# Non parametric Inference

Gibbons and chakrobarty

**Rank**

Suppose $X_\alpha$ be the $\alpha$th obsⁿ for a set of $n$, obsⁿ's; $\alpha = 1,2,\ldots, n$ from a continuous distⁿ $F_x(n)$.

$R_\alpha$ : rank of $X_\alpha$
: # of obsⁿ's $\leq X_\alpha$  or $\alpha$th smallest obsⁿ.

Due to continuity ranks are distinct with prob 1.

Rank is an ordered permutation.

Rank vector $\underline{R} = \underline{r}$ is a ordered permutation.

$$\underline{R} = (R_1 \; R_2 \; R_3 \ldots R_n)$$

Rank of 1st obs ↓ ↑ Rank of $n$th obs

$\underline{R}$ is the random vector random permutation.

**Remmark**: $Pr\{\underline{R} = \underline{r}\} = \dfrac{1}{n!}$

$$Pr\{R_\alpha = r_\alpha\} = \dfrac{1}{n} \; ; \; \alpha = 1,2,\ldots, n$$

Marginal distⁿ of rank is a discret uniform distⁿ.

**Remark**: $Pr\{R_\alpha = r_\alpha \cap R_\beta = r_\beta\} = \dfrac{1}{n(n-1)} \quad \alpha \neq \beta$

**Remark**:

$$\begin{bmatrix} E(R_\alpha) = \dfrac{n+1}{2} \\[6pt] V(R_\alpha) = \dfrac{n^2-1}{12} \\[6pt] Cov(R_\alpha, R_\beta) = -\dfrac{n+1}{12} \end{bmatrix}$$

⟨$R_\alpha$ and $R_\beta$ are not. ind⟩

⟨Also they are negatively correlated⟩

$$\left\langle Corrⁿ(R_\alpha, R_\beta) = \dfrac{-\dfrac{n+1}{12}}{\dfrac{n^2-1}{12}} = -\dfrac{1}{n-1}\right\rangle ; \boxed{n \geq 1}$$

## Linear rank statistic

Let $\underline{a} = (a_1, a_2, \ldots, a_n)$
and $\underline{b} = (b_1, b_2, \ldots, b_n)$ be two sets of coeff. (constants) based on $n$ natural number.

Let, $\underline{R} = (R_1, R_2, \ldots, R_n)$ be the random permutation of $\{1, 2, \ldots, n\}$. Then linear rank statistic is,

$$\left\langle T = \sum_{\alpha=1}^{n} a_\alpha \, b_{R_\alpha} = a_1 b_{R_1} + a_2 b_{R_2} + \cdots + a_n b_{R_n}\right\rangle$$

$a_\alpha'$s are known as regression constants and $b_{R_\alpha}'$s are scores constants.

Note that, joint dist$^n$ of rank is independent of dist$^n$ $F$ from which obs$^n$'s come, the dist$^n$ of $T$ is ind. of any $F$. Hence $T$ can be used to provide dist$^n$ free (non-parametric test)

Also $(R_1, R_2, \ldots, R_n) \overset{d}{=} (n-R_1+1, n-R_2+1, \ldots, n-R_n+1)$

■ **Mean and variance of $T$**

$$E(T) = E\left[\sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha}\right]$$

$$= \sum_{\alpha=1}^{n} a_\alpha E(b_{R_\alpha})$$

$$\therefore E(b_{R_\alpha}) = \frac{1}{n}\left[b_1 + b_2 + \cdots + b_n\right] \quad \left[\begin{matrix}\text{As the value} \\ \text{of } R_\alpha \text{ can be} \\ 1, 2, \ldots n\end{matrix}\right]$$

$$= \bar{b}$$

$$\therefore E(T) = \sum_{\alpha=1}^{n} a_\alpha \bar{b}$$

$$= n\bar{a}\,\bar{b} \quad \left[\text{where } \bar{a} = \frac{1}{n}\sum_{\alpha=1}^{n} a_\alpha\right]$$

$$V(T) = V\left(\sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha}\right)$$

$$= V\left(\sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha} - n\bar{a}\bar{b}\right)$$

$$= V\left(\sum_{\alpha=1}^{n} (a_\alpha - \bar{a}) b_{R_\alpha}\right)$$

$$= \sum_{\alpha=1}^{n} (a_\alpha - \bar{a})^2 V(b_{R_\alpha})$$

$$\qquad + \sum\sum_{\alpha \neq \beta} (a_\alpha - \bar{a})(a_\beta - \bar{a}) \text{cov}(b_{R_\alpha}, b_{R_\beta})$$

$$= V(b_{R_\alpha}) \sum (a_\alpha - \bar{a})^2 + \text{cov}(b_{R_\alpha}, b_{R_\beta}) \sum\sum_{\alpha \neq \beta}(a_\alpha - \bar{a})(a_\beta - \bar{a})$$

WLG take $\alpha = 1$,

$$V(b_{R_1}) = E(b_{R_1}^2) - \left[E(b_{R_1})\right]^2$$

$$= \frac{1}{n}\left[b_1^2 + b_2^2 + \cdots + b_n^2\right] - \bar{b}^2$$

$$= \frac{1}{n}\sum_{\alpha=1}^{n} b_\alpha^2 - \bar{b}^2$$

$$= \frac{1}{n}\sum_{\alpha=1}^{n} (b_\alpha - \bar{b})^2$$

$$V(b_{R_\alpha}) = \frac{1}{n}\sum_{\alpha=1}^{n}(b_\alpha - \bar{b})^2$$

$$\text{Cov}(b_{R_\alpha}, b_{R_\beta})$$

WLG $\alpha = 1, \beta = 2$

$$\text{Cov}(b_{R_1}, b_{R_2}) = E\left[(b_{R_1} - \bar{b})(b_{R_2} - \bar{b})\right]$$

$$= \sum_{\substack{R_1=1 \\ (\alpha)}}^{n}\sum_{\substack{R_2=1 \\ R_1 \neq R_2(\beta)}}^{n} P\left[R_1 = \alpha \wedge R_2 = \beta\right] \cdot (b_\alpha - \bar{b}) \cdot (b_\beta - \bar{b})$$

$$= \frac{1}{n(n-1)}\sum_{\substack{R_1 \neq R_2 \\ (\alpha) \quad (\beta)}}(b_\alpha - \bar{b})(b_\beta - \bar{b})$$

$$= \frac{1}{n(n-1)}\left[\sum_{R_1(\alpha)}(b_\alpha - \bar{b})\left\{\sum_{R_2(\beta)}(b_\beta - \bar{b}) - (b_\alpha - \bar{b})\right\}\right]$$

$$= \frac{1}{n(n-1)}\left[\left\{\sum_\alpha(b_\alpha - \bar{b})\right\}^2 - \sum_\alpha(b_\alpha - \bar{b})^2\right]$$

$$= -\frac{1}{n(n-1)}\sum_{\alpha=1}^{n}(b_\alpha - \bar{b})^2 = \boxed{-\frac{V(b_{R_1})}{n-1}}$$

$$\therefore V(T) = V(b_{R_1})\sum(a_\alpha - \bar{a})^2 - \frac{V(b_{R_1})}{(n-1)}\sum\sum_{\alpha \neq \beta}(a_\alpha - \bar{a})(a_\beta - \bar{a})$$

$$= \frac{1}{n}\sum_{\alpha=1}^{n}(b_\alpha - \bar{b})^2\sum_{\alpha=1}^{n}(a_\alpha - \bar{a})^2$$
$$- \frac{\sum(b_\alpha - \bar{b})^2}{n(n-1)}\sum\sum_{\alpha \neq \beta}(a_\alpha - \bar{a})(a_\beta - \bar{a})$$

