

NON-PARAMETRIC METHODS

Dr. Noah Muthai

(1)

Introduction

- In statistical inference (theory of estimation and hypothesis testing) we assume that the observations come from a distribution for which the exact form is known even though the values of some parameters are unknown e.g.

- (i) We might assume that the observations form a random sample from a poisson distribution for which the mean is unknown i.e. $x_1, x_2, \dots, x_n \sim \text{Poisson}(\lambda)$
- (ii) It might be assumed that the observations come from the Normal Distribution for which the means and variances are unknown
- In the above cases we have assumed that the observations come from certain parametric family of distributions and statistical inference must be made about the values of the parameter defining the family
- In many of the problems to be discussed i.e. non-parametric methods we shall not assume that the available observations come from a particular parametric family of distributions.

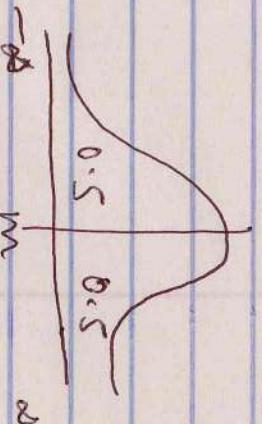
Examples

- 1) We may simply assume that the observations form a random sample from a continuous distribution without specifying the form of this distribution and further we might investigate the possibility that it is a normal distribution.
- 2) We may be interested in making inference about the value of the median of the distribution from which the sample was taken and we may assume that this is a continuous distribution.

- 3) We might be interested that two independent samples come from same distribution and we might only assume that both distributions are continuous.

Definition: Non-parametric methods: Statistical analysis which does not depend upon the knowledge of the distribution and parameters of the population are called non-parametric or distribution free methods.

The sign test

- It is the non-parametric alternative to t-test for hypothesis testing. It is so named because it counts +ve and -ve to test hypothesis.
 - This test is performed when we wish to analyse two sets of data that were gathered independently from each other.
 - To describe the test we assume that the pdf is continuous density and assume that the median denoted by m is uniquely defined by
- $$\int_m^{\infty} f(x) dx = \frac{1}{2}$$
- Note: 
- $$\int_{-\infty}^m f(x) dx = 1 - \int_m^{\infty} f(x) dx = \frac{1}{2}$$
- Let x_1, x_2, \dots, x_n be a random sample of X and consider testing the hypothesis $H_0: M = M_0$

(3)

$$\text{Vs } H_1: m > m_0 \text{ - From the definition of median it follows that when } H_0 \text{ is true then } P(X \geq m_0) = P(X - m_0 \geq 0) = \frac{1}{2} \text{ and therefore } P(X_i - m_0 \geq 0) = \begin{cases} \frac{1}{2} & i=1,2,\dots,n \\ 1 & \text{if } x_i - m_0 > 0 \end{cases}$$

$$\text{Let } Z_i = \begin{cases} 1 & \text{if } x_i - m_0 < 0 \\ 0 & \text{if } x_i - m_0 \geq 0 \end{cases}$$

- We note that Z_i has a Bernoulli distribution with $p = \frac{1}{2}$ when H_0 is true. Since the Z_i 's are independent the sum $U = \sum^n Z_i$ will be a Binomial random variable corresponding to n independent Bernoulli trials for which $p = \frac{1}{2}$ when H_0 is true.

- Under H_0 most x_i will tend to be larger than m_0 and the variance will tend to exceed the value to be expected when H_0 is true.
- At a result the right tail of the Binomial distribution should be chosen as the critical region of the test when H_0 is true i.e

$$Z_i = \begin{cases} 1 & \text{with } p = \frac{1}{2} \\ 0 & \text{when } p = \frac{1}{2} \end{cases}$$

