

Statistical Tests (Non Parametric)

Master in Data Science

September 2024

Introduction

Non- parametric tests are those tests that do not involve any parameter of the population. Some common non-parametric tests are tabulated below:

Mann-Whitney U Test

NON-PARAMETRIC TESTS			
MANN-WHITNEY U TEST			
SAMPLE SIZE	TEST STATISTIC	CRITICAL VALUE	DECISION
Small size $n_1 \leq 10$ & $n_2 \leq 10$ $n_1 + n_2 \leq 20$	$U_0 = \text{minimum of } \{U_1 \text{ and } U_2\}$ where, $U_1 + U_2 = n_1 n_2$ $U_1 = n_1 n_2 + \frac{n_1(n_1+1)}{2} - R_1$ $U_2 = n_1 n_2 + \frac{n_2(n_2+1)}{2} - R_2$	$U_{\alpha}, (n_1, n_2)$ [U Value Approach] From Mann-Whitney probability table, find p-value for (n_1, n_2) & U_0 . for two-tailed, $p\text{-value} = 2P_0$ for one-tailed, $p\text{-value} = P_0$ [p-value approach]	If $U_0 > U_{\alpha}, (n_1, n_2)$, do not reject H_0 , else reject H_0 . if $p\text{-value} > \alpha$, then we do not reject H_0 , otherwise reject H_0 .
Large size $n_1 > 10, n_2 > 10$ $n_1 + n_2 > 20$	for large sample size, the distribution of U_0 is approximated by normal dist'n with mean = $\frac{n_1 n_2}{2}$ and variance = $\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}$ Test statistic: $U_0 - \frac{n_1 n_2}{2}$ $Z = \frac{\frac{U_0 - \frac{n_1 n_2}{2}}{\sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}}}$ In case of tied observations, corrected S.d; $U_0 = \frac{n_1 n_2}{\sqrt{n(n-1)}} \left[\frac{n^3 - n}{12} - \frac{\sum (t_i^3 - t_i)}{12} \right]$ where, $n = n_1 + n_2$ t_i = number of times i^{th} rank is repeated.	Z_{α} (From Z-table)	If $Z_{\alpha} (\text{table}) > Z_{\text{calculated}} $, we accept H_0 , else reject H_0 .

Median Test

MEDIAN TEST		$\alpha \Rightarrow \text{minimum of } (n_1 \text{ and } K)$																	
SAMPLE SIZE	TEST STATISTIC	CRITICAL VALUE	DECISION																
$n_1 \leq 10$ & $n_2 \leq 10$ [Small]	'a' Test-statistic is obtained as follows: * Combine both the samples and arrange them in ascending order of magnitude such that $n = n_1 + n_2$ * Calculate median of combined sample and count no. of observations less than or equal to median in first sample. This is 'a'. $K = \frac{n_1 + n_2}{2}$	$P_0 = P(A \geq a)$ $P(A=a) = \frac{{}^{n_1}C_a {}^{n_2}C_{K-a}}{{}^{n_1+n_2}C_K}$ $* P_0 = \sum_{a=K}^K \frac{{}^{n_1}C_a {}^{n_2}C_{K-a}}{{}^{n_1+n_2}C_K}$	for two-tailed: If $2P_0 > \alpha$, accept H_0 else reject H_0 for one-tailed: If $P_0 > \alpha$, accept H_0 else reject H_0																
[Large Sample Size] $n_1 > 10$ & $n_2 > 10$	In this test, the median test is equivalent to χ^2 -test Test-statistic: χ^2 Use 2x2 Contingency table <table border="1"> <thead> <tr> <th></th><th>No of obs $\leq Md$</th><th>No of obs $> Md$</th><th>TOTAL</th></tr> </thead> <tbody> <tr> <td>Sample I</td><td>a</td><td>b</td><td>a+b</td></tr> <tr> <td>Sample II</td><td>c</td><td>d</td><td>c+d</td></tr> <tr> <td>Total</td><td>a+c</td><td>b+d</td><td>a+b+c+d</td></tr> </tbody> </table> $\chi^2 = \frac{n(ad-bc)^2}{(a+b)(c+d)(a+c)(b+d)}$		No of obs $\leq Md$	No of obs $> Md$	TOTAL	Sample I	a	b	a+b	Sample II	c	d	c+d	Total	a+c	b+d	a+b+c+d	Extract χ^2 value from χ^2 table for 1 dof and α level of significance. i.e. $\chi^2_{\alpha, 1}$	If $\chi^2_{cal} > \chi^2_{tab}$, we reject H_0 else accept H_0
	No of obs $\leq Md$	No of obs $> Md$	TOTAL																
Sample I	a	b	a+b																
Sample II	c	d	c+d																
Total	a+c	b+d	a+b+c+d																