$$= \frac{1}{n}\sum_{\alpha=1}^{n}(b_\alpha - \bar{b})^2\sum_{\alpha=1}^{n}(a_\alpha - \bar{a})^2$$
$$= \frac{\sum(b_\alpha - \bar{b})^2}{n(n-1)}\left[\sum_\alpha(a_\alpha - \bar{a})\left\{\sum_\beta(a_\beta - \bar{a}) - (a_\alpha - \bar{a})\right\}\right]$$

$$= \frac{1}{n}\sum_{\alpha=1}^{n}(b_\alpha - \bar{b})^2\sum_{\alpha=1}^{n}(a_\alpha - \bar{a})^2 + \frac{\sum(b_\alpha - \bar{b})^2}{n(n-1)}\sum(a_\alpha - \bar{a})^2$$

$$= \left[\frac{1}{n} + \frac{1}{(n-1)n}\right]\sum(a_\alpha - \bar{a})^2\sum(b_\alpha - \bar{b})^2$$

$$= \frac{1}{n-1}\sum_{\alpha=1}^{n}(a_\alpha - \bar{a})^2\sum_{\alpha=1}^{n}(b_\alpha - \bar{b})^2$$

**Asymptotic dist$^n$ of linear rank statistic :-**

i) Wald Wolfowitz cond$^n$
ii) Noether cond$^n$
iii) Hoeffding cond$^n$ 

$$\left\langle \frac{T-E(T)}{\sqrt{V_i(T)}} \xrightarrow{L} N(0,1) \right\rangle$$

**Inverse permutation**

Every permutation $P$ has its mirror image $P'$

such that, $\left\langle P \circ P' = I^+ \right\rangle = \{1, 2, \dots, n\}$

$\uparrow$ operation

That mirror image is called inverse permutation

▶ Inverse permutation in Ranking theory is called antirank $\underset{\sim}{Q} = (Q_1, Q_2 \dots Q_n)$

$$\left\langle \underset{\sim}{R} \circ \underset{\sim}{Q} = I \right\rangle$$

**Antirank**

In antiranking, number and number of the place the obs$^n$ occupies (Rank) is exchanged.

R ( 3 . 1  4  2)

Q ( 2  4  1  3)

| Number | | Place |
|--------|-----|-------|
| 3 | $\leftrightarrow$ | 1 |
| 1 | $\leftrightarrow$ | 3 |
| 1 | $\leftrightarrow$ | 2 |
| 2 | $\leftrightarrow$ | 1 |
| 4 | $\leftrightarrow$ | 3 |
| 3 | $\leftrightarrow$ | 4 |

R ( 2  3  4  5  1)

Q ( 5  1  2  3  4 )

**Remark:**

We know ranks are not independent, but the values in inverse permutation (antiranks) are ind.

Moreover the joint dist$^n$ of any permutation remains the same. Therefore the joint dist$^n$ of ranking and joint dist$^n$ of antiranking are same.

$$\left\langle \begin{array}{l} P(\underset{\sim}{R} = \underset{\sim}{r}) = \frac{1}{n!} \\ P(\underset{\sim}{Q} = \underset{\sim}{d}) = \frac{1}{n!} \end{array} \right\rangle$$

[Hence, for finding dist$^n$ of linear rank statistic anti-ranking process is more convenient.]

$Cov(R_\alpha, R_\beta)$

$= E(R_\alpha \cdot R_\beta) - E(R_\alpha)E(R_\beta)$

$= \sum\limits_{\alpha \neq \beta}^{n}\sum\limits^{n} \alpha\beta \cdot \frac{1}{n(n-1)} - \left(\frac{n+1}{2}\right)^2$

$= \frac{1}{n(n-1)}\left[(\sum \alpha)^2 - \sum \alpha^2\right] - \left(\frac{n+1}{2}\right)^2$

$= \frac{1}{n(n-1)}\left[\frac{n^2(n+1)^2}{4} - \frac{n(n+1)(2n+1)}{6}\right] - \left(\frac{n+1}{2}\right)^2$

$= -\frac{n+1}{12}$

---

**Remark**

$T = \sum\limits_{\alpha=1}^{n} a_\alpha b_{R_\alpha} \overset{d}{\equiv} \sum\limits_{\alpha=1}^{n} a_{\mathcal{D}_\alpha} b_\alpha$

$\mathcal{D}_\alpha \longrightarrow$ Antirank of the $\alpha^{th}$ obs$^n$.

$R = \begin{pmatrix} 3 & 1 & 2 \\ \underset{R_1}{} & \underset{R_2}{} & \underset{R_3}{} \end{pmatrix}$   $T = \sum\limits_{\alpha=1}^{n} a_\alpha b_{R_\alpha} = a_1 b_3 + a_2 b_1 + a_3 b_2$

$\mathcal{D} = \begin{pmatrix} 2 & 3 & 1 \\ \underset{\mathcal{D}_1}{} & \underset{\mathcal{D}_2}{} & \underset{\mathcal{D}_3}{} \end{pmatrix}$   $T' = \sum\limits_{\alpha=1}^{3} a_{\mathcal{D}_\alpha} b_\alpha = a_2 b_1 + a_3 b_2 + a_1 b_3$

**Result :**

For a linear rank statistic T, if either $a_\alpha + a_{n-\alpha+1}$
or $b_{R_\alpha} + b_{R_{(n-\alpha+1)}}$ is a constant $\forall \alpha$, the dist$^n$ of T is symmetric ~~about~~ about its mean.

**Pf :** We know that $(R_1, R_2, \ldots, R_n) \overset{d}{\equiv} (n-R_1+1, n-R_2+1, \ldots, n-R_n+1)$

Assume $b_{R_\alpha} + b_{R_{n-\alpha+1}} = k \quad \forall \alpha$ where $k$ is a constant.

$\bar{b} = \frac{k}{2}$

$\therefore T \overset{d}{\equiv} \sum\limits_{\alpha=1}^{n} a_\alpha b_{R_\alpha} \overset{d}{\equiv} \sum\limits_{\alpha=1}^{n} a_\alpha b_{n-R_\alpha+1}$

$\overset{d}{\equiv} \sum\limits_{\alpha=1}^{n} a_{\mathcal{D}_\alpha} b_{n-\alpha+1}$

$\overset{d}{\equiv} \sum\limits_{\alpha=1}^{n} a_{\mathcal{D}_\alpha} (k - b_\alpha) \quad \left[\because b_\alpha + b_{n-\alpha+1} = k\right]$

$\overset{d}{\equiv} \sum\limits_{\alpha=1}^{n} a_{\mathcal{D}_\alpha} k - \sum\limits_{\alpha=1}^{n} a_{\mathcal{D}_\alpha} b_\alpha$

$T \overset{d}{\equiv} kn\bar{a} - \sum\limits_{\alpha=1}^{n} a_\alpha b_{R_\alpha}$

$\equiv k.n\bar{a} - T \quad \forall k$

Choose, $k = 2\bar{b}$   $T \overset{d}{\equiv} 2n\bar{a}\bar{b} - T$

$\Rightarrow T - n\bar{a}\bar{b} \overset{d}{\equiv} n\bar{a}\bar{b} - T \quad$ <Hence the proof.>

# Sign Test

Let $x_1, x_2, \ldots, x_n$ be a random sample from a cont. distribution $F_X(n)$. Let $\mu$ be the median of the dist$^n$.