Recall $E[Z_i] = \sum x \cdot f(x)$ for discrete dist

$$\Rightarrow E[Z_i] = \sum Z_i \cdot f(x)$$

$$= 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2}$$

$$= \frac{1}{2}$$

$$E[X_i^2] = \sum x_i^2 \cdot f(x)$$

(4)

$$= 0 \cdot \frac{1}{2} + 1^2 \cdot \frac{1}{2} = \frac{1}{2}$$

Now $u = z_1 + z_2 + \dots + z_n$ i.e. $u = \sum z_i$

$$\begin{aligned} E[u] &= E[z_1] + E[z_2] + \dots + E[z_n] \\ &= \frac{1}{2} + \frac{1}{2} + \dots + \frac{1}{2} \end{aligned}$$

$$= n \frac{1}{2}$$

$$\text{Var}(u) = \text{Var}(z_1) + \text{Var}(z_2) + \dots + \text{Var}(z_n)$$

$$\text{But } \text{Var}(z_i) = f[z_i^2] - (E[z_i])^2$$

$$\text{and } f[z_i^2] = \frac{1}{2},$$

$$\begin{aligned} &= \frac{1}{2} - (\frac{1}{2})^2 \\ &= \frac{1}{2} - \frac{1}{4} = \frac{1}{4} \end{aligned}$$

$$\text{Therefore } \text{Var}(u) = \frac{1}{4} + \frac{1}{4} + \dots + \frac{1}{4}$$

$$= n \frac{1}{4}$$

$$\text{Thus } E[u] = n \frac{1}{2} \text{ and } \text{Var}(u) = n \frac{1}{4}$$

When θ_0 is true

- For very small values of n , it is necessary to calculate the right probabilities until a test probability of approximately α has been obtained to get the exact critical region for the test.
- For $p = \frac{1}{2}$ the Binomial distribution is approximated well by the Normal distribution for fairly small values of n .
- Therefore it is usually sufficient to use the Normal approximation for that purpose.

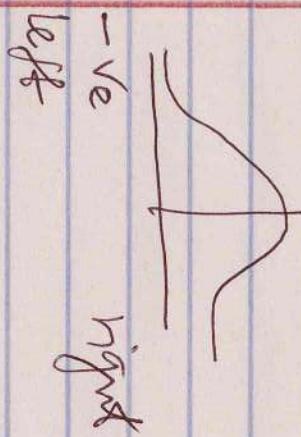
(5)

We use the standard normal variables.

$$Z = \frac{u - \mu_0 - \eta_{12}}{\sqrt{n/4}} \quad \text{for a right fail critical region } (>)$$

$$\text{and } Z = \frac{u + \mu_0 - \eta_{12}}{\sqrt{n/4}} \quad \text{for left fail}$$

$$\frac{-\sqrt{n/4}}{\sqrt{n/4}} \quad \text{critical region } (<)$$



Example

1) The following data was obtained from testing the breaking strength of ceramic tiles manufactured by a new cheaper process

20, 42, 18, 24, 22, 35, 19, 18, 26, 20, 24

32, 28, 20, 24

Suppose that experience with the old process provided a median of 25, test the hypothesis $H_0: \mu = 25$ vs $H_1: \mu < 25$

Solution

In testing $H_0: \mu = 25$ vs $H_1: \mu < 25$ we use the test statistic

$$Z = \frac{u + \mu_0 - \eta_{12}}{\sqrt{n/4}} \quad \text{and reject } H_0 \text{ for small values of } Z \text{ where}$$

(6)

$$u = \sum_{i=1}^n z_i \text{ and } z_i = x_i - m_0, \quad z_2 = x_2 - m_0 \dots$$

so we subtract 25 from each of the values to get
 $-5, 17, -7, -4, -3, 40, -6, -7, 1, -5, -4,$
 $7, -3, -5, 1$

The corresponding Z -values are:

$$0, 4, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1$$

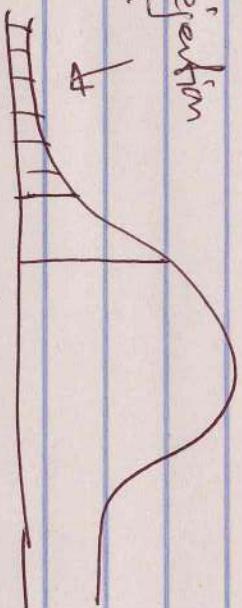
$$\text{thus } u = \sum_{i=1}^n z_i = 4$$

u is Binomial with $n = 15$, and $p = \frac{1}{2}$. Thus when H_0 is true

$$Z = \frac{u + 1/2 - 15/2}{\sqrt{15/4}} = \frac{-3}{\sqrt{3.75}} = -1.549$$

$$Z = -1.55$$

Rejection region



$$\alpha = 0.05$$

from Z -tables we find $P(Z \leq -1.65) = 0.05$. Since $-1.55 > -1.65$ we do not reject H_0 and conclude the median equals 25.

