+1-j)}{(N+1)^2(N+2)} \right\} \\
&= D\varepsilon^2 + \frac{M}{\varepsilon^2} \left\{ \frac{1}{2(N+2)} - \frac{2N+1}{6(N+1)(N+2)} \right\} = D\varepsilon^2 + \frac{M}{6\varepsilon^2(N+1)} \\
\text{E}\left[\left\{\phi\left(\frac{R_1}{N+1}\right) - \phi(U_1)\right\}^2\right] &< 2\sqrt{\frac{DM}{N+1}} \quad \left(\varepsilon^2 = \sqrt{\frac{M}{D(N+1)}} \right)
\end{aligned}$$

4.2 An asymptotically distribution-free test for dispersion based on the Jackknife-medians not necessarily equal (Miller)

Assumptions

- A1. X_1, \dots, X_m : a random sample from continuous population 1
 Y_1, \dots, Y_n : a random sample from continuous population 2
that is
 X_1, \dots, X_m : i.i.d. (independent and identically distributed)
 Y_1, \dots, Y_n : i.i.d. (independent and identically distributed)
- A2. $(X_1, \dots, X_m), (Y_1, \dots, Y_n)$: independent

Further hypothesis (Assumption)

$$H : F(t) = H\left(\frac{t - \theta_1}{\eta_1}\right) \text{ and } G(t) = H\left(\frac{t - \theta_2}{\eta_2}\right)$$

Additional assumption

A4. H has finite fourth moment.

$$E\left[\left(\frac{X_1 - \theta_1}{\eta_1}\right)^4\right] < \infty$$

Procedure

$$\bar{X}_i = \sum_{s \neq i}^m \frac{X_s}{m-1}, \quad D_i^2 = \sum_{s \neq i} \frac{(X_s - \bar{X}_i)^2}{m-2} \quad (i = 1, 2, \dots, m)$$

$$\bar{Y}_j = \sum_{t \neq j}^n \frac{Y_t}{n-1}, \quad E_j^2 = \sum_{t \neq j} \frac{(Y_t - \bar{Y}_j)^2}{n-2} \quad (j = 1, 2, \dots, n)$$

$$S_i = \log D_i^2 \quad (i = 1, \dots, m)$$

$$T_j = \log E_j^2 \quad (j = 1, \dots, n)$$

$$\bar{X}_0 = \frac{1}{m} \sum_{s=1}^m X_s, \quad S_0 = \log \left[\sum_{s=1}^m \frac{(X_s - \bar{X}_0)^2}{m-1} \right],$$

$$\bar{Y}_0 = \frac{1}{n} \sum_{t=1}^n Y_t, \quad T_0 = \log \left[\sum_{t=1}^n \frac{(Y_t - \bar{Y}_0)^2}{n-1} \right]$$

$$A_i = mS_0 - (m-1)S_i \quad (i = 1, \dots, m)$$

$$B_j = nT_0 - (n-1)T_j \quad (j = 1, \dots, n)$$

$$\bar{A} = \sum_{i=1}^m \frac{A_i}{m}, \quad \bar{B} = \sum_{j=1}^n \frac{B_j}{n}$$

$$V_1 = \sum_{i=1}^m \frac{(A_i - \bar{A})^2}{m(m-1)}, \quad V_2 = \sum_{j=1}^n \frac{(B_j - \bar{B})^2}{n(n-1)}$$

Test statistic

$$Q = \frac{\bar{A} - \bar{B}}{\sqrt{V_1 + V_2}}$$

a. One-sided upper-tail test

$$H_0 : \gamma^2 = 1 \text{ versus } H_1 : \gamma^2 > 1$$

$$Q \geq z_\alpha \Rightarrow \text{Reject } H_0,$$

b. One-sided lower-tail test

$$H_0 : \gamma^2 = 1 \text{ versus } H_2 : \gamma^2 < 1$$

$$Q \leq -z_\alpha \Rightarrow \text{Reject } H_0,$$

c. Two-sided test

$$H_0 : \gamma^2 = 1 \text{ versus } H_3 : \gamma^2 \neq 1$$

$$|Q| \geq z_{\alpha/2} \Rightarrow \text{Reject } H_0,$$

Jackknife technique

$$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} F, \quad E[X_1] = \mu, \text{Var}[X_1] = \sigma^2$$

$$\hat{\theta}_n : \text{an estimator based on } X_1, \dots, X_n$$

$$\hat{\theta}_{n-1}^{(i)} : \text{the estimator based on } X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n$$

$$\tilde{\theta}_{n,i} = n\hat{\theta}_n - (n-1)\hat{\theta}_{n-1}^{(i)} : \text{jackknife pseudo-value}$$

Tukey's conjecture

(A1) $\tilde{\theta}_{n,1}, \dots, \tilde{\theta}_{n,n}$ can be viewed as an *i.i.d.* sample

(A2) $\text{Var}[\tilde{\theta}_{n,i}] \approx \text{Var}[\sqrt{n}\hat{\theta}_n]$

Jackknife estimator

$$\tilde{\theta}_n = \frac{1}{n} \sum_{i=1}^n \tilde{\theta}_{n,i} : \text{Bias modified Jackknife estimator}$$

$$V = \frac{n-1}{n} \sum_{i=1}^n \left(\hat{\theta}_{n-1}^{(i)} - \frac{1}{n} \sum_{j=1}^n \hat{\theta}_{n-1}^{(j)} \right)^2$$

: Jackknife variance estimator for $\text{Var}[\hat{\theta}_n]$

Case of a smooth function of the sample mean

$$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} F, \quad E[X_1] = \mu, \text{Var}[X_1] = \sigma^2$$

$$\hat{\theta}_n = h(\bar{X}_n), \quad h : C^{(3)}\text{-class}, \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\hat{\theta}_n^{(i)} = h(\bar{X}_{n-1}^{(i)}), \quad \bar{X}_{n-1}^{(i)} = \frac{1}{n-1} \sum_{j \neq i}^n X_j$$

