

BY KUNAL DAS

class: BSc 6th Semester (Major)

Paper: 6.1(M)

NON-PARAMETRIC TEST

A non-parametric (N.P.) test is a test that does not depend on the particular form the basic distribution (e.g. Normal) from which the samples are drawn. In many cases an experimenter does not know the form of basic frequency function and needs statistical test techniques which are applicable regardless of the form of the density function. Another term which is often used - interchangeably with non-parametric is distribution free. In such procedures assumptions regarding the population are not necessary. The term "Non-Parametric" is due to J. Wolfowitz who used it to indicate that the population could not be specified by finite number of parameters. A non-parametric test is concerned with the form of the population and not with any parametric values. The term "distribution free" is used to denote the distribution of the test statistic and the term "non-parametric" is used to denote the type of hypothesis problem to be investigated.

According to Gibbons, a statistical technique is said to be non-parametric if it satisfies one of the following five criteria.

- (a) The data are count data of number of observations in each category.
- (b) The data are nominal in scale.
- (c) The data are ordinal in scale.
- (d) The inference does not concern a parameter.
- (e) The assumptions are general rather than specific.

Assumptions of Non-Parametric Test:

- (a) Sample observations are independent.
- (b) The variable under study is continuous.
- (c) The probability density function of a random variable is continuous.
- (d) Lower order moment exist.

Advantages and drawbacks of Non-Parametric test over parametric test: (See Fundamentals of M.S by Gupta & Kapoor)

Advantages of N.P. test

- (a) N.P. tests are readily comprehensible very simple and easy to apply and don't require complicated sample theory.
- (b) No assumption is made about the form of the frequency distribution of the general population from which the sample is taken.
- (c) N.P. technique will apply to the data which are mere classification (i.e. which are measured in nominal scale), while NP test exist to deal with such data.
- (d) Since the socio-economic data are not, in general, normally distributed, NP tests have found applications in Psychometry, Sociology and Educational Statistics.
- (e) NP tests are available to deal with the data which are given in ranks or whose seemingly numerical scores have the strength or ranks. e.g. NP test can be applied if the scores are given in grades such as A+, A-, B, A, B etc.

Drawbacks of NP test.

- (a) NP tests can be used only if the measurements are nominal or ordinal. Even in that case, if a parametric test exists it is more powerful than NP test.

2. So far, no NP test exists for testing interactions in "analysis of variance" models unless special assumptions about the additivity of the model are made.
3. NP tests are designed to test statistical hypothesis only and not for estimating the parameters.

Single Sample Problem:

Problem of Location: Let x_1, x_2, \dots, x_n be a sample of size n from some unknown distribution function (DF) $F_x(x)$. We assume that $F_x(x)$ is absolutely continuous. Let μ be a positive real number, $0 < \mu < 1$. Here an appropriate measure of location is median or the p th quantile M_p (say) the p th quantile of $F_x(x)$. One may wish to know whether the given M_0 is median of the distribution, or whether $F_x(x)$ is symmetric probability distribution. To test $H_0: M_p = M_0$ we first consider sign test.

SIGN TEST / ORDINARY SIGN TEST / SIGN TEST FOR LOCATION OF UNIVARIATE POPULATION

A random sample of N observations x_1, x_2, \dots, x_N is drawn from a population with unknown median M . Here $F_x(x)$ is assumed to be continuous around M . In other words, the assumptions are independent observations and $P_x(x=M) = 0$ or $P_x(x-M=0) = 0$. The hypothesis to be tested concerns the value of the population median

$$H_0: M = M_0$$

$$\text{against } H_1: M \neq M_0 \text{ (two-tailed)}$$

The corresponding alternative hypothesis can be one sided- or two sided on the value of M .

for any distribution which satisfies $P_0[X=M] \geq 0$,
by the definition of M , we have

$$P_r[X > M] = P_r[X < M] = \frac{1}{2}$$

Since the hypothesis here states that M_0 is the value of X which divides the area under the frequency distribution into two equal parts and symbolic representation of H_0 is equivalent to $H_0: \theta = P_r[X > M_0] = P_r[X < M_0]$

If the sample data are consistent with the hypothesized median value, on the average of half of the sample observation will lie above M_0 and half below. Thus the number of observations above M_0 which will be denoted by ' K ' can be used to test the validity of the null hypothesis. The observation will constitute a set of N independent random variables from the Bernoulli population with parameter $\theta = P_r[X > M_0]$ regardless of the population $f_X(x)$. The sampling of the random variable K thus is the Binomial probability distribution with parameter θ . The value of $\theta = \frac{1}{2} = 0.5$ as the null hypothesis is true. Since K is actually the number of positive signs (+) among the N differences $X_i - M_0$; $i=1, 2, \dots, N$ the non-parametric test based on K is called sign test.