(7)

Note:

(c) Calculation by means of Binomial distribution
With $n = 45$ and $p = \frac{1}{2}$ will give

$$\begin{aligned} f(0) &= p(u=0) = 0.0003 & P(X=x) = \\ f(1) &= p(u=1) = 0.005 & \binom{2}{x} p^x \bar{p}^{n-x} \\ f(2) &= p(u=2) = 0.032 & \binom{15}{5} \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{15-x} \\ f(3) &= p(u=3) = 0.139 \\ f(4) &= p(u=4) = 0.0417 \end{aligned}$$

$$\text{Then } P(u \leq 4) = f(0) + f(1) + f(2) + f(3) \\ = 0.0593$$

which demonstrates that the Normal approx with $\frac{1}{2}$ correction factor is excellent

(d) Instead of getting the values of Z_i in the above case we can record the size of the numbers $X_i - 2.5$ rather than the Z -values in which case we would get

$- + - - + - + - - + - -$
then the total number of the signs gives the value of u . If u is far from reason that the test is called the sign test.

Example 2

- 1) The following data in tonnes are the amounts of Sulphur oxides emitted by a large industrial plant in 40 days

8

17	15	20	29	19	18	22	25	27	9
24	20	17	6	24	14	15	23	24	26
19	23	28	19	16	22	24	17	20	13
19	10	23	18	31	13	20	17	24	14

Test the null hypothesis $H_0: \mu = 21.5$ vs $H_1: \mu < 21.5$ at $\alpha = 0.05$

Solution

$$\text{Let } Z_i = \begin{cases} 1 & \text{if } x_i - 21.5 > 0 \\ 0 & \text{if } x_i - 21.5 \leq 0 \end{cases}$$

The values of Z are

0	0	0	1	0	0	1	1	1	0
1	0	0	0	1	0	0	1	1	1
0	1	1	0	0	1	1	0	0	0
0	0	1	0	1	0	0	0	1	0

$$u = \sum Z_i = 16$$

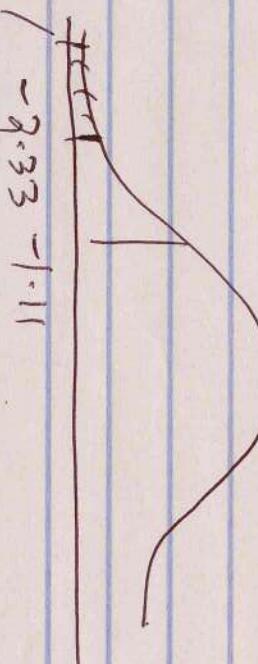
$$\begin{array}{ccccccccc} - & - & - & + & - & - & + & + & + & - \\ + & - & - & - & + & - & - & + & + & + \\ - & + & + & - & - & + & + & + & + & - \end{array}$$

②

$$\bar{z} = u + \frac{1/2 - n/2}{\sqrt{n/4}} = \frac{16 + 1/2 - 40/2}{\sqrt{40/4}} = -4.11$$

from the normal tables we have

$$P(Z \leq -2.33) = 0.01(\alpha)$$



Rejection Region

Therefore we do not reject the since $-1.11 > -2.33$ i.e. -1.11 is not in the Rejection Region

(B)

(1)

The signed-rank test (Wilcoxon signed-rank test)

- As we have seen the sign test is very easy to perform but since we utilize only the signs of the differences between the observations (M_0 (mean) in one sample case, it tends to be wasteful of information.
- An alternative non-parametric test is the Wilcoxon signed rank test is less wasteful in that it takes into account the magnitude of the difference.
- In this test we rank the difference without regard to the signs assigning rank one to the smallest difference in magnitude in absolute value rank two to the second smallest differences in magnitude absolute value and rank n to the largest difference in absolute value.
- Zero differences are ignored. If the absolute value of two or more differences are the same we assign each one the mean of the ranks, which they jointly occupy
- The signed rank test is based on T^+ the sum of the ranks assigned to the +ve difference, T^- the sum of ranks assigned to the -ve difference.
Therefore $\bar{T} = \min(T^+, T^-)$
- We note that T^+ and T^- takes the values in the interval 0 to $\frac{n(n+1)}{2}$.
- T^- and T^+ are symmetric about $\frac{n(n+1)}{4}$.
Assignment: Find out how.
- The critical value regions for various tests is given below:

(2)

Alternative hypothesis

$H_0: \mu = \mu_0$	critical region - Region
$H_1: \mu > \mu_0$	$T^- \leq T_{\alpha}$
$H_1: \mu < \mu_0$	$T^+ \leq T_{2\alpha}$

- Exampole
The following are 15 measurements of the octane rating of a certain kind of gasoline

97.5, 95.2, 97.3, 96.0, 96.8, 100.3, 97.4,
97.4, 95.3, 93.2, 99.1, 96.1, 97.6, 98.2,
98.5, 94.9

Find the signed-rank test at $\alpha=0.05$ level of significance to determine whether or not the mean octane rating of the given kind of gasoline is 98.5.

Solution

We wish to test $H_0: \mu = 98.5$ vs $H_1: \mu \neq 98.5$. Therefore we reject H_0 if $T \leq T_{0.05}$ for the appropriate value of n . We subtract 98.5 from each value and rank the absolute difference i.e (in the next page)

$$T^+ = 2 + 8 = 10 \quad (\text{sum of ranks assigned to the +ve differences})$$

$$T^- = 1 + 3 + 4 + 5 + 6 + 7 + 10 + 11 + 12 + 13 + 14 = 95$$

$$T = \min(T^+, T^-) = \min(10, 95)$$

(continue in next page)

(3)

Measurements	$x_i - \mu$	Rank
97.5	-1	4
95.2	-3.3	12
97.3	-1.2	6
96.0	-2.5	10
96.8	-1.7	7
100.3	(41.8)	2 + different
97.4	-1.1	5
95.3	-3.2	11
93.2	-5.3	14
99.1	(+0.6)	2 + difference
96.1	-2.4	9
97.6	-0.9	3
98.0.2	-0.3	1
98.5	0.0	— disregard - zero
94.9	-3.6	13

Note: — Rank absolute value starting from the smallest difference. • — Disregard zero difference

$n = 14$, $T_{0.05} = 21$ (from critical regions for Wilcoxon - sign rank test)

Since $10 < 21$, we reject H_0 .

Note: $n = 14$, disregard one value. (zero diff)
Rejection criterion is determined by the alternative hypothesis.

④

Signed-Rank test for paired data

- We have the observations $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Let μ_y be the mean of y_i 's and μ_x be the mean of x_i 's.
- We want to test $H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$ or $H_1: \mu_1 > \mu_2$ or $H_1: \mu_1 < \mu_2$. We use the table below.

Alternative hypothesis

Reject H_0 if

$$\begin{aligned} \mu_1 &\neq \mu_2 & T^- &\leq T_{\alpha} \\ \mu_1 &> \mu_2 & T^+ &\leq T_{\alpha} \\ \mu_1 &< \mu_2 & T^+ &\leq T_{2\alpha} \end{aligned}$$

For $n \geq 15$ it is considered reasonable to assume the distribution of T^+ is approximately normal. To perform the sign-rank test based on this assumption we need the following results

Theorem: The mean and variance of T^+ are

$$E(T^+) = n(n+1)/4$$

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

Proof: If the sample of size n , the ranks to be assigned are $1, 2, 3, \dots, n$, for each rank the probability it will be assigned a +ve or -ve difference are both $1/2$ when the is true

$$R_i = \begin{cases} 1 & \text{if } x_i > y_i \\ 0 & \text{if } x_i \leq y_i \end{cases} \quad \text{for } i=1, 2, \dots, n$$

(5)