Kolmogorov Smirnov Test

SAMPLE	TEST STATISTIC	CRITICAL VALUE	DECISION
<u>One Sample</u> To check if there is significant difference between observed and expected frequency	$D_0 = \text{Maximum } F_e(x) - F_o(x) $ where, $F_e(x) = \frac{c f_e}{n}$ $c f_e = \text{expected cumulative frequency}$ $f_e = \text{expected frequency}$ $= n p_i$ $= \frac{\sum f}{\text{no. of categories}}$ $F_o(x) = \frac{c f_o}{n}$; $c f_o = \text{observed cumulative frequency}$	$D_{\text{tabulated}} = D_{n, \alpha}$	If $D_0 \geq D_{n, \alpha}$, then we reject H_0 otherwise do not reject H_0 .
<u>TWO SAMPLES</u> <u>Small sample size:</u> $n_1 = n_2 \leq 40$ or $n_1 \neq n_2 \leq 20$ <u>Large sample size:</u> $n_1, n_2 > 40$ (for $n_1 = n_2$) $n_1, n_2 > 20$ (for $n_1 \neq n_2$)	$D_0 = \text{Maximum } F(x) - F_y $ where, $F(x) = \frac{c f_x}{n_1}$, $F_y = \frac{c f_y}{n_2}$ $D_0 = \text{Maximum } F(x) - F_y $ (for Two Tailed Test) $\chi^2 = 4 D_0^2 \left(\frac{n_1 n_2}{n_1 + n_2} \right)$ (for One Tailed Test)	$D_{\text{tabulated}} = D_{(n_1, n_2), \alpha}$ $\chi^2_{\text{calculated}} = \chi^2_{(\alpha, 2)}$ \uparrow degree of freedom	If $D_{(n_1, n_2), \alpha} > D_0$, we don't reject H_0 , else we do. If $D_{(n_1, n_2), \alpha} > D_0$, we do not reject H_0 , else we do. If $\chi^2_{\text{cal}} \geq \chi^2_{(\alpha, 2)}$, we reject H_0 , else we accept H_0 .

Wilcoxon and Kruskal Wallis H Test

WILCOXON MATCHED PAIR SIGN RANKED TEST

(used for small, usually $n \leq 20$ data)

SAMPLE	TEST STATISTIC	CRITICAL VALUE	DECISION
Use for small sample size ($n \leq 20$).	$T = \text{Minimum of } \{S(+), S(-)\}$ where, $S(+) = \text{Sum of ranks of difference with '+' sign}$ $S(-) = \text{Sum of ranks of difference with '-' sign.}$	$T_{\text{tabulated}} = T_{\alpha, n_e}$ $n_e = \text{effective sample size} = n - t$ $t = \text{no. of differences with zero.}$	If $T_{\alpha, n_e} \geq T$, we reject H_0 , else we accept H_0 .

KRUSKAL WALLIS H-TEST

SAMPLE	TEST STATISTIC	CRITICAL VALUE	DECISION
Small sample: $n_i \leq 5$ and $k \geq 3$	$H = \frac{12}{n(n+1)} \left(\sum_{i=1}^k \frac{R_i^2}{n_i} \right) - 3(n+1)$ Where, $n = n_1 + n_2 + \dots + n_k$ If there is tie in observations, then corrected H , $H_{\text{corr}} = \frac{H}{C.F.}$ $C.F. = \text{Correction Factor} = 1 - \frac{\sum (t_i^3 - t_i)}{n^3 - n}$	Obtain the p-value	If $p\text{-value} > \alpha$, we accept H_0 , else we reject H_0 .
Large sample: $n_i > 5$ and $k > 3$	Same test statistic as above	$\chi^2_{\text{tabulated}} = \chi^2_{\alpha, (k-1)}$ \downarrow degree of freedom	If $H < \chi^2_{\text{tab}}$, accept H_0 else reject H_0 .

Friedman F Test

FRIEDMAN F-TEST

k = no. of samples, n = size of each sample

Sample	Test Statistic	Critical Value	Decision
<u>Small Sample</u> $2 \leq n \leq 9$ and $k=3$ and $2 \leq n \leq 4$ and $k=4$	$F_r = \frac{12}{nk(k+1)} \left(\sum_{i=1}^k R_i^2 \right) - 3n(k+1)$ <p>If tied case then, then corrected value of F_r is</p> $F_{r \text{ corrected}} = \frac{F_r}{C.F.}$ $C.F. = 1 - \frac{\sum (t_i^3 - t_i)}{n(k^3 - k)}$ <p>t_i = no. of times ith rank is repeated</p>	<p>P-value is obtained from Friedman probability table.</p>	<p>If P-value > α, we accept H_0, else we reject H_0.</p>
<u>Large sample</u> $n > 5$ & $k > 3$ and	<p>Same test statistic as above</p>	$\chi^2_{\text{tabulated}} = \chi^2_{\alpha, (k-1)}$ <p>↓ degree of freedom</p>	<p>If $F_r < \chi^2_{\text{tab}}$, we accept H_0, else reject H_0.</p>