$$P\left[\ t_{n-1} > c''''\ \right] = \alpha.$$

$X \sim Bin(n, p)$

$\therefore P(X = x) = {}^nC_x\, p^x (1-p)^{n-x}, \quad x = 0, 1, 2, \cdots$

$P(X = odd)$

$= P(X = 1) + P(X = 3) + \cdots$

$= {}^nC_1\, p(1-p)^{n-1} + {}^nC_3\, p^3 (1-p)^{n-3} + \cdots$

$= p(1-p)^{n-1}\left[\ {}^nC_1 + {}^nC_3(1-p)^2 + {}^nC_5(1-p)^4 + \cdots\ \right]$

---

$X \sim Bin(n, p), \quad P(X = x) = {}^nC_x\, p^x (1-p)^{n-x}$

Let $q = (1-p)$

$(q + p)^n = {}^nC_0\, q^n p^0 + {}^nC_1\, q^{n-1} p^1$
$\qquad\qquad + \cdots + {}^nC_n\, q^0 p^n$

$= P(X = 0) + P(X = 1) + \cdots + P(X = n)$
$= \boxed{P(X = even) + P(X = odd)}$

Similarly,

$(q - p)^n = {}^nC_0\, q^n p^0 - {}^nC_1\, p^1 q^{n-1} + {}^nC_2\, p^2 q^{n-2}$
$\qquad\qquad \cdots$

$= \boxed{P(X = even) - P(X = odd)}$

$\therefore (q + p)^n - (q - p)^n = 2 \times P(X = odd)$

$\Rightarrow P(X = odd) = \dfrac{1 - (q - p)^n}{2} \quad \left[\because p + q = 1\right]$

---

Non-parametric

1. Gibbons (book).

Rank

Suppose $X_\alpha$ be the $\alpha^{th}$ observation for a set of $n$ obs. $\alpha = 1, 2, \cdots, n$, all are i.i.d. from a continuous distribution $F_X(x)$.

$\quad$ Rank of $X_\alpha$ : # of observations $\leq X_\alpha$

rank is an ordered permutation vector

$R = (R_1, R_2, \cdots, R_n)$

$\quad\downarrow\qquad\downarrow\qquad\qquad\downarrow$

Rank of $\quad$ Rank of $\qquad$ Rank of
$1^{st}$ obs $\quad$ $2^{nd}$ obs $\qquad$ $n^{th}$ obs

$P(\underline{R} = \underline{r}) = \dfrac{1}{n!}, \quad P(R_\alpha = r_\alpha) = \dfrac{1}{n}, \quad n = 1, 2, \cdots$

(Joint) $\qquad\qquad$ (Marginal)

Marginal distribution of rank is discrete uniform.

$P(R_\alpha = r_\alpha \cap R_\beta = r_\beta) = \dfrac{1}{n(n-1)} = \dfrac{1}{{}^nP_2}$

$E(R_\alpha) = \dfrac{n+1}{2}$

$\quad Var(R_\alpha) = \dfrac{n^2 - 1}{12}$

$Cov(R_\alpha, R_\beta) = \dfrac{-(n+1)}{2}$

$\therefore R_\alpha$ and $R_\beta$ are not independent

$corr(R_\alpha, R_\beta) = \dfrac{-\frac{(n+1)}{2}}{\frac{(n+1)(n-1)}{12}} = \dfrac{-6}{(n-1)}$

## Linear Rank Statistic

Let $\underset{\sim}{a} = (a_1, a_2, \cdots, a_n)$

$\underset{\sim}{b} = (b_1, b_2, \cdots, b_n)$

be two sets of coefficients based on $n$ natural number

Let $\underset{\sim}{R} = (R_1, R_2, \cdots, R_n)$ be the random permutation of $\{1, 2, \cdots, n\}$. The linear Rank Statistic is

$$T = \sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha}$$

$$= a_1 b_{R_1} + a_2 b_{R_2} + \cdots + a_n b_{R_n}$$

$a_\alpha$'s are known as <u>regression constants</u> and $b_{R_\alpha}$'s are <u>scores</u> / <u>coefficients</u>

✳ Joint distribution of ranks is independent of the distribution function from which observations are coming.

Therefore the distribution of $T$ is independent of any $F(\cdot)$.

Hence, $T$ can be used to provide distribution free test (non-parametric)

Also, $(R_1, R_2, \cdots, R_N) \overset{D}{=} (n-R_1+1, n-R_2+1, \cdots, n-R_n+1)$

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

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

$$\Rightarrow E(b_{R_\alpha}) = \frac{1}{n} \cdot b_1 + \frac{1}{n} \cdot b_2 + \cdots + \frac{1}{n} b_n$$

$$= \frac{1}{n} \sum_{i=1}^{n} b_i = \bar{b}$$

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

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

$$= \sum_{\alpha=1}^{n} a_\alpha^2 Var(b_{R_\alpha}) + \sum\sum_{\alpha \neq \beta} a_\alpha a_\beta cov(b_{R_\alpha}, b_{R_\beta})$$

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

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

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

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

$$Var(b_{R_\alpha}) = Var(b_{R_1}) \quad [\text{w.l.g fix } \alpha = 1]$$

$$Var(b_{R_1}) = E(b_{R_1}^2) - [E(b_{R_1})]^2$$

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

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

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

$$= E\left[(b_{R_\alpha} - \bar{b})(b_{R_\beta} - \bar{b})\right]$$

$$= \sum\sum_{R_\alpha \neq R_\beta}(b_{R_\alpha} - \bar{b})(b_{R_\beta} - \bar{b})\, P(R_\alpha = \alpha, R_\beta = \beta)$$

$$= \frac{1}{n(n-1)}\sum\sum_{R_\alpha \neq R_\beta}(b_\alpha - \bar{b})(b_\beta - \bar{b})$$

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

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

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

$\downarrow$ mean deviation about mean

$$= -\frac{Var(b_{R_\alpha})}{(n-1)}$$

$$\therefore Var(T)$$

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

↗ irrespective of index, same function

$$\otimes - \frac{Var(b_{R_\alpha})}{(n-1)}\sum\sum_{\alpha \neq \beta}(a_\alpha - \bar{a})(a_\beta - \bar{a})$$

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

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

$$= Var(b_{R_\alpha})\cdot\left(\frac{n}{n-1}\right)\sum_{\alpha=1}^{n}(a_\alpha - \bar{a})^2$$

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

<u>Asymptotic distribution of linear rank Statistic</u>

✓ Wald Wolforitz condition

✓ Noether condition

✓ Hoeffding condition

$$\frac{T - E(T)}{\sqrt{V(T)}} \xrightarrow{L} N(0,1).$$

$\boxed{\text{Inverse permutation}}$

Every permutation $(P)$ has its mirror image $(P')$ Such that, $\boxed{P \circ P^{-1} = I^+ = \{1, 2, \cdots, n\}}$
$\hookrightarrow$ operation

Inverse permutation in Ranking theory is called anti-rank $\underset{\sim}{D} = (D_1\ D_2\ \text{---}\ D_n)$

$$\underset{\sim}{R}\ o\ \underset{\sim}{D} = I$$

Number and Number of the place the obs. occupies is exchanged.