Thus $T^t = 1 \cdot Z_1 + 2 \cdot Z_2 + 3 \cdot Z_3 + \dots + n \cdot Z_n$
 Z_1, Z_2, \dots, Z_n are independent random variables
 having a Bernoulli distribution with $P = 1/2$ since
 $P(Z_i) = 1/2$ and $\text{Var}(Z_i) = 1/4$ for $i = 1, 2, \dots, n$
 Thus $E(T^t) = E(Z_i)$

$$= E(1 \cdot Z_1 + 2 \cdot Z_2 + \dots + n \cdot Z_n)$$

$$= 1 \cdot E(Z_1) + 2 \cdot E(Z_2) + \dots + n \cdot E(Z_n)$$

$$= 1 \cdot 1/2 + 2 \cdot 1/2 + \dots + n \cdot 1/2$$

$$= \frac{1 + 2 + \dots + n}{2} = \frac{n(n+1)}{4}$$

$$\text{Similarly } \frac{1 + 2 + \dots + n}{2} = \frac{n(n+1)}{2}$$

$$\Rightarrow \frac{n(n+1)}{2} \cdot 1/2 = \frac{n(n+1)}{4}$$

Likewise

$$\text{Var}(T^t) = \text{Var}(1 \cdot Z_1 + 2 \cdot Z_2 + \dots + n \cdot Z_n)$$

$$= 1^2 \text{Var}(Z_1) + 2^2 \text{Var}(Z_2) + \dots + n^2 \text{Var}(Z_n)$$

$$= 1^2 \cdot 1/4 + 2^2 \cdot 1/4 + \dots + n^2 \cdot 1/4$$

$$= \frac{1^2 + 2^2 + \dots + n^2}{4}$$

$$1^2 + 2^2 + \dots + n^2 = \frac{n(n+1)(2n+1)}{6}$$

$$\Rightarrow \frac{n(n+1)(2n+1)}{6} \cdot \frac{1/4}{6}$$

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

⑥

Example
 The following are the weights in pounds before and after 16 persons who started on a certain weight & reducing diet for 4 weeks:

Before	After	Before	After
147.0	137.9	166.8	158.5
183.5	178.2	131.9	130.4
232.1	219.0	150.3	149.3
161.6	163.8	197.2	189.1
197.5	193.5	159.8	159.1
206.3	207.4	121.7	123.2
177.0	180.6		
215.4	203.2		
147.7	149.0		
208.1	195.4		

Use the sign-rank test to test at $\alpha = 0.05$ level of significance whether the weight reducing diet is effective.

Solution

We are testing $H_0: \mu_1 = \mu_2$ vs. $H_1: \mu_1 > \mu_2$

We reject H_0 if $Z_{\text{calc}} \geq Z_{0.05} = 1.645$

$$\text{where } Z = \frac{T^+ - E(T^+)}{\sqrt{\text{Var}(T^+)}}$$

$$\sqrt{\text{Var}(T^+)}$$

The differences between the respective pairs are given on the next page

(7)

pair	diff	rank	
1	9.9	13	+
2	7.3	10	+
3	13.1	16	+
4	-2.2	5	
5	4.0	8	+
6	4.9	7	
7	-3.6	7	
8	12.2	14	+
9	-1.3	3	
10	12.7	15	+
11	8.3	12	+
12	-2.5	6	
13	1.0	2	+
14	8.1	11	+
15	0.7	1	+
16	-1.5	4	

N.B.: Ranker in absenk values

$$n = 16, \bar{T}t = 13 + 10 + 16 + 8 + 9 + 14 + 15 + \\ = 111$$

$$E(Tt) = \frac{n(n+1)}{4} = \frac{16(17)}{4} = 68$$

$$\text{Var}(Tt) = \frac{n(n+1)(2n+1)}{24} = \frac{16(17)(33)}{24}$$

$$Z = \frac{374}{\sqrt{374}} = 9.22$$

⑧

$$\chi^2_d = 20.05 = 1.645$$

$2.22 > 1.645$, thus we reject H_0 .

Alternatively: Reject H_0 if $T^- \leq T_{2d}$

$$T^- = 25 (5+7+3+6+4)$$

$$T_{0.10} = 35$$

$T^- = 25 < T_{0.10} = 35$, thus reject H_0 .

(c)

①

- Rank-Sum test - U-test
- The U-test is used for two samples. The other names are:
- Wilcoxon test
 - Mann-Whitney test
- We are going to test the hypothesis that we are sampling identical continuous populations against the alternative that the two populations have unequal means.
- To illustrate the procedure, suppose we want to compare the kinds of emergency flares on the basis of the following burning times.