Asymptotic normality of $\hat{\theta}_n$

$$Z_n = \sqrt{n}(\bar{X}_n - \mu) \longrightarrow N(0, \sigma^2) \quad \text{central limit theorem}$$

$$h(\bar{X}_n) = h\left(\mu + \frac{1}{\sqrt{n}}Z_n\right) = h(\mu) + \frac{1}{\sqrt{n}}h'(\mu)Z_n + \frac{1}{2n}h''(\mu)Z_n^2 + \dots$$

$$\sqrt{n}\{h(\bar{X}_n) - h(\mu)\} = h'(\mu)Z_n + \frac{1}{2\sqrt{n}}h''(\mu)Z_n^2 + \dots \longrightarrow N\left[0, \{h'(\mu)\}^2\sigma^2\right]$$

Bias of the jackknife estimator

$$\begin{aligned}\bar{X}_n &= \frac{n-1}{n}\bar{X}_{n-1}^{(i)} + \frac{1}{n}X_i = \bar{X}_{n-1}^{(i)} + \frac{1}{n}(X_i - \bar{X}_{n-1}^{(i)}) \\ &= \mu + \frac{1}{\sqrt{n-1}}Z_{n-1}^{(i)} + \frac{1}{n}(X_i - \mu) - \frac{1}{n\sqrt{n-1}}Z_{n-1}^{(i)}, \quad Z_{n-1}^{(i)} = \sqrt{n-1}(\bar{X}_{n-1}^{(i)} - \mu)\end{aligned}$$

$$\begin{aligned}h(\bar{X}_n) &= h(\mu) + h'(\mu)\left\{\frac{1}{\sqrt{n-1}}Z_{n-1}^{(i)} + \frac{1}{n}(X_i - \mu) - \frac{1}{n\sqrt{n-1}}Z_{n-1}^{(i)}\right\} \\ &\quad + \frac{1}{2}h''(\mu)\left\{\frac{1}{\sqrt{n-1}}Z_{n-1}^{(i)} + \frac{1}{n}(X_i - \mu) - \frac{1}{n\sqrt{n-1}}Z_{n-1}^{(i)}\right\}^2 \\ &\quad + \frac{1}{6}h^{(3)}(\mu)\left\{\frac{1}{\sqrt{n-1}}Z_{n-1}^{(i)} + \frac{1}{n}(X_i - \mu) - \frac{1}{n\sqrt{n-1}}Z_{n-1}^{(i)}\right\}^3 + O_p(n^{-2})\end{aligned}$$

$$\begin{aligned}h(\bar{X}_{n-1}^{(i)}) &= h(\mu) + h'(\mu)\frac{1}{\sqrt{n-1}}Z_{n-1}^{(i)} \\ &\quad + \frac{1}{2}h''(\mu)\frac{1}{n-1}\{Z_{n-1}^{(i)}\}^2 + \frac{1}{6}h^{(3)}(\mu)\frac{1}{(n-1)^{3/2}}\{Z_{n-1}^{(i)}\}^3 + O_p(n^{-2})\end{aligned}$$

$$\begin{aligned}\tilde{\theta}_{n,i} &= nh(\bar{X}_n) - (n-1)h(\bar{X}_{n-1}^{(i)}) \\ &= h(\mu) + \frac{h'(\mu)}{\sqrt{n-1}}Z_{n-1}^{(i)} + \frac{h''(\mu)}{2(n-1)}\{Z_{n-1}^{(i)}\}^2 + \frac{h^{(3)}(\mu)}{6(n-1)^{3/2}}\{Z_{n-1}^{(i)}\}^3 \\ &\quad + h'(\mu)\left\{(X_i - \mu) - \frac{1}{\sqrt{n-1}}Z_{n-1}^{(i)}\right\} \\ &\quad + \frac{h''(\mu)}{2}\left\{\frac{2}{\sqrt{n-1}}Z_{n-1}^{(i)}(X_i - \mu) + \frac{1}{n}(X_i - \mu)^2 - \frac{2}{n-1}\{Z_{n-1}^{(i)}\}^2\right\} \\ &\quad + \frac{h^{(3)}(\mu)}{6}\left\{\frac{3}{n-1}\{Z_{n-1}^{(i)}\}^2(X_i - \mu)\right\} + O_p(n^{-3/2}) \\ &= h(\mu) + h'(\mu)(X_i - \mu) + \frac{h''(\mu)}{\sqrt{n-1}}Z_{n-1}^{(i)}(X_i - \mu) \\ &\quad + \frac{h''(\mu)}{2}\left\{\frac{1}{n}(X_i - \mu)^2 - \frac{1}{(n-1)}\{Z_{n-1}^{(i)}\}^2\right\} \\ &\quad + \frac{h^{(3)}(\mu)}{2(n-1)}\{Z_{n-1}^{(i)}\}^2(X_i - \mu) + O_p(n^{-3/2})\end{aligned}$$

$Z_{n-1}^{(i)}, (X_i) : \text{independent}$

$$\Rightarrow \quad \text{E}[\tilde{\theta}_{n,i}] = h(\mu) + \frac{h''(\mu)}{2}\left\{\frac{1}{n} - \frac{1}{n-1}\right\}\sigma^2 + O(n^{-3/2}) = h(\mu) + O(n^{-3/2})$$

$$\begin{aligned}(i \neq j) &\Rightarrow \text{E}[\{\tilde{\theta}_{n,i} - h(\mu)\}\{\tilde{\theta}_{n,j} - h(\mu)\}] \\ &= \frac{\{h''(\mu)\}^2}{n-1}\text{E}[Z_{n-1}^{(i)}Z_{n-1}^{(j)}(X_i - \mu)(X_j - \mu)] + O(n^{-3/2}) = O(n^{-3/2})\end{aligned}$$

$$\begin{aligned}
\mathbb{E}[\{\tilde{\theta}_{n,i} - h(\mu)\}^2] &= \{h'(\mu)\}^2 \mathbb{E}[(X_i - \mu)^2] + \frac{h'(\mu)h''(\mu)}{n} \mathbb{E}[(X_i - \mu)^3] \\
&\quad + \frac{h'(\mu)h^{(3)}(\mu)}{n-1} \mathbb{E}[\{Z_{n-1}^{(i)}\}^2 (X_i - \mu)^2] + \frac{\{h''(\mu)\}^2}{n-1} \mathbb{E}[\{Z_{n-1}^{(i)}\}^2 (X_i - \mu)^2] \\
&\quad + O(n^{-3/2}) \quad (\kappa_3 = \mathbb{E}[(X_i - \mu)^3]) \\
&= \{h'(\mu)\}^2 \sigma^2 + \frac{h'(\mu)h''(\mu)}{n} \kappa_3 + \frac{1}{n-1} [h'(\mu)h^{(3)}(\mu) + \{h''(\mu)\}^2] \sigma^4 + O(n^{-3/2})
\end{aligned}$$