To test sign test, theoretically no problem of zero differences exists. The population was assumed to be continuous around median. But in reality zero difference can occur. The usual procedure follows here is simply to ignore such differences and reduce N .

Test Procedure:

So the null hypothesis becomes equivalent to testing

$$H_0: p = \frac{1}{2} \quad \text{since } \text{randB}(n, 0.5 = \frac{1}{2})$$

Against $H_1: p \neq \frac{1}{2}$ (whatever the case may be)
or $H_1: p < \frac{1}{2}$

The test criterion is reject H_0 if $\tau \geq \tau_{\alpha/2}$ where
 $\tau_{\alpha/2}$ is the critical value of all significance α , $\tau_{\alpha/2}$ is
the smallest integer which satisfies the conditions.

$$\sum_{\tau=0}^N \tau \cdot \binom{1}{2}^\tau \left(\frac{1}{2}\right)^{N-\tau} \leq \alpha/2$$

or $\tau \leq \tau_{\alpha/2}$ is the smallest integer such that-

$$\sum_{\tau=0}^{\tau_{\alpha/2}} \binom{N}{\tau} \left(\frac{1}{2}\right)^\tau \left(\frac{1}{2}\right)^{N-\tau} \leq \alpha/2$$

(for small sample)

for large sample: if $n \gg 25$, the normal test can be used to decide H_0 .

The Z statistic is given by

$$Z = \frac{(r + 0.5) - n \cdot \frac{1}{2}}{\sqrt{n \cdot \frac{1}{2} \cdot \frac{1}{2}}} = \frac{(r + 0.5) - \frac{n}{2}}{\sqrt{\frac{n}{2}}} \quad \text{where } r < \frac{n}{2}$$

$$= \frac{(r - 0.5) - \frac{n}{2}}{\sqrt{\frac{n}{2}}} \quad \text{where } r > \frac{n}{2}$$

 PAIRED SIGN TEST  SIGN TEST FOR LOCATION OF BIVARIATE POPULATION.

Let x_1, x_2, \dots, x_n and y_1, y_2, \dots, y_n be two random samples of same size "n" drawn from two populations with pdf $f_1(x)$ and $f_2(y)$. We want to test the null hypothesis,

$$H_0: f_1(x) = f_2(y).$$

The observations are arranged in pairs i.e. $(x_i, y_i); i=1, 2, \dots, n$ and it is assumed that

Page No. 06

each pair is observed under identical conditions.
 Now $d_{ij} = (x_i - y_j)$ is measured and only the sign (+ or -) is noted in view of the actual deviations.
 Now, under the null hypothesis the probability that the first sample exceeds the first observations of the second sample is equal to the probability that the first observation of the second sample exceeds the first observation of the first sample, and the probability of a tie is zero. So it can be written as —

$$H_0: \Pr[X - Y > 0] = \frac{1}{2} \text{ and } \Pr[X - Y < 0] = \frac{1}{2}$$

Let us define

$$v_{ij} = \begin{cases} 1 & \text{if } x_i - y_j > 0 \\ 0 & \text{if } x_i - y_j \leq 0 \end{cases}$$

So, v_{ij} is a Bernoulli variable with $p = \Pr(X_i - Y_j > 0) = \frac{1}{2}$. Since v_{ij} , $i=1, 2, \dots, n$ are independent, $V = \sum_{i=1}^n v_{ij}$, the total number of positive deviations, is a binomial variable with parameters n and $p = \frac{1}{2}$ (under H_0). Let K be the no. of positive deviations, so

$$\Pr[V \leq K] = \sum_{r=0}^K {}^n C_r p^r q^{n-r} = \left(\frac{1}{2}\right)^n \sum_{r=0}^K {}^n C_r = \frac{1}{2^n} \quad (\text{by symmetry})$$

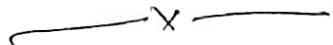
$$\therefore p = \frac{1}{2} = \alpha$$

If $\beta \leq 0.05$, then we reject H_0 at 5% level of significance. If $\beta > 0.05$ then we conclude that the data doesn't go against H_0 and so H_0 is accepted.