Brand A	14.9	11.3	13.2	16.6	17.0	14.1	15.4	13.0	16.9	
Brand B	15.2	19.8	14.7	18.3	16.2	21.2	18.9	13.2	15.3	19.4

Step 1

- Arrange the values in ascending order
- Rank them jointly

Brand A	11.3	13.0	13.2	14.1	14.9	15.4	16.6	16.9	17.0
	④	③	④	⑤	⑦	⑩	⑫	⑬	⑮

Brand B	12.2	14.7	15.2	15.3	16.2	18.3	18.9	19.4	19.8	21.2
	②	⑥	⑧	⑨	⑪	⑯	⑮	⑰	⑯	⑰

→ The values of brand A occupy ranks 1, 3, 4, 5, 7, 10, 12, 13, and 15, and B occupy 2, 6, 8, 1, 9, 11, 15, 16, 17, 18 and 19.

→ For first we assign the mean of the ranks while they 'jointly' carry

- If there is an appropriate difference between the means of the two observations most of the lower ranks are likely to go to the 1 values of one sample while most of the higher ranks will go to the values of the other sample

- The test is based on the values of W_1 where is the sum of the ranks of the first sample or W_2 which is the sum of the ranks of the second sample.

- It doesn't matter whether we use W_1 or W_2 for if there are n_1 values in the first sample and n_2 values in the second sample $W_1 + W_2$ is the sum of the first n_1 true integers i.e. $W_1 + W_2 = \frac{(n_1 + n_2)(n_1 + n_2 + 1)}{2}$

- For any pair of values of W_1 and W_2 . Therefore, the test based on W_1 and W_2 are equivalent. In actual practice we seldom use the test based on the statistics W_1 and W_2 instead we use the related statistics

$$U_1 = W_1 - \frac{n_1(n_1 + 1)}{2}$$

$$U_2 = W_2 - \frac{n_2^2(n_2 + 1)}{2}$$

$$\text{or } U = \min(U_1, U_2)$$

- The resulting tests are all equivalent to the one based on n_1 and n_2 but have the advantage that they label themselves more readily to the tables of

(3)

critical values i.e

$$u_1 + u_2 = n_1 - \frac{n_1(n_1+1)}{2} + n_2 - \frac{n_2(n_2+1)}{2}$$

$$= \frac{(n_1+n_2)(n_1+n_2+1)}{2} - \frac{n_1(n_1+1)}{2} -$$

$$= \frac{n_2(n_2+1)}{2}$$

$$= n_1^2 + n_1 n_2 + n_1 + n_2 n_1 + n_2 + n_2 - n_1^2 - n_1 - n_2^2 - n_2$$

$$= 2(u_1 + u_2) = 2n_1 n_2$$

- $\Rightarrow u_1 + u_2 = n_1 n_2$ i.e the sum of $u_1 + u_2$ is always $n_1 n_2$ and both of these random variables take on the same range of the values from 0 to $n_1 n_2$.
- We use the following table for testing $H_0: \mu_1 = \mu_2$ vs the various alternatives.

Alternative hypothesis

Reject H₀ if

$$\begin{aligned} \mu_1 &\neq \mu_2 & u_{\alpha/2} &\leq u_d \\ \mu_1 &> \mu_2 & u_d &< u_{2\alpha} \\ \mu_1 &< \mu_2 & u_d &\leq u_{2\alpha} \end{aligned}$$

where the level of significance is α for each test

- The critical values of u are such that u_d is the largest value for which $u \leq u_d$ does not

④

exceed α .

Example
1) Two brands of flares have the following burning times

	Brand A	Brand B						
	14.9	15.4						
	15.3	15.2	19.8	14.7	18.3	16.2	21.2	
	13.2	12.2	14.7	14.7	15.3	15.4		
	16.6	16.9	18.3	18.3	19.4	19.4		
	17.0	17.0	19.8	19.8	21.2	21.2		
	17.0	17.0	19.8	19.8	21.2	21.2		

Test at $\alpha = 0.05$ level of significance whether the two samples come from identical continuous populations or whether the mean burning times of B flares lie above A.