Ex:

$$R = (3 \quad 1 \quad 4 \quad 2)$$

| No. 3 | No. 1 | No. 4 | No. 2 |
|---|---|---|---|
| Pl. 1 | Pl. 2 | Pl. 3 | Pl. 4 |

| 1 | 2 | 3 | 4 |
| 3 | 1 | 4 | 2 |

$$\therefore D = (2 \quad 4 \quad 1 \quad 3) ✓$$

☆ $P = (2 \quad 3 \quad 4 \quad 5 \quad 1)$

$$D = (5 \quad 1 \quad 2 \quad 3 \quad 4) ✓$$

We know ranks are not independent.
But the values in inverse permutations are independent.
Moreover joint distribution of any permutation remains the same.
Therefore joint distribution of ranking and anti-ranking are same $= \dfrac{1}{n!}$

Hence for finding distribution of $T$, anti-ranking process is more convenient.

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

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

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

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

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

$$= \frac{1}{n(n-1)}\left[\frac{3n^2(n+1)^2 - 2n(n+1)(2n+1)}{12}\right]$$

$$= \frac{n(n+1)}{n(n-1)}\left[\frac{3n(n+1) - 2(2n+1)}{12}\right] - D$$

$$= \frac{n+1}{n-1}\left[\frac{3n^2 + 3n - 4n - 2}{12}\right] - D$$

$$= \frac{n+1}{n-1}\left[\frac{3n^2 - 3n + 2n - 2}{12}\right]$$

$$= \frac{n+1}{n-1}\left[\frac{(3n+2)(n-1)}{12}\right] - D$$

$\therefore Corr(R_\alpha, R_\beta)$

$= \frac{(n+1)(3n+2)}{12} - \frac{3(n+1)^2}{12}$

$\because Corr(R_\alpha, R_\beta) = -\frac{1}{n-1}$

$$= \frac{n+1}{12}\left[(3n+2) - 3n - 3\right]$$

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

$$T = \sum_{\alpha=1}^{3} a_{\alpha} b_{R_\alpha}$$

$\underset{\sim}{D} = $ anti-ranks

$$T' = \sum_{\alpha=1}^{3} a_{D_\alpha} b_\alpha$$

⊗ 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 distribution of $T$ is symmetric about its mean.

Proof → We know that $(R_1, R_2, \cdots, R_n)$
$$\overset{D}{\equiv} (n-R_1+1, n-R_2+1, \cdots, n-R_n+1)$$

Assume $b_{R_\alpha} + b_{R_{(n-\alpha+1)}} \equiv k$ $\forall \alpha$ where $k$ is a constant.

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

$$\therefore T \overset{D}{\equiv} \sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha} \overset{D}{\equiv} \sum_{\alpha=1}^{n} a_\alpha b_{\cancel{(n-R_\alpha+1)}}$$

$$\overset{D}{\equiv} \sum_{\alpha=1}^{n} a_{D_\alpha} b_{R_{(n-\alpha+1)}}$$

$$\overset{D}{\equiv} \sum_{\alpha=1}^{n} a_{D_\alpha} (k - b_\alpha)$$

$$\overset{D}{\equiv} \sum_{\alpha=1}^{n} a_{D_\alpha} k - \sum_{\alpha=1}^{n} a_{D_\alpha} b_\alpha$$

$$\overset{D}{\equiv} k \, n\bar{a} - \sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha}$$

$$\therefore T \overset{D}{\equiv} (k \, n \bar{a} - T), \text{ it is true for all } k.$$

Choose $k = \bar{b}$

$$\therefore T \overset{D}{\equiv} 2n \bar{a} \bar{b} - T, \text{ Hence the proof}.$$

## Sign test

Let $x_1, x_2, \cdots, x_n$ be a random sample from a continuous distribution $F_X(x)$

Let $\mu$ be the median of the distribution

We are to test
$$H_0: \mu = \mu_0 \text{ (a fixed constant)}$$
$$H_1: \mu > \mu_0$$

Let $S = \#$ of observations greater than $\mu_0$

Let us propose a linear rank statistic for testing the above

Define,
$$a_\alpha = \begin{cases} 1, & \text{if } X_\alpha > \mu_0 \\ 0 & \text{o.w} \end{cases}$$

and $b_{R_\alpha} = 1$ $\forall \alpha = 1(1)n$

$$\therefore T = \sum_{\alpha=1}^{n} a_\alpha b_{R_\alpha} = \sum_{\alpha=1}^{n} a_\alpha = S \text{ (equivalent)}$$

$$\Rightarrow E(T) = \sum_{\alpha=1}^{n} E(a_\alpha)$$

$\boxed{\text{Under } H_0}$

$$= \sum_{\alpha=1}^{n} 1 \cdot P(X_\alpha > \mu_0) = P(X_1 > \mu_0) + P(X_2 > \mu_0)$$
$$+ \cdots + P(X_n > \mu_0)$$
$$= \boxed{n/2} \quad (\tfrac{1}{2}, n \text{ times})$$

Under $H_0$

$$\text{Var}(T) = \text{Var}_{H_0}\left(\sum_{\alpha=1}^{n} a_\alpha\right)$$

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

$$= \sum_{\alpha=1}^{m} E(a_\alpha^2) - \left[E(a_\alpha)\right]^2 + \sum_{\alpha\neq\beta}\sum \text{cov}(a_\alpha, a_\beta)$$

$$= \sum_{\alpha=1}^{n} E(a_\alpha) - \left(E(a_\alpha)\right)^2 + \sum_{\alpha\neq\beta}\sum \text{cov}(a_\alpha, a_\beta)$$

$\downarrow$ Binary

$$= \frac{n}{4} + \sum_{\alpha\neq\beta}\sum E(a_\alpha a_\beta) - E(a_\alpha) E(a_\beta)$$

$$= \frac{n}{4} + \sum_{\alpha\neq\beta}^{n}\sum^{n} P(X_\alpha > M_0 \cap X_\beta > M_0) - \left(\frac{1}{2}\cdot\frac{1}{2}\right)$$

$$= \frac{n}{4} + \sum_{\alpha\neq\beta}^{n}\sum^{n} P(X_\alpha > M_0) \cdot P(X_\beta > M_0) - \frac{1}{4}$$

Since the sample is random

$$= \frac{n}{4} + \sum_{\alpha\neq\beta}^{n}\sum^{n} \left(\frac{1}{4} - \frac{1}{4}\right) = \boxed{\frac{n}{4}}$$

Therefore $S$ is a particular case of linear rank statistic.

Also; it is a symmetric linear rank statistic around mean $\frac{n}{2}$.

---

Critical function of sign test

Let us propose a test function

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

$c$ and $\gamma$ are determined from the size condition.