Let $Z_n = \sqrt{n}(\bar{X}_n - \mu)$ then

$$\begin{aligned}
\hat{\theta}_n &= h(\bar{X}_n) = h(\mu) + \frac{h'(\mu)}{\sqrt{n}} Z_n + \frac{h''(\mu)}{2n} Z_n^2 + \frac{h^{(3)}(\mu)}{6n} Z_n^3 + O_p(n^{-3/2}) \\
\text{Var}[\sqrt{n}\hat{\theta}_n] &= \{h'(\mu)\}^2 \sigma^2 + \frac{h'(\mu)h''(\mu)}{n} \kappa_3 \\
&\quad + \frac{1}{n} [h'(\mu)h^{(3)}(\mu) + \frac{1}{2}\{h''(\mu)\}^2] \sigma^4 + O(n^{-3/2})
\end{aligned}$$

So

$$\text{Var}[\sqrt{n}\hat{\theta}_n] - \text{Var}[\tilde{\theta}_n] = O(n^{-1})$$

Asymptotic normality of $\tilde{\theta}_n$

$$\begin{aligned}
\bar{X}_{n-1}^{(i)} &= \bar{X}_n - \frac{1}{n-1}(\bar{X}_n - X_i) \\
nh(\bar{X}_n) - (n-1)h(\bar{X}_{n-1}^{(i)}) &= h(\bar{X}_n) - h'(\bar{X}_n)(\bar{X}_n - X_i) - \frac{h''(\bar{X}_n)}{2(n-1)}(\bar{X}_n - X_i)^2 + O_p(n^{-2}) \\
\tilde{\theta}_n &= \frac{1}{n} \sum_{i=1}^n \tilde{\theta}_{n,i} = h(\bar{X}_n) - \frac{h''(\bar{X}_n)}{2n} S_n + O_p(n^{-2}), \quad S_n = \frac{1}{n-1} \sum_{i=1}^n (\bar{X}_n - X_i)^2 \\
\sqrt{n}\{\hat{\theta}_n - h(\mu)\} - \sqrt{n}\{\tilde{\theta}_n - h(\mu)\} &= O_p(n^{-1/2}) \quad (\because \bar{X}_n \xrightarrow{P} \mu, S_n \xrightarrow{P} \sigma^2)
\end{aligned}$$

4.3 A distribution-free rank test for either location or dispersion (Lepage)

Assumptions

A1. X_1, \dots, X_m : a random sample from continuous population 1

Y_1, \dots, Y_n : a random sample from continuous population 2

that is

X_1, \dots, X_m : i.i.d. (independent and identically distributed)

Y_1, \dots, Y_n : i.i.d. (independent and identically distributed)

A2. $(X_1, \dots, X_m), (Y_1, \dots, Y_n)$: independent

Further hypothesis (Assumption)

$$H : F(t) = H\left(\frac{t - \theta_1}{\eta_1}\right) \text{ and } G(t) = H\left(\frac{t - \theta_2}{\eta_2}\right)$$

Testing problem

$$H_0 : \theta_1 = \theta_2 \text{ and } \eta_1 = \eta_2 \text{ versus } H_1 : \theta_1 \neq \theta_2 \text{ and/or } \eta_1 \neq \eta_2$$

Test statistic

R_i : rank of Y_i in $X_1, \dots, X_m, Y_1, \dots, Y_n$

$W = \sum_{i=1}^n R_i$: Wilcoxon rank sum statistic ,

$C = \sum_{i=1}^n \left\{ \frac{N+1}{2} - \left| \frac{N+1}{2} - R_i \right| \right\}$: Ansari-Bradley scale statistic

$$D = (W^*)^2 + (C^*)^2, \quad \text{where } W^* = \frac{W - E_0[W]}{\sqrt{\text{Var}_0[W]}} \text{ and } C^* = \frac{C - E_0[C]}{\sqrt{\text{Var}_0[C]}}$$

Procedure

$$D \geq d_\alpha \Rightarrow \text{reject } H_0$$

d_α is obtained based on the fact that

(R_1, \dots, R_n) is a random sample from $\{1, 2, \dots, N\}$ without replacement

Large-sample approximation

$$D \longrightarrow \chi_2^2, \quad (\min\{m, n\} \rightarrow \infty)$$

$$D \geq \chi_{2,\alpha}^2 \Rightarrow \text{reject}$$

χ_2^2 : chi-square distribution with degree of freedom 2

$\chi_{2,\alpha}^2$: upper 100 α percent point of χ_2^2

\vdots

Let

$$T_W = \frac{m}{N} \sum_{i=1}^n U_i - \frac{n}{N} \sum_{i=1}^m U_{n+i}$$

$$T_C = \frac{m}{N} \sum_{i=1}^n \phi(U_i) - \frac{n}{N} \sum_{i=1}^m \phi(U_{n+i}) \quad \phi(u) = |1 - 2u|.$$

where $U_i = F^{-1}(X_i)$ $i = 1, \dots, m$; $U_i = F^{-1}(Y_i)$ $i = 1, \dots, n$

Then from the proof of Theorem in section 4.1,

$$W^* - \frac{T_W}{\sqrt{\text{Var}(T_W)}} \xrightarrow{P} 0, \quad C^* - \frac{T_C}{\sqrt{\text{Var}(T_C)}} \xrightarrow{P} 0$$

Let a_1, a_2 : any constants, and

$$T_a = a_1 T_W + a_2 T_C = \frac{m}{N} \sum_{i=1}^n \{a_1 U_i + a_2 \phi(U_i)\} - \frac{n}{N} \sum_{i=1}^m \{a_1 U_{i+n} + a_2 \phi(U_{i+n})\}$$