Solution

We want to test $H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 < \mu_2$
Arrange the values in ascending order (rank them jointly)

Brand A	11.3	13.0	13.2	14.1	14.9	15.4	16.6	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0
	①	③	④	⑤	⑦	⑩	⑫	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯
	16.9	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0	17.0
	⑬	⑭	⑮	⑯	⑰	⑱	⑲	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳
Brand B	13.2	14.7	15.2	15.3	16.2	18.3	18.3	19.4	19.8	21.2	21.2	21.2	21.2	21.2	21.2	21.2	21.2	21.2
	②	⑥	⑧	⑨	⑪	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯	⑯
	13.2	14.7	15.2	15.3	16.2	18.3	18.3	19.4	19.8	21.2	21.2	21.2	21.2	21.2	21.2	21.2	21.2	21.2
	⑰	⑱	⑲	⑲	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳	⑳

We reject H_0 if $u_1 \leq u_{\alpha}$

(5)

thus $w_1 = 1 + 3 + 4 + 5 + 7 + 10 + 12 + 13 + 14$
 $= 69$

$$\Rightarrow u_1 = w_1 = \frac{w_1(n_1+n_2)}{2} = \frac{69 - 9(10)}{2} = 24$$

$$\text{for } n_1 = 9, n_2 = 10$$

$$u_{2d} = u_{0.10} = 24$$

thus we conclude since $24 \leq 24$ we reject the null hypothesis that on average brand A flakes as a mean value time that is less than that of brand B.

- When both n_1 and n_2 are greater than 8, it is considered reasonable to assume that the distribution of u_{1d} and u_{2d} u_1 and u_2 can be approximated closely by the normal distribution. To perform the given rank sum test on the basis of this assumption, we need the following result

Theorem: Under the null hypothesis, the mean and variance of u_{1d} and u_{2d} u_1 and u_2 are

$$1) E[u_1] = E[u_2] = \frac{n_1 n_2}{2} \text{ and}$$

$$2) \text{Var}(u_1) = \text{Var}(u_2) = \frac{n_1 n_2 (n_1 + n_2 + 1)}{12}$$

⑥

Example
 The following are the number of minutes it took
 random samples of 15 men and 12 women to
 complete a written test given for their renewal
 of their driving licences

Men: 9.9, 7.4, 8.9, 9.1, 7.7, 9.7, 11.8, 9.2,
 10.0, 10.2, 9.5, 10.8, 8.0, 11.0, 7.5,

Women: 8.6, 10.9, 9.8, 10.7, 9.4, 10.3, 7.3, 11.5,
 7.6, 9.3, 8.8, 9.6

Use the t -test based on the normal approxi-
 mation to test $H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$
 Where μ_1 and μ_2 are the average amount of
 time it takes men and women to complete
 the test respectively

Solution

Re-arrange the data in ascending order.

7.3, 7.4, 7.5, 7.6, 7.7, 8.0, 8.6, 8.8, 8.9
9.0, 9.1, 9.2, 9.3, 9.4, 9.5, 9.6, 9.7, 9.8, 9.9
10.0, 10.2, 10.3, 10.7, 10.8, 10.9, 11.0, 11.5, 11.8

The underlined values are for men

$$n_1 = 15$$

$$n_2 = 12$$

⑦

$$\begin{aligned} \text{Ranks for men} \Rightarrow n_1 &= 2 + 3 + 5 + 6 + 9 + \\ 10 + 11 + 14 + 16 + 18 + 19 + 20 + 23 + 25 + \\ 27 &= 208 \end{aligned}$$

$$\begin{aligned} \text{Ranks for women} \Rightarrow n_2 &= 1 + 4 + 7 + 8 + 12 + \\ 13 + 15 + 17 + 24 + 22 + 24 + \\ &= 170 \end{aligned}$$