Remark

For sign test, if zero difference occurs, for any $X_\alpha = M_0$, for continuous distribution assumptions, this does not create any problem as $Pr(X_\alpha = M_0) = 0$.

But in practical, zero difference can be avoided by ignoring and dropping them, simultaneously. Hence, $n$ is minimized.

Result
(sign test is UMP)

Propose a test function

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

For any c.d.f $F_X(x)$, with p.d.f $f_X(x)$

define, $f(x) = F(0) f^-(x) + (1 - F(0)) f^+(x)$ $\quad --- (1)$

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

against $H_1: \mu > \mu_0$.

$F_X(\mu_0) = \frac{1}{2}$

Let us make a transformation $X_\alpha - \mu_0$
such that $\mu_0 = 0$ and the corresponding
test will be

$H_0: \mu' = 0 \Longrightarrow \underline{F_X(0) = \frac{1}{2}}$

$H_1: \mu' > 0 \quad (\text{where } F_X(0) = \frac{1}{2})$

$F_X(0) < \frac{1}{2}$

In (1) $\quad f^- = \frac{f(x)}{F(0)}$ if $X \leq 0$

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

Under $H_0$, $\boxed{f_0(x) = \frac{1}{2} f_0^-(x) + \frac{1}{2} f_0^+(x)}$

Under $H_1$, $f_1(x) = F_1(0) f_1^-(x) + (1 - F_1(0)) f_1^+(x)$

Let us reframe the test function as
per $f^+, f^-$.

$\phi(x) = \begin{cases} 1 & \text{if } \prod\limits_{a=1}^{n} f_1(x_i) > k \prod\limits_{a=1}^{n} f_0(x_i) \\ \gamma & \prod\limits_{a=1}^{n} f_1(x_i) = k \prod\limits_{a=1}^{n} f_0(x_i) \\ 0 & \text{o.w} \end{cases}$

$\Downarrow$

$\phi(x) = \begin{cases} 1 & \prod\limits_{\alpha \neq (\alpha, \alpha_2 \dots \alpha_s)} f_1^-(x_\alpha) \prod\limits_{\alpha = (\alpha, \alpha_2 \dots, \alpha_s)} f_1^+(x_\alpha) > k \prod\limits_{\alpha \neq (\alpha_1 \dots \alpha_s)} f_0^-(x_i) \prod\limits_{\alpha = (\alpha_1 \dots \alpha_s)} f_0^+(x_i) \\ \gamma & = \\ 0 & < \end{cases}$

$\Rightarrow \phi(x) = \begin{cases} 1 & \text{if } \prod \dfrac{f_1(x_\alpha)}{F_1(0)} \prod\limits_{\alpha(\alpha_1 - \alpha_s)} \dfrac{f_1(x_\alpha)}{1 - F_1(0)} \\ \gamma & > k \prod\limits_{\alpha} \dfrac{f_0(x_\alpha)}{1/2} \prod\limits_{\alpha} \dfrac{f_0(x_\alpha)}{1/2} \\ & \quad '' \quad = \quad '' \\ 0 & \quad '' \quad < \quad '' \end{cases}$

$\Rightarrow \phi(x) = \begin{cases} 1 & \dfrac{\prod\limits_{\alpha \neq (\alpha_1 - \alpha_s)} f_1(x_\alpha) \prod\limits_{\alpha = (\alpha_1 - \alpha_s)} f_1(x_\alpha)}{(F_1(0))^{n-s} (1 - F_1(0))^{s}} > 2^n \cdot k \prod\limits_{\alpha \neq (\alpha_1 - \alpha_s)} f_0(x_\alpha) \prod\limits_{\alpha = (\alpha_1 - \alpha_s)} f_0(x_\alpha) \\ \gamma & \quad '' \quad = \quad '' \\ 0 & \quad '' \quad < \quad '' \end{cases}$

$\Rightarrow \phi(x) = \begin{cases} 1 & \left(\dfrac{F_1(0)}{1 - F_1(0)}\right)^{s} > k^* \prod\limits_{\alpha=1}^{n} \dfrac{f_0(x_\alpha)}{f_1(x_\alpha)} \\ \gamma & \quad '' \quad = \quad '' \\ 0 & \quad '' \quad < \quad '' \end{cases}$

$\boxed{k^* = \dfrac{2^n k}{(F_1(0))^n}}$

$\therefore$ For fixed $x_1, x_2, \ldots, x_n$, $\prod \dfrac{f_0}{f_1}$ is a constant

$$\phi(s) = \begin{cases} 1 & \left(\dfrac{F_1(0)}{1-F_1(0)}\right)^s > K^{**} \text{ monotonically} \\ \gamma & " = " \text{ increasing } f^n \\ & \qquad \text{ of } s \\ 0 & " < " \end{cases}$$

$\therefore \phi(x)$ can be written, in terms of $S$, # of positive observations.

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

As the above satisfies N-P test construction, it is a UMP test.

## Homework (Practical)

(Problems on Non-parametric inferences)

⊕ Suppose that each of 13 randomly chosen female registered voters was asked to indicate if she was going to vote for candidate A or B in an upcoming election. The result shows that 9 of the subjects preferred A. Is this sufficient evidence to conclude that candidate A is preferred to B by female voters?

---

Draw the power curve taking at least 8 points.

---

## Solution

We have a population of 13 female registered voters.

In an upcoming election, they have ~~an election~~ option to vote for candidate A or B.

Let us consider $p$ as the probability to vote for candidate 'A'.

We test,

$H_0: p = \dfrac{1}{2}$  vs  $H_1: p > \dfrac{1}{2}$

We have

test statistic $S = 9$

As test is constructed as

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

under $H_0$, $S_{H_0} \sim Bin\left(13, \dfrac{1}{2}\right)$

$\gamma$ and $k_\alpha$ are to be determined from size condition,

$E_{H_0}(\phi(s)) = 0.05 \; (\alpha)$

$\Rightarrow P_{H_0}(s > k_\alpha) + \gamma\, P_{H_0}(S = k_\alpha) = 0.05$

$\Rightarrow \gamma = \dfrac{0.05 - P(S > k_\alpha)}{P(S = k_\alpha)}$

$= \dfrac{P_{H_0}(S \leq k_\alpha) - 0.95}{P_{H_0}(S = k_\alpha)}$

We construct the table as follows

Under H₀, $S \sim Bin(13, \frac{1}{2})$

| $k_\alpha$ | $P(S=k_\alpha)$ | $P(S \le k_\alpha)$ |
|---|---|---|
| 0 | 0.00012 | 0.00012 |
| 1 | 0.00159 | 0.00171 |
| 2 | 0.00952 | 0.01123 |
| 3 | 0.03491 | 0.04614 |
| 4 | 0.08728 | 0.13342 |
| 5 | 0.1571 | 0.29052 |
| 6 | 0.20947 | 0.49999 |
| 7 | 0.20947 | 0.70946 |
| 8 | 0.1571 | 0.86656 |
| 9 | 0.08728 | 0.95384 → The point to be randomised |
| 10 | 0.03491 | 0.98875 |

We take $k_\alpha = 9$ for randomization
and we see

$$\gamma = \frac{P(S \le k_\alpha) - 0.95}{P(S=k_\alpha)}$$

$$= \frac{P(S \le 9) - 0.95}{P(S=9)}$$

$$= \frac{0.95384 - 0.95}{0.08728}$$

$$= 0.044$$

∴ The test is constructed as,-

$$\phi(s) = \begin{cases} 1 & s > 9 \\ 0.044 & s = 9 \\ 0 & s < \end{cases}$$

Since, $S = 9$, we reject the null hypothesis with rejection probability $\gamma = 0.044$.