For large samples ($n \geq 30$), we may regard V to be asymptotically normal with $E[V] = np = \frac{n}{2}$ and $\sqrt{V} = \sqrt{npq} = \sqrt{\frac{n}{4}}$

$$\therefore Z = \frac{V - E[V]}{\sqrt{V}} = \frac{V - \frac{n}{2}}{\sqrt{\frac{n}{4}}} \sim N(0, 1)$$

and the decision will be taken on the basis of normality test for null hypothesis H_0 .



WILCOXON SIGNED RANK TEST/WILCOXON SIGNED RANKSUM TEST

The ordinary sign test utilises only the signs of difference between each observation and the hypothesized median M_0 . Here the magnitude of these observations relative to M_0 are ignored. Assuming that such information is available, a test statistic which takes into account the individual relative magnitude might be expected to give better performance. If we are willing to make assumption that the parent population is symmetric, the Wilcoxon signed rank test provides an alternative test of location which is affected by both the magnitude and signs of these differences.

Test procedure:

Let the sample x_1, x_2, \dots, x_n be drawn from a population with CDF $F(x)$. For testing $H_0: F(m) = \frac{1}{2}$ against $H_1: F(m) \neq \frac{1}{2}$
 or $H_1: F(m) > \frac{1}{2}$
 or $H_1: F(m) < \frac{1}{2}$

case 1: for sample size $n > 25$.

Let us suppose that $m=0$. So the distn is symmetric about the origin. We have $F(-x) = 1 - F(x)$ and $f(-x) = f(x)$
 then test statistic $W \sim N\left(0, \frac{n(n+1)(2n+1)}{6}\right)$ for testing
 $H_0: m=0$ for which $T = \frac{W}{\sqrt{\frac{n(n+1)(2n+1)}{6}}} \sim N(0,1)$

case 2: for sample size $n > 25$

for testing $H_0: m=m_0$ i.e. $H_0: F(m=m_0) = \frac{1}{2}$
 we have $W \sim N\left(\frac{n(n+1)}{4}, \frac{n(n+1)(2n+1)}{24}\right)$

for which $T = \frac{W - n(n+1)/4}{\sqrt{\frac{n(n+1)(2n+1)}{6}}} \sim N(0,1)$

~~INR~~ Case III Sample size $n < 25$

For testing $H_0: m = m_0$ we first find out the differences $d_i = x_i - \mu_0$ under H_0 . We can assume that the values of d_i are independent and come from a population symmetrical about zero, we then find $|d_i|$, the absolute difference. These absolute differences are then arranged in ascending order and are accordingly ranked. The absolute difference with zero values are ignored. Let the number of observations be now $n_1 < n$ i.e. if the tied ranks are eliminated. In case of a tie in rank, each of the tied values are given the average value of the rank. Let T^+ be the sum of ranks for positive d_i 's and T^- be the sum of ranks for negative d_i 's, Then $T^+ + T^- = \frac{n_1(n_1+1)}{2}$

The null distribution of T^+ and T^- are identical and symmetrical about the value $\frac{n_1(n_1+1)}{2}$. The smaller (minimum) of the two values T^+ and T^- is compared with the table value (critical value for T) in the Wilcoxon signed rank test for a given level of significance for n_1 observations and accordingly about the null hypothesis is taken.

of AH . Hyp. $H_1: m > m_0$ Then we reject H_0 if $T^- \leq T_{\alpha}$

$H_1: m < m_0$ Then we reject H_0 if $T^+ \leq T_{\alpha}$

$H_1: m \neq m_0$ Then we reject H_0 if $T^- \leq T_{\alpha/2}$ or $T^+ \leq T_{\alpha/2}$

— X —

Example 1: A drug was injected to a fresh group of 10 rats every day. The scientist -in-charge of the experiment made a claim that not more than 3 rats showed an increase in blood pressure on an average. The following increase in blood pressure was noticed in the following number of rats in the last 10 days after the drug was administered:

2, 4, 5, 1, 6, 3, 2, 1, 7 and 8.

Solution : Here the null hypothesis is $H_0 : m = 3$ tested against $H_1 : m > 3$. To perform the calculations we construct the following table:

x_i	$d_i = x_i - 3$	Signs
2	-1	-
4	1	+
5	2	+
1	-2	-
6	3	+
3	0	ignored
2	-1	-
1	-2	-
7	4	+
8	5	+

Thus, $n = 9$, and the number of plus signs $= x = 5$. Under the null hypothesis $X \sim B\left(9, \frac{1}{2}\right)$.

$$\begin{aligned}
 \text{So, } P(X \geq 5) &= 1 - P(X < 5) = 1 - \sum_{x=0}^4 {}^9C_x \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{9-x} \\
 &= 1 - \left[\left(\frac{1}{2}\right)^9 + 9 \times \left(\frac{1}{2}\right)^9 + 36 \times \left(\frac{1}{2}\right)^9 + 84 \times \left(\frac{1}{2}\right)^9 + 126 \times \left(\frac{1}{2}\right)^9 \right] \\
 &= 1 - \left(\frac{1}{2}\right)^9 (1 + 9 + 36 + 84 + 126) = \frac{1}{2} = 0.5
 \end{aligned}$$

Thus, $P(X \geq 5) > 0.05$ (level of significance α). So, H_0 is accepted and thus the claim made by the scientist is true.