Then

$$\begin{aligned} E[T_a] &= 0 \text{ and } \text{Var}[T_a] = \frac{mn}{12N}(a_1^2 + a_2^2) \\ \frac{T_a}{\sqrt{\text{Var}[T_a]}} &= \frac{a_1 T_W^* + a_2 T_C^*}{\sqrt{a_1^2 + a_2^2}} \rightarrow N(0, 1) \quad (\min\{m, n\} \rightarrow \infty) \\ T_W^* &= \frac{T_W}{\sqrt{\text{Var}(T_W)}}, \quad T_C^* = \frac{T_C}{\sqrt{\text{Var}(T_C)}}, \end{aligned}$$

Hence

$$\begin{pmatrix} T_W^* \\ T_C^* \end{pmatrix} \rightarrow N_2 \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right]$$

4.4 A distribution-free test for general differences in two populations (Kolmogorov–Smirnov)

Assumptions

A1. $X_1, \dots, X_m \stackrel{i.i.d.}{\sim} F$, $Y_1, \dots, Y_n \stackrel{i.i.d.}{\sim} G$

A2. $(X_1, \dots, X_m), (Y_1, \dots, Y_n)$: independent

Testing problem

$H_0 : F(t) = G(t)$ for any t versus $H_1 : F(t) \neq G(t)$ for at least one t

Test statistic

$$J = \frac{mn}{d} \max_{-\infty < t < \infty} \{|F_m(t) - F_n(t)|\},$$

where

d : greatest common divisor of m and n

$$F_m(t) = \frac{1}{m} \#\{X_i | X_i \leq t, i = 1, \dots, m\}$$

(the empirical distribution function for X)

$$G_n(t) = \frac{1}{n} \#\{Y_i | Y_i \leq t, i = 1, \dots, n\}$$

(the empirical distribution function for Y)

Procedure

$$J \geq j_\alpha \Rightarrow \text{reject } H_0$$

Under H_0 the conditional distribution of (X_1, \dots, X_m) can be viewed as the one of random sample from $Z_{(1)} \leq Z_{(2)} \leq \dots \leq Z_{(N)}$ without replacement.

Deriving J

data

No.	1	2	3	4	5	6	7	8	9	10
X	3	3	5	7	9					
Y	3	4	4	6	7	8	10	10	11	12

$Z_{(1)} \leq Z_{(2)} \leq \dots \leq Z_{(N)}$: ordered statistic of $X_1, \dots, X_m, Y_1, \dots, Y_n$

i	$Z_{(i)}$	$F_5(Z_{(i)})$	$G_{10}(Z_{(i)})$	$ F_5(Z_{(5)}) - G_{10}(Z_{(i)}) $
1	3	2/5	1/10	3/10
2	3	2/5	1/10	3/10
3	3	2/5	1/10	3/10
4	4	2/5	3/10	1/10
5	4	2/5	3/10	1/10
6	5	3/5	3/10	3/10
7	6	3/5	4/10	2/10
8	7	4/5	5/10	3/10
9	7	4/5	5/10	3/10
10	8	4/5	6/10	2/10
11	9	5/5	6/10	4/10
12	10	5/5	8/10	2/10
13	10	5/5	8/10	2/10
14	11	5/5	9/10	1/10
15	12	5/5	10/10	0

$$\Rightarrow J = \frac{10 \times 5}{5} \frac{4}{10} = 4$$

Large-sample approximation

$$J^* = \sqrt{\frac{mn}{N}} \max_{-\infty < t < \infty} \{|F_m(t) - F_n(t)|\}$$

$$P_0(J^* < s) \rightarrow \sum_{k=-\infty}^{\infty} (-1)^k e^{-2k^2 s^2}, 0 \text{ for } s >, \leq 0 \quad (\min\{m, n\} \rightarrow \infty)$$

5 The One-Way Layout

Data

Treatments			
1	2	...	k
X_{11}	X_{12}	...	X_{1k}
X_{21}	X_{22}	...	X_{2k}
\vdots	\vdots		\vdots
X_{n_11}	X_{n_22}	...	X_{n_kk}

Assumptions

A1. $X_{11}, X_{21}, \dots, X_{n_11}, \dots, X_{1k}, \dots, X_{n_kk}$: independent

A2. $X_{1j}, X_{2j}, \dots, X_{n_jj} \stackrel{i.i.d.}{\sim} F_j$

A3. F_1, \dots, F_k belong to the same location family of distributions:

$$F_j(t) = F(t - \tau_j), \quad -\infty < t < \infty, \quad j = 1, \dots, k$$

τ_j : unknown parameter

F : continuous distribution function with unknown median θ

One-way layout model

$$X_{ij} = \theta + \tau_j + e_{ij} \quad i = 1, \dots, n_j, \quad j = 1, \dots, k$$

τ_j : treatment j effect

e_{ij} : i.i.d. , median[e_{ij}] = 0

Hypothesis

$$H_0 : \tau_1 = \dots = \tau_k$$

5.1 A distribution-free test for general alternatives (Kruskal-Wallis)

Testing problem

$$H_0 \quad \text{versus} \quad H_1 : \tau_1, \dots, \tau_k \text{ are not all equal}$$

Test statistic

$$H = \frac{12}{N(N+1)} \sum_{j=1}^k n_j \left(R_{.j} - \frac{N+1}{2} \right)^2 = \left(\frac{12}{N(N+1)} \sum_{j=1}^k \frac{R_j^2}{n_j} \right) - 3(N+1),$$

$$N = \sum_{j=1}^k n_j, \quad R_j = \sum_{i=1}^{n_j} r_{ij}, \quad R_{.j} = \frac{R_j}{n_j}, \quad j = 1, \dots, k$$

r_{ij} : rank of X_{ij} in all observations

Procedure

$$H \geq h_\alpha \Rightarrow \text{reject } H_0$$

h_α is derived based on the fact that

Under H_0 , all $N!/(\prod_{j=1}^k n_j!)$ assignments of n_1 ranks to the treatment 1 observations, n_2 ranks to the treatment 2 observations, \dots , n_k ranks to the treatment k observations, are equally likely.