We are to test, $H_0 : \mu = \mu_0$ (a fixed constant)

$H_1 : \mu > \mu_0$ (median of $F_X(n) > \mu_0$)

Let, $S =$ no. of positive obs$^n$ $(x_i - \mu_0 > 0)$

Let us propose a linear rank statistic for testing the above,

Define, 
$$\left[ a_\alpha = \begin{cases} 1 & \text{if} \quad x_\alpha > \mu_0 \\ 0 & \quad 0.\omega \end{cases} \right]$$

and, $b_{R_\alpha} = 1 \quad \forall \alpha$

$$\left\langle T = \sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha} = \sum_{\substack{\alpha=1 \\ x_\alpha > \mu_0}}^{n} a_\alpha = S \right\rangle$$

$$E_{H_0}(T) = E_{H_0}\left( \sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha} \right)$$

$$= E_{H_0}\left( \sum_{\alpha=1}^{n} a_\alpha \right)$$

$$= \sum_{\alpha=1}^{n} E_{H_0}(a_\alpha)$$

$$= \sum_{\alpha=1}^{n} 1 \cdot P_{H_0}\left[ x_\alpha > \mu_0 \right]$$

$$= \frac{n}{2}$$

$\left[ E_{H_0}[a_\alpha] = \frac{1}{2} \right]$

$$V_{H_0}(T) = V_{H_0}\left( \sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha} \right)$$

Now,

$V(a_\alpha)$

$$= V_{H_0}\left( \sum_{\alpha=1}^{n} a_\alpha \right)$$

$$= E[a_\alpha^2] - [E(a_\alpha)]^2$$

$$= \sum_{\alpha=1}^{n} V(a_\alpha) + \sum\sum_{\alpha \neq \beta} cov(a_\alpha, a_\beta)$$

$$= 1 \cdot P[x_\alpha > \mu_0] - \{1 \cdot P[x_\alpha > \mu_0]\}^2$$

$$= \frac{1}{2} - \frac{1}{4} = \frac{1}{4} \quad \ldots (i)$$

$$= \frac{n}{4} \quad [\text{Using (i) and (ii)}]$$

$cov(a_\alpha, a_\beta)$

$$= E(a_\alpha a_\beta) - E(a_\alpha) E(a_\beta)$$

$$= 1 \cdot P[x_\alpha > \mu_0 \cap x_\beta > \mu_0] - \frac{1}{4}$$

$$= P[x_\alpha > \mu_0] \cdot P[x_\beta > \mu_0] - \frac{1}{4}$$

$$= \frac{1}{2} \cdot \frac{1}{2} - \frac{1}{4} \quad [\because x_\beta \text{ and } x_\beta \text{ are ind}]$$

$$= 0 \quad \ldots (ii)$$

$\left\langle \text{Therefore } S \text{ is a particular case of linear rank statistic} \right\rangle$

$\left\langle \text{Also it is a symmetric linear rank statistic around mean } \frac{n}{2} \right\rangle$

## Critical fun of sign test:

Let us propose a test fun; $\phi(s)$

$$\phi(s) = \begin{cases} 1 & \text{if} & s - n/2 > c \\ \gamma & & s - n/2 = c \\ 0 & & o.w \end{cases}$$

$c$ and $\gamma$ are determined from size condn.

**Remark:** For sign test, if zero difference occurs for any $x_\alpha = \mu_0$, under cont. distn assumpn, this does not create any problem as $P[x_\alpha = \mu_0] = 0$.

But in practical, zero difference can be avoided by ignoring them and reducing the sample size simultaneously.

**Result:**

Sign test is a UMP test.

**proof:**

Propose a test fun as,

$$\phi(x) = \begin{cases} 1 & \text{if} & \prod_{i=1}^{n} f_1(x_i) > k \prod_{i=1}^{n} f_0(x_i) \\ \gamma & & \prod f_1 = k \prod f_0 \\ 0 & & \prod f_1 < k \prod f_0 \end{cases}$$

For any c.d.f $F_x(x)$ with p.d.f $f_x(x)$, define

$$f(x) = F(0) f^-(x) + (1 - F(0)) f^+(x) \quad \ldots \text{(1)}$$

For testing, $H_0 : \mu = \mu_0$

against $H_1 : \mu > \mu_0$

$F_x(\mu_0) = 1/2$. Let us make a transformation $x_\alpha - \mu$ such that $\mu_0 = 0$ and the corresponding test will be,

$H_0 : \mu' = 0 \equiv F_x(0) = 1/2$

$H_1 : \mu' > 0$ where, $\langle F_x(0) = 1/2 \rangle$

$\equiv F_x(0) < 1/2$

In (1), $f^-(x) = \begin{cases} f(x), & \text{if } x \leq 0 \\ \dfrac{f(x)}{F(0)}, & \text{if } x > 0 \\ 0 \end{cases}$

and, $f^+(x) = \begin{cases} 0 & \text{if } x < 0 \\ \dfrac{f(x)}{1 - F(0)} & \text{if } x > 0 \end{cases}$

Under $H_0$, $l_0(m) = \frac{1}{2} l_0^-(m) + \frac{1}{2} l_0^+(m)$, $\langle l_0(m) = $ p.d.f under $H_0 \rangle$

$[F_x(0) = \frac{1}{2}]$

Under $H_1$, $l_1(m) = F_1(0) l_1^-(m) + (1 - F_1(0)) l_1^+(m)$.

Let us reframe the test fun as per $l^+, l^-$.

$$\phi(m) = \begin{cases} 1 & \text{if } \boxed{H_1} \quad \frac{\prod l_1(n_\alpha) \prod l_1(n_\alpha)}{\alpha \notin (\alpha_1, \alpha_2 \dots \alpha_s) \; \alpha \in (\alpha_1, \alpha_2, \dots \alpha_s)} \\ \gamma & = \quad > k \; \frac{\prod l_0^-(n_\alpha) \cdot \prod l_0^+(n_\alpha)}{\alpha \notin (\alpha_1, \dots, \alpha_s) \; \alpha \in (\alpha_1, \dots, \alpha_s)} \\ 0 & < \end{cases}$$

$$\Rightarrow \phi(m) = \begin{cases} 1 & \text{if } \prod_{\alpha \notin (\alpha_1, \dots \alpha_s)} \frac{l_1(n_\alpha)}{F_1(0)} \prod_{\alpha \in (\alpha_1, \dots \alpha_s)} \frac{l_1(n_\alpha)}{1 - F_1(0)} \\ \gamma & = \quad > k \prod_{\alpha \notin} \frac{l_0(n_\alpha)}{\frac{1}{2}} \prod_{\alpha \in} \frac{l_0(n_\alpha)}{\frac{1}{2}} \\ 0 & < \end{cases}$$

$$\Rightarrow \phi(m) = \begin{cases} 1 & \frac{\prod l_1(n_\alpha)}{[F_1(0)]^{n-s}} \cdot \frac{\prod l_1(n_\alpha)}{(1 - F_1(0))^{s}} \\ \gamma & = \quad > 2^n k \cdot \prod_{\alpha \notin} l_0(n_\alpha) \prod_\alpha l_0(n_\alpha) \\ 0 & < \end{cases}$$

$$\Rightarrow \phi(m) = \begin{cases} 1 & \frac{\prod_\alpha l_1(n_\alpha)}{\prod l_0(n_\alpha)} > 2^n k (F_1(0))^{n-s} (1 - F(0))^{s} \\ \gamma & = \\ 0 & < \end{cases}$$

$$\Rightarrow \phi(m) = \begin{cases} 1 & \left(\frac{F_1(0)}{1 - F_1(0)}\right)^{s} > k^* \cdot \frac{\prod_{\alpha=1}^{n} l_0(n_\alpha)}{\prod_{\alpha=1}^{n} l_1(n_\alpha)} \\ \gamma & = \\ 0 & < \end{cases}$$

$\therefore$ For fixed $x_1, x_2, \dots x_n$, $\prod \frac{l_0}{l_1}$ is a constant.

$$\phi(m) = \begin{cases} 1 & \left(\frac{F_1(0)}{1 - F_1(0)}\right)^{s} > k^{**} \\ \gamma & = k^{**} \\ 0 & < k^{**} \end{cases}$$

∴ $\phi(m)$ can be written in terms of $S = $ no. of positive des's.