Hence, the evidence is not so conclusive.

## Power function

⊛ Why sign test is a non-parametric test?

Sign test is considered as a non-parametric test because

⟶ (H.W) (1) No pre-assumption of specific distribution
(2) Handles ordinal data and deals with median

## Wilcoxon sign rank test

Let $X_1, X_2, \ldots, X_n$ be the random sample from a continuous c.d.f $F(\cdot)$ with median $\mu$

We are to test

$$H_0 : \mu = \mu_0$$

Consider the difference $D_\alpha = X_\alpha - \mu_0$
Clearly, the differences are distributed symmetrically under $H_0$.

$$F_D(-c) = P(D_\alpha \le -c) = P(D_\alpha > c)$$
$$= 1 - F_D(c)$$

With the assumption of a continuous population, zero or tied difference can be avoided by dropping them.

Next we order absolute $D_\alpha$'s

i.e, $|D_\alpha|$'s increasingly.

$$|D_1|, |D_2|, \ldots, |D_n|.$$

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

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

② $T^+ + T^- = \dfrac{n(n+1)}{2}$

---

From the linear model, solving the normal equation we get

~~$H = (X^T X)^{-1} X^T \theta X$~~

~~$H^T =$~~

---

$T^+$ and $T^-$ 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 $T^+ = \sum\limits_{\alpha=1}^{n} Z_\alpha R_\alpha^+$

where $Z_\alpha = \begin{cases} 1 & \text{if } D_\alpha > 0 \equiv X_\alpha > M_0 \\ 0 & \text{o.w} \end{cases}$ (check if L. Rank statn or not)

Similarly, $T^- = \sum\limits_{\alpha=1}^{n} (1 - Z_\alpha) R_\alpha^+$

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

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

$= 2 \sum\limits_{\alpha=1}^{n} Z_\alpha R_\alpha^+ - \dfrac{n(n+1)}{2}$

$\boxed{\text{Difference}}$

Sign test considers only the directions while W-sign rank test considers not only directions but also the magnitude of observation

Under $H_0$, $Z_1, Z_2, \ldots, Z_n$ are iid random variables with $P(Z_\alpha = 1) = \frac{1}{2}$

because $P(Z_\alpha = 1) = P(X_\alpha > M_0) = \frac{1}{2}$

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

**Proof:** $P(Z_\alpha = 1 \cap |D_\alpha| \leq x)$

$= P(D_\alpha > 0 \cap -x \leq D_\alpha \leq x)$

$= P(0 < D_\alpha \leq x)$

$= F_D(x) - F_D(0)$

$= F_D(x) - \frac{1}{2}$  $\left[\text{under } H_0 \atop P(D_\alpha \leq 0) = \frac{1}{2}\right]$

$= \frac{1}{2}[2F_D(x) - 1]$

$= P(Z_\alpha = 1) \cdot P(-x < D_\alpha < x)$

$Z_\alpha$'s and $|D|$'s are independent

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

$= P(Z_\alpha = 1) P(|D_\alpha| \leq x)$

∴ Therefore $R_\alpha^+$ are the f^n of $|D_\alpha|$. ⟹ Any function of $Z_\alpha$'s and $|D_\alpha|$'s are ind.

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

---

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

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

$= \sum E(Z_\alpha) \cdot E(R_\alpha^+)$

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

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

---

- Determination of rejection region by $T^+$

To determine the rejection region, the probability distribution of $T^+$ has to be determined under $H_0$.

$H_0: M = M_0 \equiv H_0: P(X_\alpha > M_0) = \frac{1}{2}$

(left) $M < M_0 \equiv \pi > \frac{1}{2}$

(right) $M > M_0 \equiv \pi < \frac{1}{2}$

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. $\left[2^n \text{ possibility} \atop (2 \times 2 \times \cdots 2)\right]$.

Each of this distinguishable arrangements is equally likely under $H_0$. Then the null distribution of $T^+$ is. $P(T^+=t) = \dfrac{u(t)}{2^n}$

where, $u(t)$: # of ways to assign '+' and '−' sign on the $n$ integers $(1, 2, \ldots, n)$ ranks of

Such that the sum of positive obs. is $t$.

Every assignment has a conjugate assignment interchanging + sign to − sign (vice-verra).

Ex: $n=3$

| $T^+$ | Ranks associated to t | | Values of $P(T^+=t)$ |
|---|---|---|---|
| 0 | − | − − − | $P(T^+=0)=1/8$ |
| 1 | 1 (+) | + − − | $P(T^+=1)=1/8$ |
| 2 | 2 (+) | − + − | $P(T^+=2)=1/8$ |
| 3 | 3, 1 1,2 | − − +, + + − | $P(T^+=3)=2/8$ |
| 4 | 1, 3 | (+ − +) | $P(T^+=4)=1/8$ |
| 5 | 2, 3 | − + + | $P(T^+=5)=1/8$ |
| 6 | 1, 2, 3 | + + + | $P(T^+=6)=1/8$ |

☐ Distribution of $T^+$ is symmetric.

☑ conjugate pair ✗ (proves symmetricity) always exists

$n=4$

| $T^+$ | Ranks associated to t | $P(T^+=t)$ |
|---|---|---|
| 0 | − − − − | $\dfrac{u(t)}{2^n} = 1/16$ |
| 1 | + − − − | 1/16 |
| 2 | − + − − | 1/16 |
| 3 | + + − −, − − + − | 2/16 |
| 4 | + − + −, − − − + | 2/16 |
| 5 | − + + −, + − − + | 2/16 |
| 6 | + − − +, + + + − | 2/16 |
| 7 | − − + +, + + − + | 2/16 |
| 8 | + − + + | 1/16 |
| 9 | − + + + | 1/16 |
| 10 | + + + + | 1/16 |

H.W
Using any arbitrary points show that

$T^+$ is symmetric ($n=4$) around its mean, 5.