Example for deriving the exact null distribution

$$k = 3, n_1 = n_2 = n_3 = 2, \quad \text{Number of rank assignments} = \frac{6!}{2!2!2!} = 90$$

$$H = \frac{12}{6(6+1)} \frac{R_1^2 + R_2^2 + R_3^2}{2} - 3(6+1) = \frac{A}{7} - 21, \quad A = R_1^2 + R_2^2 + R_3^2$$

	I		II		III		
(a)	1	2	3	4	5	6	$A = 179, H = 4.57$
(b)	1	2	3	5	4	6	$A = 173, H = 3.71$
(c)	1	2	3	6	4	5	$A = 171, H = 3.43$
(d)	1	3	2	4	5	6	$A = 173, H = 3.71$
(e)	1	3	2	5	4	6	$A = 165, H = 2.57$
(f)	1	3	2	6	4	5	$A = 161, H = 2.00$
(g)	1	4	2	3	5	6	$A = 171, H = 3.43$
(h)	1	4	2	5	3	6	$A = 155, H = 1.14$
(i)	1	4	2	6	3	5	$A = 153, H = 0.857$
(j)	1	5	2	3	4	6	$A = 161, H = 2.00$
(k)	1	5	2	4	3	6	$A = 153, H = 0.857$
(l)	1	5	2	6	3	4	$A = 149, H = 0.286$
(m)	1	6	2	3	4	5	$A = 155, H = 1.14$
(n)	1	6	2	4	3	5	$A = 149, H = 0.286$
(o)	1	6	2	5	3	4	$A = 147, H = 0$

$$(I, II, III) \Rightarrow (I, III, II), (II, I, III), (II, III, I), (III, I, II), (III, II, I)$$

h	0	0.286	0.857	1.14	2.00	2.57	3.43	3.71	4.57
$\Pr\{H = h\}$	$\frac{1}{15}$	$\frac{2}{15}$	$\frac{2}{15}$	$\frac{2}{15}$	$\frac{2}{15}$	$\frac{1}{15}$	$\frac{2}{15}$	$\frac{2}{15}$	$\frac{1}{15}$

Large-sample approximation

$$H \longrightarrow \chi_{k-1}^2, \quad (\min\{n_1, \dots, n_k\} \rightarrow \infty)$$

$$H \geq \chi_{k-1, \alpha}^2 \Rightarrow \text{reject } H_0$$

Proof

$$S_n := \sum_{j=1}^k \sum_{i=1}^{n_j} \frac{a_j}{n_j} r_{ij}, \quad a_1, \dots, a_k : \text{constants}$$

Then

$$\frac{S_N - E[S_N]}{\sqrt{\text{Var}[S_N]}} \rightarrow N(0, 1)$$

\therefore Theorem in section 4.1 can be applied since

$$c_{ij}^{(N)} := \frac{N+1}{n_j} a_j \quad (i = 1, \dots, n_j; j = 1, \dots, k)$$

$$\bar{c}^{(N)} = \frac{1}{N} \sum_{j=1}^k \sum_{i=1}^{n_j} c_{ij}^{(N)} = \frac{N+1}{N} \sum_{j=1}^k a_j = \frac{N+1}{N} A \quad (\text{say})$$

$$\frac{\max_{i,j} (c_{ij}^{(N)} - \bar{c}^{(N)})^2}{\sum_{j=1}^k \sum_{i=1}^{n_j} (c_{ij}^{(N)} - \bar{c}^{(N)})^2} = \frac{\max_j (\frac{a_j}{n_j} - \frac{A}{N})^2}{\sum_{j=1}^k n_j (\frac{a_j}{n_j} - \frac{A}{N})^2} < \frac{1}{n_{j_*}} \rightarrow 0, \quad j_* = \operatorname{argmax}_j (\frac{a_j}{n_j} - \frac{A}{N})^2$$

$$E[r_{ij}] = \frac{1}{N} \sum_{l=1}^N l = \frac{N+1}{2} = \mu \quad (\text{say}),$$

$$\text{Var}[r_{ij}] = \frac{1}{N} \sum_{l=1}^N l^2 - \mu^2 = \frac{(N+1)(N-1)}{12} = \sigma^2 \quad (\text{say}),$$

$$\begin{aligned} \text{Cov}[r_{ij}, r_{i'j'}] &= \frac{1}{N(N-1)} \sum_{l \neq l'} ll' - \mu^2 \\ &= \frac{N}{N-1} \left\{ \left(\frac{1}{N} \sum_{l=1}^N l \right)^2 - \frac{1}{N^2} \sum_{l=1}^N l^2 \right\} - \mu^2 \\ &= \frac{N}{N-1} \left\{ \mu^2 - \frac{1}{N} (\sigma^2 + \mu^2) \right\} - \mu^2 = -\frac{1}{N-1} \sigma^2 \end{aligned}$$

$$E[S_N] = \left(\sum_{j=1}^k a_j \right) \mu$$

$$\text{Var}[S_N] = \sum_{j=1}^k \text{Var} \left[\sum_{i=1}^{n_j} \frac{a_j}{n_j} r_{ij} \right] + \sum_{j \neq j'}^k \text{Cov} \left[\sum_{i=1}^{n_j} \frac{a_j}{n_j} r_{ij}, \sum_{i=1}^{n_{j'}} \frac{a_{j'}}{n_{j'}} r_{ij'} \right]$$

$$\begin{aligned} \text{Var} \left[\sum_{i=1}^{n_j} \frac{a_j}{n_j} r_{ij} \right] &= \left(\frac{a_j}{n_j} \right)^2 \left\{ \sum_{i=1}^{n_j} \text{Var}[r_{ij}] + \sum_{i \neq i'}^{n_j} \text{Cov}[r_{ij}, r_{i'j}] \right\} \\ &= \left(\frac{a_j}{n_j} \right)^2 \left\{ n_j - \frac{n_j(n_j-1)}{N-1} \right\} \sigma^2 \end{aligned}$$

$$\text{Cov} \left[\sum_{i=1}^{n_j} \frac{a_j}{n_j} r_{ij}, \sum_{i=1}^{n_{j'}} \frac{a_{j'}}{n_{j'}} r_{ij'} \right] = \sum_{i=1}^{n_j} \sum_{i'=1}^{n_{j'}} \frac{a_j a_{j'}}{n_j n_{j'}} \left(-\frac{\sigma^2}{N-1} \right)$$