Hence,

$$\phi(s) = \begin{cases} 1 & s > k' \\ \gamma & s = k' \\ 0 & s < k' \end{cases}$$

As the above satisfies N-P test construc⁺
it is a UMP test.

(Practical) < Problems on Non-par. Inference)

Suppose that each of 13 randomly choosen female registerd
voters was asked to indicate if she is going to vote for
candidate A or candidate in the upcoming election. The
result shows that 9 of the subjects prefers A. Is
this sufficient evidence to conclude that candidate A
is prefered to B by female voters.
                Draw the power curve taking at least

8 points.

$\Rightarrow$  $\underset{H_0}{S \sim} Bin(13, \frac{1}{2})$          $S = 9$

We test,
$H_0 : p = \frac{1}{2}$
ag. $H_1 : p > \frac{1}{2}$

| $k_\alpha$ | $P(S = k_\alpha)$ | $P(S \le k_\alpha)$ |
|---|---|---|
| 6 | | |
| 7 | 0.20947 | 0.70906 |
| 8 | 0.1571 | 0.867 |
| 9 | 0.08728 | 0.954 → It is the point. |
| 10 | 0.03491 | 0.989 |

Test is constructed as,
$$\phi(s) = \begin{cases} 1 & s > k_\alpha \\ \gamma & s = k_\alpha \\ 0 & s < k_\alpha \end{cases}$$

$\therefore \gamma = \dfrac{P(S \le 9) - 0.95}{P(S = 9)}$   Here, $k_\alpha = 9$

$= \dfrac{0.954 - 0.95}{0.08728}$

$= 0.044$

The test is constructed as
$$\phi(s) = \begin{cases} 1 & s > 9 \\ 0.044 & s = 9 \\ 0 & o.w \end{cases}$$

● Why sign test is a non-parametric test?

Here $S \underset{H_0}{\sim} Bin(n, \frac{1}{2})$ but $X_i$'s are cont.

∴ Dist$^n$ of $S$ does not depend on the parent pop$^n$.

Wilcoxon Signed-rank Test: (Test of loca$^n$)

Let $X_1, X_2, \ldots, X_n$ be the r.s. from a cont. c.d.f $F(\cdot)$ and with median $\mu$.

We are to test, $H_0 : \mu = \mu_0$.

First consider the difference $D_\alpha = X_\alpha - \mu_0$; Clearly, the

differences are distributed symmetrically under $H_0$.

$$F_D(-c) = Pr(D_\alpha \leq -c) = Pr(D_\alpha > c) = 1 - F_D(c)$$

With the assumption of a cont. pop$^n$, zero or tied

differences can be avoided by dropping them.

Next we order absolute $D_\alpha$'s, i.e $|D_\alpha|$'s inc. ly

(from smallest to the largest), $|D_1|, |D_2|, \ldots, |D_n|$

Then we rank $|D_\alpha|$'s.

The test statistic is $T^+ =$ Sum of ranks for positive obs$^n$ $(D_\alpha > 0)$.

$T^- =$ Sum of ranks for negative obs$^n$ $(D_\alpha < 0)$

● $T^+ + T^- =$ Sum of all possible ranks $= 1 + 2 + \ldots + n = \frac{n(n+1)}{2}$

$$\left\langle T^+ = \frac{n(n+1)}{2} - T^- \right\rangle \to \text{It means } T^+ \text{ and } T^- \text{ are lin. related}$$

Tests based on $T^+$ only, $T^-$ only or $T^+ - T^-$ are

all equivalent.

Let us define the rank of $|D_\alpha|$, $R_\alpha^+$. $T^+$ is a linear rank statistic.

Redefine, $T^+ = \sum_{\alpha=1}^{n} Z_\alpha R_\alpha^+$ where, $Z_\alpha = \begin{cases} 1 & \text{if } D_\alpha > 0 \equiv X_\alpha > \mu_0 \\ 0 & \text{o.w} \end{cases}$

⟨Check it is a linear rank statistic⟩

Similarly. $T^- = \sum (1 - Z_\alpha) R_\alpha^+$

$$\therefore T^+ - T^- = \sum_{\alpha=1}^{n} Z_\alpha R_\alpha^+ - \sum_{\alpha=1}^{n} (1 - Z_\alpha) R_\alpha^+$$

$$= \sum_{\alpha=1}^{n} Z_\alpha R_\alpha^+ + \sum_{\alpha=1}^{n} Z_\alpha R_\alpha^+ - \sum_{\alpha=1}^{n} R_\alpha^+$$

$$= 2 \sum_{\alpha=1}^{n} Z_\alpha R_\alpha^+ - \frac{n(n+1)}{2}$$

● Difference b/w sign test and Wilcoxon rank sum test.

Sign test considers only the directions. while W-sign test

considers not only directions but also the magnitude of the obs$^n$'s.

- Under $H_0$ . $Z_1, Z_2, \ldots, Z_n$ are iid s.v with $P(Z_\alpha = 1) = \frac{1}{2}$

  because $P(Z_\alpha = 1) = Pr(X_\alpha > \mu_0) = \frac{1}{2}$

  ($x's$ are independent so $Z_\alpha's$ are also independent)

- $(Z_1, Z_2, \ldots, Z_n)$ are ind. of $(R_1^+, R_2^+, \ldots, R_n^+)$

<u>proof</u>

$$P\left(Z_\alpha = 1 \wedge |D_\alpha| \leq n\right)$$
$$\underset{\uparrow}{\text{arbitrary pt.}}$$

$$= P(0 < D_\alpha \leq n) = F_D(n) - F_D(0) = F_D(n) - \frac{1}{2} \quad [\text{under } H_0]$$
$$\underset{\uparrow}{\text{dist}^n \text{ fn of } D} = \frac{1}{2}\left[2F_D(n) - 1\right]$$

$$= Pr(Z_\alpha = 1) \cdot Pr(-n < D_\alpha < n)$$

$Z_\alpha's$ and $|D|'s$ are ind. $= Pr(Z_\alpha = 1) \cdot P(|D_\alpha| \leq n)$

Now the $R_\alpha^+$'s are the ranks of $|D_\alpha|$'s.

$\therefore R_\alpha^+$'s are the fn of $|D_\alpha|$

$\therefore Z_\alpha's$ and $R_\alpha^+$'s are ind.