**⊛** The educational testing service(ETS) reports that the $75^{th}$ percentile for scores of the GRE examinations is 693 in a ~~year~~ certain year. A random sample of 15 freshmen majoring in statistics report their GRE scores as

690, 750, 680, 700, 660, 710, 720, 730, 650, 670, 740, 730, 660, 750, 690

Are the scores of students majoring in statistics consistent with the $75^{th}$ percentile value?

$$\left[\left\{S - \frac{np}{2}\right\} > k_{\alpha/2}\right]$$
$$_{or} \ S - np < k'_{\alpha/2}$$

---

**⊛** Using any arbitrary points show that, $T^+$ is symmetric about its mean for, $n=4$

$$\Rightarrow Pr\left(T^+ > \frac{n(n+1)}{4} + c\right)$$
$$Pr\left(\frac{n(n+1)}{2} - T^+ < \frac{n(n+1)}{2} - \frac{n(n+1)}{4}\right) - c\right)$$
$$Pr\left(T^- < \frac{n(n+1)}{4}\right) - c\right)$$
$$= Pr\left(T^+ < \frac{n(n+1)}{4} - c\right) \left[\begin{array}{l} \therefore T^+ \text{ and } T^- \text{ are} \\ \text{linearly related,} \\ \text{and identically distributed} \end{array}\right]$$

$\therefore$ $T^+$ has its distribution symmetric about mean $\forall n$

---

Practical-2

Educational Testing Agency reports that $75^{th}$ percentile scores of GRE exams is 693.

Let $p$ be the probability GRE scores lies in the percentile range.

We test
$$H_0 : p = \frac{3}{4} \quad vs \quad H_1 : p \neq \frac{3}{4}$$

We have a sample of GRE scores of size 15.

Under $H_0$, the no. of observations satisfying

$$P(X_i < 693) = \frac{3}{4}, \quad S_{H_0} \sim Bin\left(15, \frac{3}{4}\right)$$

$$\therefore P_{H_0}(S = 8) = \binom{15}{8}\left(\frac{3}{4}\right)^8 \cdot \left(\frac{1}{4}\right)^{15-8}$$

Now,

| Obs − 693 | Sign |
|---|---|
| 690 − 693 | − |
| 750 − 693 | + ✓ |
| 680 − 693 | − |
| 700 − 693 | + |
| 660 − 693 | − |
| 710 − 693 | + |
| 720 − 693 | + |
| 730 − 693 | + |
| 650 − 693 | − |
| 670 − 693 | − |
| 740 − 693 | + |
| 730 − 693 | + |
| 660 − 693 | − |
| 750 − 693 | + |
| 690 − 693 | − |

$\therefore$ the no. of positive differences = the value of test statistic

$$\boxed{S = 8}$$

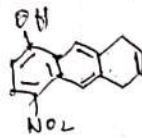Fix $\alpha =$ (level of significance)
$$= 0.1$$

let $k_{\alpha/2}$ and $k'_{\alpha/2}$ be two constants constructing rejection region

From the size condition,

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

When $k_{\alpha/2} = 7$ we find that

$$P_{H_0}(s \leq 7) \leq 0.05$$

Also,

$$\sum_{s=k'_{\alpha/2}}^{15} \binom{15}{s}\left(\frac{3}{4}\right)^{s}\left(\frac{1}{4}\right)^{15-s} \leq 0.05$$

finding $k'_{\alpha/2} = 14$ satisfies the inequality
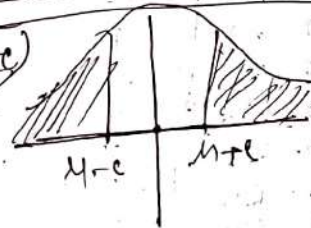
Hence the test is constructed by

$$\phi(s) = \begin{cases} 1 & \text{if } s \leq 7 \text{ or } s \geq 14 \\ \\ 0 & \text{o.w} \end{cases}$$

As we have $s=8$, we can conclude that, students are majoring in Statistics has scores consistent with the percentile value (Accepting $H_0$)

⑧ $P(X > \mu + c) = P(X \leq \mu - c)$

~~$\phi(\mu) > X(c)$~~

$$= P(X - \mu > c)$$

---

Prove that $T^{+}$ is symmetric.

In the construction of $T^{+}$, every assignment has a conjugate assignment with plus and minus sign interchange. Since we defined $Z_{\alpha} = \begin{cases} 1 & \text{if } X_{\alpha} > M_0 \\ \\ 0 & \text{if } X_{\alpha} < M_0 \end{cases}$

Conjugate variable of $Z_{\alpha}$,

$$1 - Z_{\alpha} = \begin{cases} 1 & \text{if } X_{\alpha} < M_0 \\ \\ 0 & \text{if } X_{\alpha} > M_0 \end{cases}$$

∴ Therefore the value of $T^{+}$ for those conjugate assignments will be

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

Since every assignment occurs with equal prob. $\frac{1}{2^n}$ this implies that

$T^{+}$ is symmetric.

$T_{conj}^{+} \to T^{-}$

∴ $T^{+}$ is symmetric about $\frac{n(n+1)}{4}$.

$$T_{conj}^{+} = \frac{n(n+1)}{2} - T_{org}^{+}$$

$$T_{conj}^{+} - \frac{n(n+1)}{4} = \frac{n(n+1)}{4} - T_{org}^{+}$$

## Result

$T^+$ and $T^-$ are identically distributed

$$\left[\begin{array}{c} P(x>\mu+c) \\ = P(x<\mu-c) \end{array}\right]$$

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

$$\doteq P\left[\frac{n(n+1)}{4} - T^+ \geq c - \frac{n(n+1)}{4}\right] \qquad \left[\because T^+ \text{ is symmetric}\right]$$

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

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

## Result (Rejection region will be on the left-side)

Since it is more convenient to work with the smaller term, so we use $T^+$ or $T^-$ accordingly. If $t_\alpha$ is the critical point such that $P(T \leq t_\alpha) = \alpha$ the rejection region for different alternative will be as follows:

| $H_1$ | Interpretation |
|---|---|
| $H_1 : \mu > \mu_0$ | $T^+$ higher . to reject $T^-$ lower , $P(T^- \leq t_\alpha) = \alpha$ |
| $H_1 : \mu < \mu_0$ | $T^+$ lower, to reject $T^-$ higher $P(T^+ \leq t_\alpha) = \alpha$ |
| $H_1 : \mu \neq \mu_0$ | we reject $T^+ \leq t_{\alpha/2}$ or $T^- \leq t_{\alpha/2}$ |

※ For every choice of $n$ and $\alpha$, the cut off point may not be found in Wilcoxon-signed rank test.

Therefore

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

(ii) The critical point is not found does not imply that the test is invalid.

● For paired (bivariate) observations, both sign test and Wilcoxon-signed rank test can be applied — by constructing the test on the differences $D_\alpha^* = X_\alpha - Y_\alpha$ as the univariate observation.

Practical (H.W)

In a marketing research test, 15 adult males were asked to shape one side of their face with a brand A razor and the otherside of their face with a brand B razor and state their preferred razor.

12 men preferred brand A. Find the p-value for the alternative for preferring brand A is greater than 0.05.