$$\text{Var}[S_N] = (a_1, \dots, a_k)' B \begin{pmatrix} a_1 \\ \vdots \\ a_k \end{pmatrix} \sigma^2,$$

$$B = (b_{jl}), \quad b_{jl} = \begin{cases} \frac{1}{n_j} - \frac{1}{N-1} \frac{n_j-1}{n_j} & (j=l) \\ -\frac{1}{N-1} & (j \neq l) \end{cases}$$

$$B = \frac{N}{N-1} \begin{pmatrix} \frac{1}{n_1} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{n_k} \end{pmatrix} - \frac{1}{N-1} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} (1, \dots, 1)$$

$$= \frac{N}{N-1} \begin{pmatrix} \frac{1}{\sqrt{n_1}} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sqrt{n_k}} \end{pmatrix} \left\{ I_k - \frac{1}{N} \begin{pmatrix} \sqrt{n_1} \\ \vdots \\ \sqrt{n_k} \end{pmatrix} (\sqrt{n_1}, \dots, \sqrt{n_k}) \right\} \begin{pmatrix} \frac{1}{\sqrt{n_1}} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sqrt{n_k}} \end{pmatrix}$$

Define

$$\mathbf{Z}_N = \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} r_{i1} \\ \vdots \\ \frac{1}{\sqrt{n_k}} \sum_{i=1}^{n_k} r_{ik} \end{pmatrix}, \quad D_N = \begin{pmatrix} \frac{1}{\sqrt{n_1}} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sqrt{n_k}} \end{pmatrix}, \quad \mathbf{n} = \begin{pmatrix} \sqrt{n_1} \\ \vdots \\ \sqrt{n_k} \end{pmatrix} \quad \text{and} \quad \Pi_N = I_k - \frac{1}{N} \mathbf{n} \mathbf{n}'$$

then

$$S_N = \mathbf{a}' D_N \mathbf{Z}_N,$$

$$E[S_N] = \frac{N+1}{2} \mathbf{a}' D_N \mathbf{n},$$

$$\text{Var}[S_N] = \frac{N(N+1)}{12} \mathbf{a}' D_N \Pi_N D_N \mathbf{a}.$$

Define

$$\begin{aligned}\Pi_N &= I_k - \frac{1}{N} \mathbf{n} \mathbf{n}' = \begin{matrix} H_N & H'_N \\ k \times (k-1) \end{matrix} \\ \mathbf{U}_N &= \sqrt{\frac{12}{N(N+1)}} H'_N \left(\mathbf{Z}_N - \frac{N+1}{2} \mathbf{n} \right) = \sqrt{\frac{12}{N(N+1)}} H'_N \mathbf{Z}_N \\ \Rightarrow \mathbf{U}_N &\rightarrow N_{k-1}(\mathbf{0}, I_{k-1}) \\ (\because \mathbf{s}' \mathbf{U}_N &= S_N - \mathbb{E}[S_N] \text{ with } \mathbf{a} = \sqrt{\frac{12}{N(N+1)}} D_N^{-1} H_N \mathbf{s}, \text{ and } \text{Var}[\mathbf{s}' \mathbf{U}_N] = \mathbf{s}' \mathbf{s}) \\ \Pi_N \mathbf{Z}_N &= \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} r_{i1} \\ \vdots \\ \frac{1}{\sqrt{n_k}} \sum_{i=1}^{n_k} r_{ik} \end{pmatrix} - \frac{1}{N} \mathbf{n} \sum_{j=1}^k \sum_{i=1}^{n_j} r_{ij} = \begin{pmatrix} \sqrt{n_1} (R_{.1} - \frac{N+1}{2}) \\ \vdots \\ \sqrt{n_k} (R_{.k} - \frac{N+1}{2}) \end{pmatrix} \\ \mathbf{U}'_N \mathbf{U}_N &= \frac{12}{N(N+1)} \mathbf{Z}'_N \Pi_N \mathbf{Z}_N = \frac{12}{N(N+1)} \sum_{j=1}^k n_j \left(R_{.j} - \frac{N+1}{2} \right)^2 = H\end{aligned}$$

Hence

$$H = \mathbf{U}_N \mathbf{U}_N \rightarrow \chi_{k-1}^2$$

Note

$$\begin{aligned}\mathbf{Z}'_N \Pi_N \mathbf{Z}_N &= \mathbf{Z}'_N \mathbf{Z}_N - \frac{1}{N} (\mathbf{Z}'_N \mathbf{n})^2 = \sum_{j=1}^k \frac{1}{n_j} \left(\sum_{i=1}^{n_j} r_{ij} \right)^2 - \frac{1}{N} \left(\frac{N(N+1)}{2} \right)^2 \\ H &= \left(\frac{12}{N(N+1)} \sum_{j=1}^k \frac{R_j^2}{n_j} \right) - 3(N+1)\end{aligned}$$

5.2 A distribution-free test for ordered alternatives (Jonckheere–Terpstra)

Testing problem

$$H_0 \text{ versus } H_2 : \tau_1 \leq \tau_2 \leq \cdots \leq \tau_k, \text{ with at least one strict inequality}$$

Test statistic

$$\begin{aligned}J &= \sum_{a=1}^{k-1} \sum_{b=a+1}^k U_{ab} \\ U_{ab} &= \sum_{i=1}^{n_a} \sum_{j=1}^{n_b} \phi(X_{ia}, X_{jb}), \quad 1 \leq a < b \leq k \\ \phi(x, y) &= \begin{cases} 1 & \text{if } x < y \\ 0 & \text{otherwise} \end{cases}\end{aligned}$$

Procedure

$$J \geq j_\alpha \Rightarrow \text{reject } H_0$$

Example for deriving the exact null distribution

$$R_{ij} : \text{rank of } X_{ij} \Rightarrow U_{ab} = \sum_{i=1}^{n_a} \sum_{j=1}^{n_b} \phi(X_{ia}, X_{jb}) = \sum_{i=1}^{n_a} \sum_{j=1}^{n_b} \phi(R_{ia}, R_{jb})$$

$$k = 3, n_1 = n_2 = 1, n_3 = 2$$

	I	II	III	J
(a)	1	2	3 4	5
(b)	2	1	3 4	4
(c)	1	3	2 4	4
(d)	3	1	2 4	3
(e)	1	4	2 3	3
(f)	4	1	3 2	2
(g)	2	3	1 4	3
(h)	3	2	1 4	2
(i)	2	4	1 3	2
(j)	4	2	1 3	1
(k)	3	4	1 2	1
(l)	4	3	1 2	0

j	0	1	2	3	4	5
$\Pr\{J = j\}$	$\frac{1}{12}$	$\frac{2}{12}$	$\frac{3}{12}$	$\frac{3}{12}$	$\frac{2}{12}$	$\frac{1}{12}$