Example 2: Suppose that we want to test the hypothesis that the median body length (θ) of frogs of a particular variety is $\theta_0 = 6.9$ cms against the alternative hypothesis $\theta_0 \neq 6.9$ cms with $\alpha = 0.05$ on the basis of the following measurements.

6.3, 5.8, 7.7, 8.5, 5.2, 6.7, 7.3, 5.6, 8.3, 7.7, 8.2, 6.0, 6.8, 6.9, 7.3, 7.0, 7.1, 6.6, 7.4

Solution : We set up the following null hypothesis $H_0 : \theta = 6.9$ to be tested against the alternative hypothesis $H_1 : \theta \neq 6.9$.

Let us put '+' for values in the series greater than 6.9 for values in the series less than 6.9 and 0 for values in the series equal to 6.9 Thus we get

-,-,+,-,-,+,-,+,-,-,0,-,+,+,-,+

Thus, number of positive sign = 10, number of negative sign = 9 and so $n = 19$.

x_i	$d_i = x_i - 6.9$	$ d_i $	R_i
6.3	-0.6	0.6	9.5
5.8	-1.1	1.1	14
7.7	0.8	0.8	11.5
8.5	1.7	1.7	18
5.2	-1.7	1.7	19
6.7	-0.2	2	3.5
7.3	0.4	0.4	6.5
5.6	-1.3	1.3	15.5
8.3	1.4	1.4	17
7.7	0.8	0.8	11.5
8.2	1.3	1.3	15.5
6.0	-0.9	0.9	13
6.8	-0.1	0.1	1.5
6.3	-0.6	0.6	9.4
6.9	0	0	ignored
7.3	0.4	0.4	6.5
7.0	0.1	0.1	1.5
7.1	0.2	0.2	3.5
6.6	-0.3	0.3	5
7.4	0.5	0.5	8

Thus, the sum of ranks with positive values of d_i 's is $T^+ = 99.5$. Similarly, the sum of ranks with negative values of d_i 's $T^- = 90.5$. Thus, $T = \min(T^+, T^-) = \min(99.5, 90.5) = 90.5$.

Here, since one of the observation is ignored so we have $n = 19$ and thus for a two sided alternative we find the table value of $T_{\alpha/2} = 46$ for $n = 19$ for $\alpha = 0.05$ level of significance. Since $T_n < T$ so we accept the null hypothesis and conclude that $H_0 : \theta = 6.9$ is true.

Example 3: A medical representative visited 12 doctors in a town. In order to meet the doctor he had to wait for 25, 10, 15, 20, 17, 11, 30, 27, 36, 40, 5 and 26 minutes respectively. However, the senior sales representative earlier claimed that the doctor kept him waiting for more than 20 minutes on an average. Using Wilcoxon Signed rank test verify the claim made by the senior sales representative at 5% level of significance.

Solution : Here the null hypothesis is

$$H_0 : \mu = 20 \text{ minutes}$$

tested against the alternative hypothesis

$$H_1 : \mu > 20 \text{ minutes.}$$

To calculate the test statistic we construct the following table:

x_i	$x_i - 20$	$d_i = x_i - 20 $	Ranks
25	5	5	2.5
10	10	10	7.5
15	-5	5	2.5
20	0	0	ignored
17	-3	3	1
11	-9	9	6
30	10	10	7.5
27	7	7	5
36	16	16	10
40	20	20	11
5	-15	15	9
26	6	6	4

Here, the sum of positive ranks is equal to $T^+ = 40$ and the sum of negative ranks is equal to $T^- = 26$. Also, $T = \min(T^+, T^-) = 26$. The effective sample size is $n = 11$.

Since the alternative sample size is of the form $H_1 : m > m_0$ so the test statistic is T^- . Here $T^- = 26 > T_{0.05} = 14$. Thus H_0 is accepted and so it can be concluded that the average waiting time of the sales representative for the doctor is 20 minutes.