- $$E(Z_\alpha) = \frac{1}{2} \quad , \quad V(Z_\alpha) = \frac{1}{4}$$

$$E(T^+) = E\left(\sum_{\alpha=1}^{n} Z_\alpha R_\alpha^+\right)$$

$$= \sum_{\alpha=1}^{n} E(Z_\alpha) R_\alpha^+$$

$$= \frac{1}{2} \sum_{\alpha=1}^{n} R_\alpha^+ = \frac{1}{2}(1 + 2 + \ldots + n) = \frac{n(n+1)}{4}$$

$$V(T^+) = \frac{n(n+1)(2n+1)}{24}$$

- <u>Determination of rejection region by $T^+$</u>

  To determine the rejection region, the prob. dist$^n$ of $T^+$ has to be determined under $H_0$.

  $$H_0: \mu = \mu_0 \equiv H_0: \underbrace{Pr(X_\alpha > \mu_0)}_{\pi} = 0.5$$

  $$H_1: \mu \leq \mu_0 \equiv \pi > 0.5$$

  $$H_1: \mu > \mu_0 \equiv \pi < 0.5$$

  $$\mu \neq \mu_0 \equiv \pi \neq 0.5$$

  $\therefore$ The extreme values of $T^+$ are zero and $\frac{n(n+1)}{2}$

  Since $T^+$ is completely determined by $Z_\alpha's$ the sample space can be considered to be the set of all possible $n$-tuples $\{Z_1, Z_2, \ldots, Z_n\}$ with components, either 0 or 1 forming $2^n$ possible possibilities (arrangements)

  Each of these distinguishable arrangements is equally

likely under H₀. Then the null distⁿ T⁺ is,

$$P(T^+ = t) = \frac{u(t)}{2^n}, \quad \text{where } u(t) \text{ is the number}$$

of ways to assign + and − sign on the n-integer $(1,2,...,n)$ values of such that the sum of the positive obsⁿ's is t.

Every assignment has a conjugate assignment. Interchanging + sign to − sign and − to +.

Ex. $n = 3$, $x_1, x_2, x_3$
R   1, 2, 3        $T^+ = \begin{matrix}0\\1\\2\\3\\4\\5\\6\end{matrix}$

| $T^+$ | Ranks associated to t | prob. value of u(t) |
|---|---|---|
| 0 | — | $P(T^+=0) = \frac{1}{8}$ |
| 1 | 1 (+1) | $P(T^+=1) = 1/8$ |
| 2 | 2 (+) | $P(T^+=2) = 1/8$ |
| 3 | 3 (+); 1,2 | $P(T^+=3) = 2/8$ |
| 4 | 1, 3 | $1/8$ |
| 5 | 2, 3 | $1/8$ |
| 6 | 1, 2, 3 | $1/8$ |

▪ Distⁿ of T⁺ is symmetric
▪ Conjugate pair always exists

Ex. $n = 4$, $x_1, x_2, x_3, x_4$
R   1   2   3   4

| $T^+$ | Ranks associated to + | Value of P(u(t)) |
|---|---|---|
| 10 | 1, 2, 3, 4 | $1/26$ |
| 9 | 2, 3, 4 | $1/16$ |
| 8 | 1, 3, 4 | $1/16$ |
| 7 | 4, 3 ; 1, 2, 4 | $2/16$ |
| 6 | 4, 2 ; 1, 2, 3 | $2/16$ |
| 5 | 2,3,1,4  2,3 ; 1,4 | $2/16$ |
| 4 | 1, 3 ; 4 | $2/16$ |
| 3 | 1, 2 ; 3 | $2/16$ |
| 2 | 2 | $1/16$ |
| 1 | 1 | $1/16$ |
| 0 | 0 | $1/16$ |

H.W. Q) Using any arbitrary kl (n=4) show that T⁺ is symmetric around its mean 5.

H.W. An educationᵃˡ testing service reports that the 75th percentile for scores of the GRE is 693 in a certain year. A r.s of 15 freshmen majoring in Stat. report

their GRE scores as 690, 750, 680, 700, 660, 710, 720, 730, 650, 670, 740, 730, 607, 750 and 690, are the scores of students majoring in Stat consistent with the $75^{th}$ percentile value.

$$\left( \begin{array}{l} H_0 = \pi = 3/4 \\ H_g : \pi \ne 3/4 \end{array} \right)$$

$$\therefore \pi = P_n(x > 693) = \frac{75}{100} = 3/4$$

1) Mean $= \frac{n(n+1)}{4}$

$$\therefore P\left(T^+ > \frac{n(n+1)}{4}\right) \qquad \left\langle T^+ + T^- = \frac{n(n+1)}{2} \quad \cdots(i) \right\rangle$$

$$= P\left(\frac{n(n+1)}{2} - T^+ < \frac{n(n+1)}{2} - \frac{n(n+1)}{4}\right)$$

$$= P\left(T^- < \frac{n(n+1)}{4}\right)$$

$$= P\left(T^- < \frac{T^+ + T^-}{2}\right)$$

$$= P(T^- < T^+)$$

$$= P\left(\frac{n(n+1)}{2} - T^+ < T^+\right)$$

$$P\left(T^+ > \frac{n(n+1)}{4}\right)$$

$$= P\left(T^+ > \frac{T^+ + T^-}{2}\right)$$

$$= P(T^+ > T^-)$$

$$P\left(T^+ > \frac{n(n+1)}{4}\right)$$

$$P\left(T^+ = \frac{n(n+1)}{4} + t\right)$$
$$= P\left(T^+ = \frac{T^+ + T^-}{2} + t\right)$$
$$= P\left(\frac{T^+}{2} = \frac{T^-}{2} + t\right)$$
$$= P\left(\frac{T^+}{2} = \frac{n(n+1)}{4} - \frac{T^+}{2} + t\right)$$
$$\Rightarrow P(T^+ = T^- + t)$$

$$P\left(T^+ = \frac{n(n+1)}{4} + t\right)$$
$$\Rightarrow P\left(T^+ = T^- + 2t\right)$$
$$\Rightarrow$$
$$\therefore T^+ - \frac{n(n+1)}{4} = T^+ - T^-$$
$$\Rightarrow T^+ - \frac{n(n+1)}{4} = \frac{-n(n+1)}{2}$$
$$\Rightarrow T^+ = \frac{n(n+1)}{4}$$
$$\Rightarrow T^+ - \frac{n(n+1)}{2}$$
$$= \frac{n(n+1)}{2} - T^+$$

$$P\left(\frac{n(n+1)}{4} - T^- = t\right)$$

$$P\left(T^+ > \frac{n(n+1)}{4}\right)$$
$$= P\left(2T^+ > \frac{n(n+1)}{2}\right)$$
$$= P\left(T^+ - \frac{n(n+1)}{4} > \frac{n(n+1)}{4} - T^+\right)$$
$$\therefore P\left(T^+ - \frac{n(n+1)}{4} > 0\right) = P\left(T^+ - \frac{n(n+1)}{4} > \frac{n(n+1)}{4} - T^+\right)$$

2) $H_0 : p = 3/4$ og. $H_1 : p \neq 3/4$

$$P(X_i < 693) = 3/4$$

$$S_{H_0} \sim \text{Bin}(15, \, 3/4)$$

| $X_i - 693$ | Sign |
|---|---|
| -3 | - |
| 57 | + |
| -13 | - |
| 7 | + |
| -33 | - |
| 17 | + |
| 27 | + |
| 37 | + |
| -43 | - |
| -23 | - |
| 47 | + |
| 37 | + |
| -33 | - |
| 57 | + |
| +3 | - |

Number of $X_i$'s greater than 693 is 8.