Large-sample approximation

Under H_0

$$\mathbb{E}[J] = \frac{N^2 - \sum_{j=1}^k n_j}{4},$$

$$\text{Var}[J] = \frac{N^2(2N + 3) - \sum_{j=1}^k n_j^2(2n_j + 3)}{72}$$

$$\frac{J - \mathbb{E}[J]}{\sqrt{\text{Var}[J]}} \rightarrow N(0, 1) \quad (\min\{n_1, \dots, n_k\} \rightarrow \infty)$$

Proof

$$\mathbb{E}[J] = \sum_{a=1}^{k-1} \sum_{b=a+1}^k \mathbb{E}[U_{ab}] = \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a n_b \mathbb{E}[\phi(X, Y)] = \frac{N^2 - \sum_{a=1}^k n_a^2}{2} \mathbb{E}[\phi(X, Y)]$$

$$J^* = \sum_{a=1}^{k-1} \sum_{b=a+1}^k U_{ab}^*$$

$$U_{ab}^* = n_b \sum_{i=1}^{n_a} \phi_{10}(X_{ia}) + n_a \sum_{j=1}^{n_b} \phi_{01}(X_{jb})$$

$$\phi_{10}(x) = \mathbb{E}[\phi(X, Y)|X = x], \quad \phi_{01}(y) = \mathbb{E}[\phi(X, Y)|Y = y]$$

Formulas

$X, Y, Z \stackrel{i.i.d.}{\sim} F$

$$\begin{aligned} \sigma_{10}^2 &= \text{Var}[\phi_{10}(X)], \quad \sigma_{01}^2 = \text{Var}[\phi_{01}(X)], \quad \sigma_{11}^2 = \text{Cov}[\phi_{10}(X), \phi_{01}(X)] \\ \Rightarrow \text{Cov}[\phi(X, Y), \phi(X, Z)] &= \sigma_{10}^2, \quad \text{Cov}[\phi(X, Y), \phi(Z, Y)] = \sigma_{01}^2, \\ \text{Cov}[\phi(X, Y), \phi(Y, Z)] &= \sigma_{11}, \quad \text{Cov}[\phi(X, Y), \phi(Z, X)] = \sigma_{11}, \\ \text{Cov}[\phi(X, Y), \phi_{10}(X)] &= \sigma_{10}^2, \quad \text{Cov}[\phi(X, Y), \phi_{10}(Y)] = \sigma_{11}, \\ \text{Cov}[\phi(X, Y), \phi_{01}(X)] &= \sigma_{11}, \quad \text{Cov}[\phi(X, Y), \phi_{01}(Y)] = \sigma_{01}^2, \end{aligned}$$

$$\begin{aligned} \text{Var}[U_{ab}] &= n_a n_b (\sigma^2 - \sigma_{10}^2 - \sigma_{01}^2) + n_a n_b^2 \sigma_{10}^2 + n_a^2 n_b \sigma_{01}^2 \\ \text{Cov}[U_{ab}, U_{ac}] &= n_a n_b n_c \sigma_{10}^2, \quad (b \neq c) \\ \text{Cov}[U_{ab}, U_{bc}] &= n_a n_b n_c \sigma_{11}, \\ \text{Cov}[U_{ab}, U_{cb}] &= n_a n_b n_c \sigma_{01}^2, \quad (a \neq c) \end{aligned}$$

$$\begin{aligned} \text{Cov}[U_{ab}, U_{ab}^*] &= n_a n_b^2 \sigma_{10}^2 + n_a^2 n_b \sigma_{01}^2, \\ \text{Cov}[U_{ab}, U_{ac}^*] &= n_a n_b n_c \sigma_{10}^2, \quad (b \neq c) \\ \text{Cov}[U_{ab}, U_{bc}^*] &= \text{Cov}[U_{ab}, U_{ca}^*] = n_a n_b n_c \sigma_{11} \\ \text{Cov}[U_{ab}, U_{cb}^*] &= n_a n_b n_c \sigma_{01}^2, \quad (a \neq c) \end{aligned}$$

$$\begin{aligned} \text{Var}[U_{ab}^*] &= n_a n_b^2 \sigma_{10}^2 + n_a^2 n_b \sigma_{01}^2 \\ \text{Cov}[U_{ab}^*, U_{ac}^*] &= n_a n_b n_c \sigma_{10}^2, \quad (b \neq c) \\ \text{Cov}[U_{ab}^*, U_{bc}^*] &= n_a n_b n_c \sigma_{11}, \\ \text{Cov}[U_{ab}^*, U_{cb}^*] &= n_a n_b n_c \sigma_{01}^2, \quad (a \neq c) \end{aligned}$$

$$\begin{aligned} \text{Var}[J - J^*] &= \sum_{a=1}^{k-1} \sum_{b=a+1}^k \text{Var}[U_{ab} - U_{ab}^*] \\ &\quad + \sum_{a=1}^{k-1} \sum_{b=a+1}^k \sum_{a'=1}^{k-1} \sum_{b'=a'+1}^k \text{Cov}[(U_{ab} - U_{ab}^*), (U_{a'b'} - U_{a'b'}^*)] \\ &\quad \text{with } (a, b) \neq (a', b') \\ &= \left\{ \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a n_b \right\} (\sigma^2 - \sigma_{10}^2 - \sigma_{01}^2) \end{aligned}$$