H₀: A and B are equally likely

*The preference is more / medium shifts.*

$H_0: M = M_0$ (proportion of finding choice of A)

$H_1: M > M_0$

Let $S$ be sample statistic is no. of adults preferring A.

$S = 12$

We use sign test, sign test where

$$S \underset{H_0}{\sim} Bin\left(15, \tfrac{1}{2}\right)$$

p-value
---
$P(\text{Reject } H_0 \mid H_0 \text{ is true})$

$= P\left(S \leq 12 \mid P = \tfrac{1}{2}\right)$      $\alpha = 0.11 \text{ (say)}$

Reject...

$= P\left(S - \ldots \mid P = \tfrac{1}{2}\right) \geq 12$

$= P(S - \ldots \geq S_\alpha)$

---

A study 5 years ago reported that liniant amount of sleep by American adults in 7.5 hours. out of 24 hours. A current sample of 8 adults reported their average amount of sleep per 24 hrs. as 7.2, 8.3, (5.6) 7.4, 7.8, 5.2, 9.1, 5.8 hours.

Use the most apt test to determine whether American adults sleep less today than 5 years ago.   $\alpha = 0.05$, $n = 8$

| | | |
|---|---|---|
| -0.3 → 0.3 | 2.5 | |
| 0.8 → 0.8 | 4 | |
| -1.9 → 1.9 | 7 | $T^+ = 11.5$ |
| -0.1 → 0.1 | 1 | |
| 0.3 → 0.3 | 2.5 | $T^- = 24.5$ |
| -2.3 → 2.3 | 8 | We can't reject |
| 1.6 → 1.6 | 5 | $\dfrac{8(8+1)}{2}$ |
| -1.7 → 1.7 | 6 | |

$n = 8$, $\alpha = 0.05$

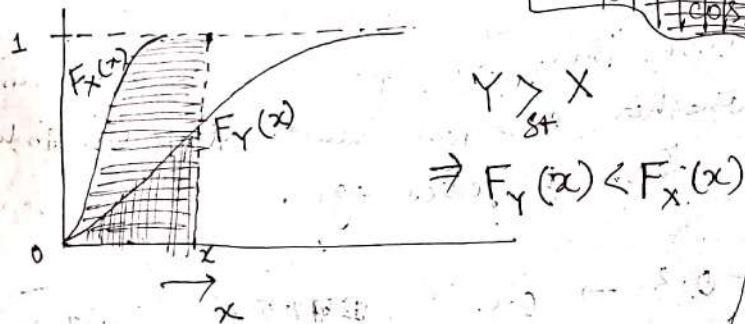We reject if $T^+ \leq \ldots$      $\Rightarrow 36$

## Two population Median test

Let $X_1, X_2, \ldots, X_{n_1}$ come from a continuous population with c.d.f $F_X(x)$

And $Y_1, Y_2, \ldots, Y_{n_2}$ come from another independent continuous population with c.d.f $F_Y(x)$.

The r.v $Y$ is called stochastically larger to $X$ if $Y$ takes some probability for higher values while $X$ takes that probability for lower value.



$$Y \underset{st}{>} X$$

$$\Rightarrow F_Y(x) < F_X(x)$$

to match the area covered by $F_X(i)$, $x$ has to be increased.

Remark: Two population non-parametric location test is based on the idea of equality of two medians ($\mu_X$ and $\mu_Y$).

$$Y \underset{st}{>} X \Rightarrow \mu_Y > \mu_X \Rightarrow F_Y(x) < F_X(x)$$

• We are to test

$H_0: \mu_X = \mu_Y$

$H_1: \begin{matrix} \mu_X < \mu_Y \\ \mu_X > \mu_Y \end{matrix} \equiv \begin{matrix} Y \underset{st}{>} X \\ \mu_X + \mu_Y = Y \underset{st}{\leq} X \end{matrix}$

## Mann-Whitney Test

M-W U test is a special choice of testing the above where it is assumed that two populations are differed by a location parameter $\theta$. i.e

$$F_X(x) = F_Y(x+\theta)$$

Therefore, $H_0: \mu_X = \mu_Y$          $H_0: \theta = 0$

against $H_1: \mu_Y > \mu_X \equiv H_1: \theta > 0$

(Analogous to test).

For testing the above, we check how many of $Y$ sample observations are less than $X$ observations in the combined sample.

Define, $D_{ij} = \begin{cases} 1 & \text{if } Y_j < X_i \\ 0 & \text{O.W} \end{cases} \quad \begin{matrix} i = 1(1)n_1 \\ j = 1(1)n_2 \end{matrix}$

∴ U-statistic is

$$U = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} D_{ij} \Rightarrow \text{The no. of times } Y \text{ preceeds } X.$$

Clearly small value of $U$ reject $H_0$. Therefore the test based on $U$ will be a left-tailed test.

❋ $\boxed{E(U), \text{Var}(U)}$ ⇒ Is to be calculated.

Assume $P(Y_j < X_i) = \pi$ (say)    $i = 1(1)n_1$
$j = 1(1)n_2$

$E(D_{ij}) = \pi$

$\therefore E(U) = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} E(D_{ij}) = n_1 n_2 \pi$

$Var(D_{ij})$

$= E(D_{ij})\left[1 - E(D_{ij})\right]$

$= \pi(1-\pi)$

$D_{ij}$'s are independent variable with $D_{lk}$
as $Y_j < X_i$ and $Y_k < X_l$ are two
independent events [paired obs. are
different].

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

Let $P(Y_j < X_i$ and $Y_k < X_i)$

$= P(\max(Y_j, Y_k) < X_i)$

$= \int_{-\infty}^{\infty} \left[F_Y(x)\right]^2 dF_X(x) = \pi_1$

Also $P(Y_j < X_i$ and $Y_j < X_k)$

$= P(\min(X_i, X_k) > Y_j)$

$= \int_{-\infty}^{\infty} (1 - F_X(x))^2 dF_Y(x) = \pi_2$

$Cov(D_{ij}, D_{ik})$

$= E(D_{ij} D_{ik}) - E(D_{ij}) E(D_{ik})$

$= \pi_1 - \pi^2$

Similarly, $Cov(D_{ij}, D_{kj})$

$= \pi_2 - \pi^2$

$\therefore Var(U)$    k is a dummy notation there.