∴ Here $\langle S = 8 \rangle$

Let, $\langle \alpha = 0.1 \rangle$

We have to find $K_{\alpha/2}$ and $K'_{1-\alpha/2}$ such that it holds,

$$\sum_{s=0}^{K_{\alpha/2}} \binom{15}{s} \left(\frac{3}{4}\right)^s \left(\frac{1}{4}\right)^{n-s} \leq 0.05 \quad \dagger i)$$

and, $$\sum_{s=K'_{1-\alpha/2}}^{15} \binom{15}{s} \left(\frac{3}{4}\right)^s \left(\frac{1}{4}\right)^{n-s} \leq 0.05 \quad \cdots (ii)$$

respectively.

i) holds for, $K_{\alpha/2} = 7$ and (ii) holds for, $K'_{1-\alpha/2} = 14$

The test construction will be,

$$\phi(s) = \begin{cases} 1 & s \leq 7 \text{ and } s \geq 14 \\ \\ 0 & 0.\omega \end{cases}$$

∴ $H_0$ is accepted, i.e scores of students majoring in Stats consistent with the percentile value.

---

1) **Prove that $T^+$ is symmetric:**

In the construct. of $T^+$, every assignment has a conjugate assignment with plus and minus sign interchanged. Since we defined,

$$Z_\alpha = \begin{cases} 1 & x_\alpha > \mu_0 \\ 0 & x_\alpha < \mu_0 \end{cases}$$

Conjugate variable of $Z_\alpha$ will be, $(1 - Z_\alpha)$

∴ The value of $T^+$ for those conjugate assignments will be

$$T^+_{adj} = \sum_{\alpha=1}^{n} R^+_\alpha (1-z_\alpha) = \frac{n(n+1)}{2} - \sum_{\alpha=1}^{n} R^+_\alpha z_\alpha = \frac{n(n+1)}{2} - T^+_{org}.$$

Since every assignment occurs with equal prob. $\frac{1}{2^n}$, it implies that $T^+$ is symmetric around its mean $\frac{n(n+1)}{4}$

$$\Rightarrow T^+_{adj} - \frac{n(n+1)}{4} = \frac{n(n+1)}{4} - T^+_{org}.$$

**Result**: $T^+$ and $T^-$ are identically distributed.

$$\Rightarrow P[T^+ \geq c]$$

$$= P\left[T^+ - \frac{n(n+1)}{4} \geq c - \frac{n(n+1)}{4}\right]$$

$X$ is sym.
$P(X > \mu + c)$
$= P(X \leq \mu - c)$
$\Rightarrow P(x - \mu > c)$
$= P(x - \mu < -c)$
$= P(\mu - x > c)$

$$= P\left[\frac{n(n+1)}{4} - T^+ \geq c - \frac{n(n+1)}{4}\right]$$

$$= P\left[\frac{n(n+1)}{2} - T^+ \geq c\right] \quad [\because T^+ \text{ is symmetric}]$$

$$= P[T^- \geq c]$$

$T^+$ and $T^-$ follow the identical prob. dist$^n$.

**Remark**: Since it is more convenient to work with the smaller rank sum, so we use $T^+$ or $T^-$ accordingly. If $t_\alpha$ is the critical pt. such that (rejection region will be on the left hand side) $P(T \leq t_\alpha) = \alpha$, the rejection region for different alt. will be as follows,

| $H_1$ | Interpretation |
|---|---|
| $H_1: \mu > \mu_0$ | under $H_1$, $T^+$ will be higher, $T^-$ will be smaller $P(T^- \leq t_\alpha) = \alpha$ ↑ Critical point. we reject $H_0$ $T^- \leq t_\alpha$ |
| $H_1: \mu < \mu_0$ | $P(T^+ \leq t_\alpha) = \alpha$ we reject $H_0$ if $T^+ \leq t_\alpha$ |
| $H_1: \mu \neq \mu_0$ | $T^+ \leq t_{\alpha/2}$ or $T^- \leq t_{\alpha/2}$ |

**Result** $n = 3$

| $T^+$ | |
|---|---|
| 0 | 1/8 |
| 1 | 1/8 |
| 2 | 1/8 |
| 3 | 1/4 |
| 4 | 1/8 |
| 5 | 1/8 |
| 6 | 1/8 |

$H_0: \mu = \mu_0$
$H_1: \mu < \mu_0$
$P(T^+ \leq t_\alpha) = 0.05$

For every choice of $n$ and $\alpha$ the cut off pt. may not be found in Wilcoxon cont. signed rank test.

Therefore (i) choice of $\alpha$ is essential before constructing the test.

(ii) The critical pt. is not bound does not imply the test is invalid.

[ANS]

**mark:** For paired obs$^n$'s both sign test and Wilcoxon signed rank test can be applied by constructing the test on the differences, $Q_i^* = X_\alpha - Y_\alpha$ as the univariate obs$^n$.

**Problem 1** In a marketing research test 15 adult males ⓐ were asked to shave one side of their face with a brand A razor and the other side " " " " " " " B " and state their preferred razor. 12 male preferred preferred brand A find the p-value for the alt. that the prob. of preferring the brand A is greater than ~~0.05~~ 0.5.

$H_0$: A and B are equally preferable

$\equiv H_0 : \pi = \frac{1}{2}$

ag. $H_1$ : A is more preferable

$\equiv H_1 : \pi > \frac{1}{2}$

Let $S$ be the Sample statistic is no. of adults who preferred brand A.

$$S = 12$$

We use sign test where under $H_0$,

$$S \sim Bin(15, \frac{1}{2})$$

$$\text{p-value} = Pr[S \leq 12 / H_0]$$

$H_0$: A and B are equally preferable
That means if $\mu_0$ is the median of the pop$^n$. prefers brand A

$\Rightarrow F_{\mu_0}(x) \sim \frac{1}{2}$, under, $H_0$, $H_0 : \mu = \mu_0$

when brand preference of A is more than 50% median should be shifted to the right of $\mu_0$.

Then alt. hypo. is $H_1 : \mu > \mu_0$

We reject $H_0$ is $S - \frac{n}{2} > S_\alpha$

where, $P_{H_0}\left[S - n/2 \ge S_\alpha\right] \le 0.1$

$\Rightarrow \sum\limits_{S'_\alpha}^{n} \binom{n}{s}\left(\frac{1}{2}\right)^n \le 0.1$

**Problem (H.W)**

A study of 5 years ago reported that median amount of sleep by American adults is 7.5 hours out of 24 hours. A current sample of 8 adults reported their avg. amount of sleep per 24 hours as, 7.2, 8.3, 5.6, 7.4, 7.8, 5.2, 9.1 and 5.8 hours. Use the most appropriate stat test to determine whether American adults sleeps less today than they did 5 years ago.

$\Rightarrow$

| $D(x_\alpha - \mu_0)$ | $|D|$ | $R_i^+$ | |
|---|---|---|---|
| - 0.3 | 0.3 | 2.5 | $T^+ = 4 + 2.5 + 5$ |
| 0.8 | 0.8 | 4 | cal = 11.5 |
| - 1.9 | 1.9 | 7 | $T^- = 16.5$ |
| - 0.1 | 0.1 | 1 | |
| 0.3 | 0.3 | 2.5 | |
| - 2.3 | 2.3 | 8 | |
| 1.6 | 1.6 | 5 | |
| - 1.7 | 1.7 | 6 | |

we reject $H_0$ if,
$$\left\langle T^+ < T_{\frac{\alpha}{2}}\right\rangle$$
$\langle 11.5 \rangle$

Accept $H_0$ that american adults equally today as they did for 5 years ago.