$$\begin{aligned}
\text{Var}[J] &= \sum_{a=1}^{k-1} \sum_{b=a+1}^k \text{Var}[U_{ab}] + \sum_{a=1}^{k-1} \sum_{b=a+1}^k \sum_{\substack{a'=1 \\ (a,b) \neq (a',b')}}^{k-1} \sum_{b'=a'+1}^k \text{Cov}[U_{ab}, U_{a'b'}] \\
&= \left\{ \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a n_b \right\} (\sigma^2 - \sigma_{10}^2 - \sigma_{01}^2) + 2 \left\{ \sum_{a=1}^{k-2} \sum_{b=a+1}^{k-1} \sum_{c=b+1}^k n_a n_b n_c \right\} \sigma_{11} \\
&\quad + \left\{ \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a n_b^2 + 2 \sum_{a=1}^{k-2} \sum_{b=a+1}^{k-1} \sum_{c=b+1}^k n_a n_b n_c \right\} \sigma_{10}^2 \\
&\quad + \left\{ \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a^2 n_b + 2 \sum_{a=1}^{k-2} \sum_{b=a+1}^{k-1} \sum_{c=b+1}^k n_a n_b n_c \right\} \sigma_{01}^2 \\
\text{Var}[J^*] &= 2 \left\{ \sum_{a=1}^{k-2} \sum_{b=a+1}^{k-1} \sum_{c=b+1}^k n_a n_b n_c \right\} \sigma_{11} \\
&\quad + \left\{ \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a n_b^2 + 2 \sum_{a=1}^{k-2} \sum_{b=a+1}^{k-1} \sum_{c=b+1}^k n_a n_b n_c \right\} \sigma_{10}^2 \\
&\quad + \left\{ \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a^2 n_b + 2 \sum_{a=1}^{k-2} \sum_{b=a+1}^{k-1} \sum_{c=b+1}^k n_a n_b n_c \right\} \sigma_{01}^2
\end{aligned}$$

$$\phi_{10}(x) = \mathbb{E}[\phi(x, Y)] = 1 - F(x), \quad \phi_{01}(y) = \mathbb{E}[\phi(X, y)] = F(y)$$

$$\mathbb{E}[\phi(X, Y)] = \mathbb{E}[\phi_{10}(X)] = \mathbb{E}[\phi_{01}(X)] = \frac{1}{2} \quad \because F(X) \sim \mathcal{U}(0, 1)$$

$$\mathbb{E}[J] = \frac{N^2 - \sum_{a=1}^k n_a^2}{4}$$

$$\sigma_{10}^2 = \sigma_{01}^2 = \frac{1}{12}$$

$$\sigma_{11} = \mathbb{E}[\{1 - F(X)\}F(X)] - \frac{1}{4} = -\frac{1}{12}$$

$$\sigma^2 = \mathbb{E}[\phi(X, Y)^2] - \frac{1}{4} = \mathbb{E}[\phi(X, Y)] - \frac{1}{4} = \frac{1}{4}$$

$$N^3 = \sum_{a=1}^k \sum_{b=1}^k \sum_{c=1}^k n_a n_b n_c = 6 \sum_{\substack{a=1 \\ a < b < c}}^k \sum_{b=1}^k \sum_{c=1}^k n_a n_b n_c + 3 \sum_{\substack{a=1 \\ a < b}}^k \sum_{b=1}^k (n_a^2 n_b + n_a n_b^2) + \sum_{a=1}^k n_a^3$$

$$\text{Var}[J] = \frac{1}{12} \left\{ \frac{N^2 - \sum_{a=1}^k n_a^2}{2} + \frac{N^3 - \sum_{a=1}^k n_a^3}{3} \right\} = \frac{N^2(2N + 3) - \sum_{a=1}^k n_a^2(2n_a + 3)}{72}$$

$$\text{Var}[J^*] = \frac{N^3 - \sum_{a=1}^k n_a^3}{36}, \quad \text{Var}[J - J^*] = \frac{N^2 - \sum_{a=1}^k n_a^2}{24}$$

$$\Rightarrow \frac{J - \mathbb{E}[J]}{\sqrt{\text{Var}[J]}} - \frac{J^* - \mathbb{E}[J^*]}{\sqrt{\text{Var}[J^*]}} \xrightarrow{P} 0$$

$$\begin{aligned}
U_{ab}^* &= n_b \sum_{i=1}^{n_a} \phi_{10}(X_{ia}) + n_a \sum_{j=1}^{n_b} \phi_{01}(X_{ib}) \\
J^* &= \sum_{a < b} U_{ab}^* = \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_b \sum_{i=1}^{n_a} \phi_{10}(X_{ia}) + \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a \sum_{i=1}^{n_b} \phi_{01}(X_{ib}) \\
&= \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_b \sum_{i=1}^{n_a} \phi_{10}(X_{ia}) + \sum_{a=2}^k \sum_{b=1}^{a-1} n_b \sum_{i=1}^{n_a} \phi_{01}(X_{ia}) \\
&= \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_b \sum_{i=1}^{n_a} \{1 - F(X_{ia})\} + \sum_{a=2}^k \sum_{b=1}^{a-1} n_b \sum_{i=1}^{n_a} F(X_{ia}) \\
&= \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a n_b + \sum_{a=1}^k \left\{ - \sum_{b=a+1}^k n_b + \sum_{b=1}^{a-1} n_b \right\} \sum_{i=1}^{n_a} F(X_{ia}) \\
c_{ia} &= \left\{ - \sum_{b=a+1}^k n_b + \sum_{b=1}^{a-1} n_b \right\}, \quad (i = 1, \dots, n_a; a = 1, \dots, k) \\
\bar{c}_{ia} &= \frac{1}{N} \sum_{a=1}^k \sum_{i=1}^{n_a} c_{ia} = \frac{1}{N} \left\{ - \sum_{a=1}^{k-1} \sum_{b=a+1}^k n_a n_b + \sum_{a=2}^k \sum_{b=1}^{a-1} n_a n_b \right\} = 0 \\
\sum_{a=1}^k \sum_{i=1}^{n_a} c_{ia}^2 &= \frac{N^3 - \sum_{a=1}^k n_a^3}{3}, \quad |c_{ia}| < N \\
\frac{\max_{i,a} c_{ia}^2}{\sum_{a=1}^k \sum_{i=1}^{n_a} c_{ia}^2} &\leq \frac{3N^2}{N^3 - \sum_{a=1}^k n_a^3} \rightarrow 0
\end{aligned}$$