$= Var\left(\sum_{i=1}^{n_1} \sum_{j=1}^{n_2} D_{ij}\right)$

$= \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} Var(D_{ij}) + \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \sum_{\substack{h,k \\ h \neq k}} Cov(D_{ij}, D_{hk})$

$(h,k) \neq (i,j)$
not together

$= n_1 n_2 \pi(1-\pi) + \sum_{i=1}^{n_1} \sum \sum_{1 \le j \neq k \le n_2} Cov(D_{ij}, D_{ik})$

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

$= n_1 n_2 \pi(1-\pi) + n_1 n_2 (n_2 - 1)(\pi_1 - \pi^2)$

$\qquad + n_2 n_1 (n_1 - 1)(\pi_2 - \pi^2)$

$= n_1 n_2 \left[\pi(1-\pi) + (n_2 - 1)(\pi_1 - \pi^2) + (n_1 - 1)(\pi_2 - \pi^2)\right]$

$$\therefore \text{Var}(U)$$

$$= n_1 n_2 \left[ 2\pi + (n_2-1)\pi_1 + (n_1-1)\pi_2 \right. $$
$$\left. - (n_1+n_2-1)\pi^2 \right]$$

where $N = n_1 + n_2$

$$\therefore \text{Var}\left(\frac{U}{n_1 n_2}\right) \to 0 \quad \text{as} \quad \begin{array}{c} n_1 \to \infty \\ n_2 \to \infty \end{array}$$

$$\therefore E\left(\frac{U}{n_1 n_2}\right) = \pi$$

$\therefore \dfrac{U}{n_1 n_2}$ is a consistent estimator of $\pi$.

$$= P(Y_j < X_i)$$

$$\begin{array}{c} i=1(1)n_1 \\ j=1(1)n_2 \end{array}$$

**⑩** Show that

$\left(1 - \dfrac{U}{n_1 n_2}\right)$ is consistent for $(1 - \pi)$

## Discrete Distribution of $U$.

For $n_1$ X obs. and $n_2$ Y obs. there are $\binom{n_1+n_2}{n_1}$ arrangements by X and Y in combined sample. For every particular arrangement $Z$.

There exist one conjugate arrangement as if $Z$ denotes a set of X and Y written from smallest to largest, then its

conjugate arrangement may be proposed from largest to smallest. (conjugate arrangement: how many X follow Y).

If $U$ be the value of an arrangement then the prob. dist. of its conjugate arrangement will be the same and that value is.

$$U' = \sum_{i=1}^{n_1}\sum_{j=1}^{n_2}(1 - D_{ij})$$

**Ex.** $n_1=4, n_2=5$

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

The p.m.f of $U$ is

$$P(U=u) = \frac{r_u}{\binom{n_1+n_2}{n_1}}$$

$r_u$: No of distinguishable arrangements for which r.v $U$ takes the value $u$.

Find out $E(U)$ and $V(U)$ under $H_0$

**⊙** For $H_1:$ $Y \overset{st}{>} X$ or $\mu_Y > \mu_X$

we reject $H_0$ is $U < U_\alpha$ (tab. value at $\alpha^{th}$ level of significance)

On the basis of discrete distribution and $n_1, n_2$ the table is available.

For, $H_1: Y \overset{st}{\leq} X$ or $\mu_Y < \mu_X$ we reject $H_0$ if $U' < U_\alpha$

Construct conjugate arrangement and statistic $U'$

⊛ For tied case

$$D_{ij} = \begin{cases} 1 & Y_j < X_i \\ 0.5 & Y_j = X_i \\ 0 & \text{o.w} \end{cases}$$

We might have fraction part in U for this.

①

---

⊛ The 2000 census statistics for Alabama district give the % changes in population between 1990 and 2000 for each 67 countys. There are two types of countys — rural and non-rural according to the population size less than 25000. Below is the data 9 rural and 7 non-rural countyies.

Rural : 1.1, -21.7, -16.3, -11.3, -10.4, -7, -2.0, 1.9, 6.2

Non-rural: -2.4, 9.9, 14.2, 18.4, 20.1, 23.1, 70.4

Use Mann-Whitney U test to test equal population change.

---

→ Let the population change of rural county come from a continuous distribution with c.d.f $F_Y$ and median $M_Y$. In contrast, non-rural come from $F_X$ and median $M_X$

$H_0: M_X = M_Y$

$H_1: M_X \neq M_Y$

i) Arranged combined sample

-20.7, -16.3, -11.3, -10.4, -7, -2.4, -2.0, 1.1  1.9  6.2
  Y     Y     Y     Y     Y     X     Y     Y    Y   Y

  9.9    14.2    18.4    20.1    23.1    70.4
  X      X       X       X       X       X

$U = 5 + 8 + 8 + 8 + 8 + 8 + 8 = 59$

$U' = 6 + 6 + 6 + 7 + 7 + 7 + 7 + 7 + 7 = $

$63 - 54 = 4$
$(7 \times 9)$

$U_{\alpha/2} = 9$

$n_1 = 9, \quad n_2 = 7$

$\alpha = 0.02$

$\boxed{U_{0.01, 7, 9} = 9}$

$U \cancel{<} U_{\alpha/2}$   $U' \leq U_{\alpha/2}$

rejection from one side

## Pearson's goodness of fit test
### (approximately $\chi^2$)

$$\sum_{i=1}^{n} \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{n-1}$$

## Kolmogorov - Smirnov Goodness of fit test

### Empirical cumulative distribution function

Definition: Empirical cdf of a random sample of size $n$ is denoted by $F_n(x)$ where $F_n(x) = \dfrac{\# \leq x}{n}$

$$= \frac{\text{no. of observations} \leq x}{n}$$

☒ ecdf is an estimate of $F(x)$ {e. df}.

☒ Distribution of $F_n(x)$

Define an indicator function $\delta_j(x) = \begin{cases} 1, & x_j \leq x \\ 0, & x_j > x \end{cases}$

$$P\left(\delta_j(x) = 1\right) = P\left(x_j \leq x\right) = F_x(x)$$

∴ $\delta_j(x)$ is a Bernoullian trial with $F_x(x)$

∴ $n F_n(x) = \sum_j \delta_j(x) \rightarrow$ Sum of $n$ 'ber' trials

∴ $E\left(n F_n(x)\right) = \sum_j 1 \cdot F_x(x) = n F_x(x)$

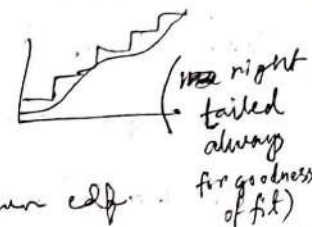∴ $n F_n(x) \sim Bin(n, F_x(x))$

∴ $F_n(x)$ is unbiased for $F_x(x)$.

is ecdf a consistent estimator of $F_x(x)$

## K.S test statistic


( is right tailed always for goodness of fit)

$H_0: F_x(x) = F_0(x) \leftarrow$ known cdf

$$D = \sup_{-\infty < x < \infty} \left| F_n(x) - F_0(x) \right|$$

☒ $D$ is distribution free test statistic
$D$ takes pairwise probability based on N.
For $n$, the cutoff value of $D$

☒ for various values $n$, the cut off value of $D$ is available in table.

$D_{cal} > D_{tab}$, we reject $H_0$.