**Application of Wilcoxon signed rank test in paired obs:**

Given the r.s of n-pairs for any specific individual $(x_1, Y_1), (x_2, Y_2), \ldots, (x_n, Y_n)$
Suppose $(x_i, Y_i)$; $i = 1(1)n$ come from a cont. dist" having the jt prob. dist" fun $F_{x,Y}(n, y) = P_r(x \le n \wedge Y \le y)$
Further suppose the differences b/w x and Y be $D = x - Y$ followed by $\mu_d$
In order to test, $H_0 : \mu_d = \mu_0$ where $\mu_d = median(x - Y)$ we have to form n differences $D_i = x_i - Y_i - \mu_0$. Remember by differencing x and Y, we convert a bivariate random variable to a univariate one. $\langle \mu_d : \mu_x - \mu_Y \rangle$
First are rank absolute value of $D_i$ and thereafter construct the Wilcoxon signed rank test statistic, based on the positive obs or negative obs.

$$T_\phi^+ = \sum_{\phi_i > 0} R_{\phi,i}^+ \quad \langle R_{\phi,i}^+ = \text{Rank of } |\phi_i| \rangle.$$

Rejection criterion remains the same as before.

**Practical**

1) A large company was distributed disturbed about the number of person-hours lost per month due to plant accidents and instituted an extensive industrial safety program. The data below show the number of person-hours lost in a month at each of 8 diff. plant before and after the safety program was implemented. Has the safety program being effective in reducing time lost from accidents.

| Plant: | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|--------|-----|------|------|------|------|------|------|------|
| (x) Before: | 51.2 | 46.5 | 24.1 | 10.2 | 66.3 | 92.1 | 30.3 | 49.2 |
| (y) After: | 45.8 | 41.3 | 15.8 | 11.1 | 58.5 | 70.3 | 31.6 | 35.4 |

⇒ Suppose person hour loss before and after safety program is denoted a bivariate random variable $(x, y)$
(Assume $(x, y)$ coming from a cont. dist$^n$ $F_{X,Y}(x,y)$.
Let us assume that $(\mu_x, \mu_y)$ is the bivariate median of $F_{X,Y}(x,y)$.

we are to test. $H_0: \mu_x = \mu_y$

$H_1: \mu_x > \mu_y$ --- (1)

Transform $\phi = x - y$. Assume $\mu_\phi$ be the median of the dist$^n$ of $\phi$.

(1) can be rewritten as;

$H_0: \mu_\phi = 0$
$H_1: \mu_\phi > 0$

Under alt. hypothesis rank of positive obs$^n$ will be higher, resulting $T^+$ larger and $T^-$ smaller simultaneously.

$\langle T^- = 3 \rangle$

$\langle T^+ = 33 \rangle$

∴ We reject $H_0$ if $T^- < t_\alpha$

where $T^-$ being the tabular value.

$t_\alpha \approx 2$ at $\alpha = 0.01$.

$3 > 2$, we fail to reject $H_0$.

| No. | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-----|-----|-----|-----|-----|-----|------|-----|------|
| \|x-y\| | 5.4 | 5.2 | 8.3 | 0.9 | 6.8 | 21.8 | 1.3 | 13.8 |
| $R^+$ | 4 | 3 | 6 | 1 | 5 | 8 | 2 | 7 |

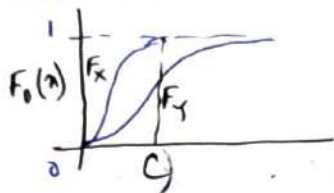Safety program is not effective in reducing time loss from accident.

2) Reducing High BP by diet requires reduc$^n$ of Na intake. Listed below are the avg. Na contents of 5 ordinary foods in processed from and natural for for equivalent quantities. do you see any difference b/w the median of processed and natural food.

| Natural food: | Corn of the cob (2) | Chicken (63) | Ground sirloin (60) | Beans (3) | Fresh tuna (40) |
|---|---|---|---|---|---|

| Processed food: | Canned Corn (251) | Fried Chicken (1220) | All beef biscoti (461) | Canned beans (300) | Canned tuna (409) |
|---|---|---|---|---|---|

Two pop$^n$ median test

Let $X_1, X_2, \ldots, X_{n_1}$ come from a cont. pop$^n$ $F_X(n)$.
and $Y_1, Y_2, \ldots, Y_{n_2}$ " " another ind. cont. pop$^n$ $F_Y(n)$.
The r.v. $Y$ is called stochastically larger to $x$ if $Y$ takes same prob prob. for higher values while $X$ takes that prob. for lower values.



$$\left\langle \begin{array}{c} Y >_{st} X \\ \Rightarrow F_Y(c) < F_X(c) \end{array} \right\rangle$$

**Remark** Two pop$^n$ non-parametric location test is based on the idea of equality of two medians ($\mu_X$ and $\mu_Y$)

$$\boxed{Y \gtrsim_{st} X \equiv \mu_X > \mu_Y}$$

We are to test,
$H_0: \mu_X = \mu_Y$  $H_1$: $\mu_X < \mu_Y \equiv Y \gtrsim_{st} X$
$\mu_X > \mu_Y \equiv Y \lesssim_{st} X$
$\mu_X \neq \mu_Y$

## Mann-Whiteny Test

M-W U test is a special choice of testing the above, where it is assumed that, two pop$^n$'s are differed by a loca$^n$ parameter $\theta$, i.e.

$$\langle F_X(n) = F_Y(n+\theta) \rangle$$

$H_0: \mu_X = \mu_Y$
eg. $H_1: \mu_Y > \mu_X$      is analogous of writing
$Y \gtrsim_{st} X$   $\left\langle \begin{array}{c} H_0: \theta = 0 \\ H_1: \theta > 0 \end{array} \right\rangle$

For testing the above, we check how many of $Y$ sample obs$^n$'s are _less than_ $X$ obs$^n$'s in the combined sample.

Define, $\partial_{ij} = \begin{cases} 1 & \text{if } Y_j < X_i \begin{cases} i=1(1)n_1 \\ j=1(1)n_2 \end{cases} \\ 0 & \text{o.w} \end{cases}$

$\therefore$ U-statistics is $\left\langle U = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \partial_{ij} \right\rangle$

$\times$ Number of times $Y$ precedes $x$ $\rangle$

$\langle$ Clearly small value of U reject $H_0 \rangle$
~~Therefore~~ Therefore the test based on U will be a left tail test.

## E(U) and V(U)

Assume $P[Y < X] = \pi$ (say)
$E(\partial_{ij}) = \pi$
$E(U) = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \pi = n_1 n_2 \pi$

$\partial_{ij}$'s are ind. variable with $\partial_{jk}$ as $P[Y_j < X_i \wedge Y_k < X_l]$ is just the product of $P[Y_j < X_i] \cdot P[Y_k < X_l]$

But $D_{ij}$'s are not independent for common subscript.

Let,
$$P\left[Y_j \leq x_i \text{ and } Y_k < x_i\right]$$
$$= \int_{-\infty}^{\infty} [F_Y(m)]^2 \, dF_X(m) = \Pi_1 \text{ (say)}$$

and,
$$P\left[Y_j < x_i \text{ and } Y_j < x_k\right]$$
$$= P\left[x_i \text{ and } x_k > Y_j\right]$$
$$= \int_{-\infty}^{\infty} (1 - F_X(m))^2 \, dF_Y(m) = \Pi_2 \text{ (say)}$$

$\therefore \left[Cov(D_{ij}, D_{kj}) = E(D_{ij} \cdot D_{ik}) - E(D_{ij}) \cdot E(D_{ik})\right.$
$$= \Pi_1 - \Pi^2$$

$\therefore \langle Cov(D_{ij}, D_{kj}) = \Pi_2 - \Pi^2 \rangle \langle Cov(D_{ij}, D_{lk}) = 0 \rangle$

$$V(U) = V\left(\sum_i^{n_1} \sum_j^{n_2} D_{ij}\right)$$

$$= \sum_i \sum_j V(D_{ij}) + \sum_i^{n_1} \sum_j^{n_2} \sum_l^{n_1} \sum_k^{n_2} Cov(D_{ij}, D_{lk})$$
$$\quad\quad i \neq l, j \neq k$$

$$+ \sum_i^{n_1} \sum_j^{n_2} \sum_k^{n_2} Cov(D_{ij}, D_{ik})$$
$$\quad i=1 \; 1 \leq j \neq k \leq n_2$$

$$+ \sum_{l \leq i \neq k \leq n_1} \sum \sum_{j=1}^{n_2} Cov(D_{ij}, D_{kj})$$

$$= n_1 n_2 \Pi(1-\Pi) + 0 + n_1(n_2-1)_2(\Pi_1 - \Pi^2) + n_2(n_1-1)(\Pi_2 - \Pi^2) n_1$$

$$= n_1 n_2 \left[\Pi(1-\Pi) + (n_2-1)(\Pi_1 - \Pi^2) + (n_1-1)(\Pi_2 - \Pi^2)\right]$$

$$= n_1 n_2 \left[\Pi - \Pi^2 + (n_2-1)\Pi^2 + (n_2-1)\Pi_1 + (n_1-1)\Pi_2 - (n_1-1)\Pi^2\right]$$

$$= n_1 n_2 \left[\Pi - \Pi^2(1 + (n_2-1) + (n_1-1)) + (n_2-1)\Pi_1 + (n_1-1)\Pi_2\right]$$

$$= n_1 n_2 \left[\Pi - (N-1)\Pi^2 + (n_2-1)\Pi_1 + (n_1-1)\Pi_2\right]$$

$\langle V\left(\frac{U}{n_1 n_2}\right) = 0 \text{ as } n_1 \to \infty, n_2 \to \infty \rangle$  $\boxed{N = n_1 + n_2}$

$\langle E\left(\frac{U}{n_1 n_2}\right) = \Pi \rangle$ $\therefore \langle \frac{U}{n_1 n_2} \text{ is a consistent estimator} \rangle$
$$P[Y < x] = \Pi$$

**Discrete dist$^n$ of U** For $n_1$ X obs and $n_2$ Y obs. There are $\binom{n_1+n_2}{n_1}$ arrangements bs X and Y in combined sample. For every particular arrangement Z, $\exists$ one conjugate arrangement as if Z denotes a set of X and Y written from smallest to largest. Then its conjugate

arrangement z' may be proposed from largest to smallest (conjugate arrangement: how many X follow Y). If U be an arrangement then the prob. dist$^n$ of its conjugate arrangement value will be the same, and that value is $\left\langle U' = \sum_i \sum_j (1 - \partial_{ij}) \right\rangle$

Ex: $n_1 = 4, n_2 = 5$

$$\binom{4+5}{5} = 126$$

| arrangement | | $P(U=u)$ |
|---|---|---|
| X X Y X X Y Y Y Y | $\begin{matrix} 4 \\ 0 \end{matrix}$ | $1/126$ |
| X X Y X X Y Y Y Y | $\left.\begin{matrix} 2 \\ 2 \end{matrix}\right\}$ | $2/126$ |
| X X X Y Y X X Y Y | | |
| X X X Y X Y Y Y Y | $1$ | $1/126$ |

∴ The pmf of U is $\left\langle P[U=u] = \dfrac{r_u}{\binom{n_1+n_2}{n}} \right\rangle$ distinguishable

where, $r_u$ is the number of arrangements, for which r.v U takes the value u.

**W** $\langle$ Find out $E(U)$ and $V(U)$ under $H_0\rangle$

Remark i) For alt. hypothesis, $H_1 : Y \underset{st}{\geqslant} X \equiv H_1 : \mu_Y > \mu_X$

we reject $H_0$ if $\boxed{U < u_\alpha}$

where $u_\alpha$ be the tabular value at $\alpha$ level of significance.

For, $H_1 : Y \underset{st}{\leqslant} X \equiv H_1 : \mu_Y < \mu_X$

we reject $H_0$ if $\boxed{U' < u_\alpha}$ $\left[ U' = \sum_i \sum_j (1 - \partial_{ij}) \right]$

For, $H_1 : \mu_X \neq \mu_Y$

we reject $H_0$ if $\boxed{U < u_{\alpha/2}}$ on $\boxed{U' < u_{\alpha/2}}$

Remark ii) For tied case,

$$\partial_{ij} = \begin{cases} 1 & \text{if } Y_j < X_i \\ 0.5 & Y_j = X_i \\ 0 & Y_j > X_i \end{cases}$$

**Practical** The 2000 census statistics for Alabamma, if the % changes in pop$^n$ b/w 1990 and 2000 for each of the district 67 counties. There are 2 types of counties — rural and non rural acc$^n$ to the pop$^n$ size < 25000. Below is the data of 9 rural and 7 non-rural counties on % of pop$^n$ change

Row Rural : 1.1, -21.7, -16.3, -11.3, -10.4, -7.0, -2.0, 1.9, 6.2
Non " : -2.4, 9.9, 14.2, 18.4, 20.1, 23.1, 70.4 the null hypo.
Use Mann-Whitney U test for testing the equal pop$^n$ change.

⇒ Let the pop$^n$ change of rural county comes from a cont. dist with c.d.f $F_Y$ median $\mu_Y$ & th cmtrst and the pop$^n$ change of non " " " " " " " " $F_X$ and median $\mu_X$. We are to test $H_0 : \mu_X = \mu_Y$ — $H_1 : \mu_X \neq \mu_Y$

Arrange the combined sample,

Y Y Y   YY X Y Y Y Y X X X X X X

$$U = 59, \quad U' = 4$$

$\alpha = 0.02 \quad \langle U_{tabular} = 9 \rangle$

$U' = 4 < U_{tab} = 9$

$\therefore$ We reject the null hyp i.e $H_0$

**ol**

Under $H_0$ we know that, $\pi = \frac{1}{2}$

$$\therefore E(U) = \frac{n_1 n_2}{2}$$

Under $H_0$,
$$\pi_1 = \int_{-\infty}^{\infty} [F_x(n)]^2 \, dF_x(n)$$

$$= \frac{[F_x(n)]^3}{3} \bigg|_{-\infty}^{\infty}$$

$$= \frac{1}{3}$$

$$\pi_2 = \int_{-\infty}^{\infty} [1 - F_x(n)]^2 \, dF_x(n)$$

$$= -\frac{[1 - F_x(n)]^3}{3} \bigg|_{-\infty}^{\infty}$$

$$= \frac{1}{3}$$

$\therefore$ Putting the value of $\pi_1$ and $\pi_2$ in the expression of $V(U)$ we will get that,

$$V(U) = n_1 n_2 \left[ \frac{1}{2} - (N-1)\frac{1}{4} + \frac{(n_2-1)}{3} + \frac{(n_2-1)}{3} \right]$$

$$= n_1 n_2 \left[ \frac{6 - 3N + 3 + 4n_2 - 4 + 4n_1 - 4}{12} \right]$$

$$= n_1 n_2 \left[ \frac{n_1 + n_2 + 1}{12} \right]$$