### [ Cramer-Von-Mises
$$D = \int_{-\infty}^{\infty} \left(F_n(x) - F_0(x)\right)^2 dF_0(x)$$

## Two sample - KS Test

$$X: x_1, x_2, x_3, \ldots x_{n_1}$$
$$Y: y_1, y_2, y_3, \ldots y_{n_2}$$

are they coming from same dist?

$$D_{n_1, n_2} = \sup \left| F_{n_1}(x) - F_{n_2}(x) \right|$$

# Practical

1) 10 students take a test and their scores are as follows (out of hundred)

95, 80, 40, 52, 60, 80, 82, 58, 65, 50

Test the null hypothesis that the c.d.f of the proportion of right answer a student gets on the test is

$$F_0(x) = \begin{cases} x^2(3-2x), & 0 \le x < 1 \\ 1, & x \ge 1 \end{cases}$$

| X | P(x) | Ordered $P(x)$ | $F_n(x)$ | $F_0(x)$ | $|F_n(x)-F_0(x)|$ |
|---|---|---|---|---|---|
| 95 | 0.95 | 0.4 | 0.1 | 0.352 | 0.252 |
| 80 | 0.8 | 0.5 | 0.2 | 0.5 | 0.3 |
| 40 | 0.4 | 0.52 | 0.3 | 0.53 | 0.23 |
| 52 | 0.52 | 0.58 | 0.4 | 0.618 | 0.219 |
| 60 | 0.6 | 0.6 | 0.5 | 0.648 | 0.148 |
| 80 | 0.8 | 0.65 | 0.6 | 0.718 | 0.118 |
| 82 | 0.82 | 0.8 | 0.8 | 0.896 | 0.096 |
| 58 | 0.58 | 0.8 | 0.8 | 0.896 | 0.014 |
| 65 | 0.65 | 0.82 | 0.9 | 0.914 | 0.007 |
| 50 | 0.50 | 0.95 | 1 | 0.993 | |

$$D_{cal} = \sup |F_n(x) - F_0(x)|$$
$$= 0.3$$

$\alpha = 0.05$

$D_\alpha = 0.409$

$D_{cal} < D_{tab}$

⇒ Failed to reject $H_0$

∴ The proportion of right answer is coming from $F_0(x)$

---

2) A random sample of 12 persons are interviewed to estimate median annual gross income in a certain economically depressed town. Use the most apt. test for the null hypothesis that income data are standard normally distributed.

9800, 10200, 9300, 8700, 15200, 6800, 8600, 9600, 12200, 15500, 11600, 7200.

$H_0 : F(x) = \phi(x)$
$H_1 : F(x) \ne F(x)$

$$\frac{X - \bar{X}}{S_x} \qquad S_x^2 = \frac{1}{n-1}\Sigma(x_i - \bar{x})^2$$

③ two mutually ind. r.s each of size 8, were generated, one from the standard normal distribution and another from the chi-square distr. with df 18.

The resulting data are as follows

Norm: $-1.91$ $-1.22$ $-0.96$ $-0.72$ $0.14$ $0.82$ $1.45$ $1.86$

$\chi^2$: $4.90$ $7.25$ $8.04$ $14.10$ $18.30$ $21.21$ $23.10$ $28.12$

Do you believe they are coming from the same dist?

⟶ Convert all to standard form

$$\frac{(x^2) - 18}{\sqrt{36}}$$

Norm: $-1.91$ $-1.22$ $-0.96$ $-0.72$ $0.14$ $0.82$ $1.45$ $1.86$

$\chi^2$: $-2.183$ $-1.79$ $-1.66$ $-0.65$ $0.05$ $0.535$ $0.85$ $1.687$

We would test whether two sets of obs. are coming from the same dist.

So, $H_0: F_1(x) = F_2(x)$   $F_1$ is standard normal
$H_1: F_1(x) \neq F_2(x)$   $F_2$ is standard $\chi^2$ d F

---

First we combine both the samples increasingly

| Norm samp (1st) | Chi Sq sample (2nd) | $F_{n_1}(x)$ | $F_{n_2}(x)$ | $|F_{n_1}(x)-F_{n_2}(x)|$ |
|---|---|---|---|---|
| | $-2.183$ | 0 | 1/8 | 1/7 |
| $-1.91$ | $-1.79$ | 1/2 | 2/8 | 1/7 |
| $-1.22$ | $-1.66$ | 2/8 | 3/8 | |
| $-0.96$ | $-0.65$ | 3/8 | 4/7 | |
| $-0.72$ | $0.05$ | 4/8 | 5/8 | |
| $0.14$ | $0.535$ | 5/8 | 6/8 | |
| $0.82$ | $0.85$ | 6/8 | 7/8 | |
| $1.45$ | $1.687$ | 7/8 | 8/8 | |
| $1.86$ | | 8/8 | 8/8 | 0 |

| | $F_{n_1}(x)$ | $F_{n_2}(x)$ | $|F_1(x)-F_2(x)|$ |
|---|---|---|---|
| $-2.183\ (2^{nd})$ | 0 | 1/8 | 1/8 |
| $-1.91\ (1^{st})$ | 1/8 | 1/8 | 0 |
| $-1.79\ (2^{nd})$ | 1/8 | 2/8 | 1/8 |
| $-1.66\ (2^{nd})$ | 1/8 | 3/8 | [1/4] |
| | | | 1/8 |
| $-1.22\ (1^{st})$ | 2/8 | 3/8 | 0 |
| $-0.96\ (1^{st})$ | 3/8 | 3/8 | 1/2 |
| $-0.72\ (1^{st})$ | 4/8 | 4/8 | 0 |
| $-0.65\ (2^{nd})$ | 4/8 | 4/8 | 1/8 |
| $0.05\ (2^{nd})$ | 5/8 | 5/8 | 0 |
| $0.14\ (1^{st})$ | 5/8 | 6/8 | 1/8 |
| $0.535\ (2^{nd})$ | 6/8 | 6/8 | 0 |
| $0.82\ (1^{st})$ | 6/8 | 7/8 | 1/8 |
| $0.85\ (2^{nd})$ | 7/8 | 7/8 | 0 |
| $1.45\ (1^{st})$ | 7/8 | 8/8 | 1/8 |
| $1.687\ (2^{nd})$ | 8/8 | 8/8 | 0 |

$\max |F_{n_1}(x) - F_{n_2}(x)| = \frac{1}{4}$

$n_1 n_2 D = 8 \times 8 \times \frac{1}{4} = \boxed{16}$

From table for

$n_1 = 8, \quad n_2 = 8$

if $n_1 n_2 D = \cancel{0.283} \longrightarrow$ p-value is

$= 32$ has p-value $\boxed{=0.283}$ increasing

for $n_1 n_2 D = 16$ if we fix $\alpha = 0.05$,

p-value $> 0.283$ $\cancel{0.05}$

$\Rightarrow$ We fail to reject $H_0$.

\# From the table $n_1 n_2 D = 32$

yields the p-value as $0.283$

Therefore for $n_1 n_2 D_{cal} = 16$

the p-value will be greater than

$0.283$. So if we fix $\alpha = 0.05$,

we fail to reject $H_0$.

$\therefore$ This two obs. are coming from the

same distribution.