$$u_1 = n_1 - \frac{n_1(n_1+1)}{2}$$

$$= 208 - \frac{15(16)}{2} = 88$$

$$u_2 = n_2 - \frac{n_2(n_2+1)}{2}$$

$$= 170 - \frac{12(13)}{2} = 92$$

$$\Rightarrow u = \min(u_1, u_2) = \min(88, 92)$$

$$\text{Thus } E[u] = \frac{n_1 n_2}{2} = \frac{15 * 12}{2} = 90$$

$$\text{Var}(u) = \frac{n_1 n_2 (n_1 + n_2 + 1)}{12}$$

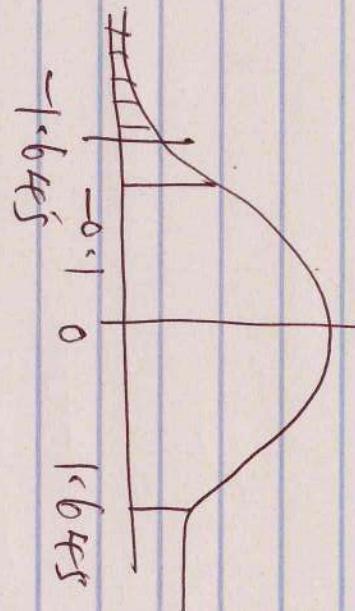
$$= \frac{15 * 12 (15 + 12 + 1)}{12}$$

$$= 420$$

⑧

$$Z = \frac{U - E[U]}{\sqrt{\text{Var}(U)}} = \frac{88 - 90}{\sqrt{420}} \approx -0.1$$

$$Z_\alpha = 20.05 = 1.645 / \text{H}_0: \mu_1 = \mu_2 \\ \text{H}_1: \mu_1 \neq \mu_2 \Rightarrow \text{2-sided}$$



since $-1.645 < -0.1$ we do no reject the null hypothesis (-0.1 is not in the rejection region)

Exercises

- 1) The following are the amount of time in mins which took a random sample of 20 technicians to perform a task

18.1, 20.3, 18.3, 15.6, 22.5, 16.8, 17.6, 16.9, 18.2, 17.0, 19.3, 16.5, 19.5, 18.6, 20.0, 18.8, 19.1, 17.6, 18.5, 18.0

Use the sign test at $\alpha = 0.05$ to test the null hypothesis that the measurement constitute a random sample from a continuous pdf with mean $\mu = 19.4$ mins against the 2 sided alternative $\mu \neq 19.4$ mins. Use the Binomial test.

- 2) Rewrite problem 1 using the normal approximation to the Binomial. Use rank-sign-rank test also.

⑤

④

Rank sum test : The H-test (Kruskal-Wallis test)

- The H-test / Kruskal-Wallis test is a generalization of the rank-sum test to the case where we test the null hypothesis that K-samples come from identical continuous distributions.
- As in the U-test the data are ranked jointly from low to high as though they came from one sample.
- Let R_i be the sum of the ranks of the values of the i th sample.
- We base the test on the statistic

$$H = \frac{12}{n(n+1)} \sum_{i=1}^k n_i \left(\frac{R_i}{n_i} - \frac{n+1}{2} \right)^2 - 0$$

Therefore the H-statistic is proportional to the weighted mean of the squared difference $(R_i/n_i - \frac{n+1}{2})^2$ where R_i/n_i is the mean

rank for sample i and $(n+1)/2$ is the mean rank of all the data. It follows that the null hypothesis must be rejected for large values of H .

- Expression 1 can be simplified as

$$(R_i/n_i - \frac{n+1}{2})^2 = \frac{R_i^2}{n_i^2} - 2 \left(\frac{R_i}{n_i} \right) \left(\frac{n+1}{2} \right) + \left(\frac{n+1}{2} \right)^2$$

$$= \frac{R_i^2}{n_i^2} - \frac{R_i}{n_i} \left(n+1 \right) + \frac{(n+1)^2}{4}$$

$$= n_i \left(\frac{R_i}{n_i} - \frac{n+1}{2} \right)^2 = \frac{R_i^2}{n_i} - R_i (n+1) + \frac{n_i(n+1)^2}{4}$$

(2)

$$= \sum_{i=1}^K n_i \left(\frac{R_i}{n_i} - \frac{n+1}{2} \right)^2$$

$$= \sum_{i=1}^K R_i^2 / n_i - (n+1) \sum_{i=1}^K R_i + \frac{(n+1)^2}{4} \sum_{i=1}^K n_i$$

But

$$\sum_{i=1}^K n_i = n$$

$$= \sum_{i=1}^K R_i^2 / n_i - (n+1) \cdot n \left(\frac{n+1}{2} \right)^2 + \frac{(n+1)^2}{4} n - \textcircled{*}$$

$$= \sum_{i=1}^K R_i^2 / n_i - n \left(\frac{n+1}{4} \right)^2$$

Substituting in ⑥ we find that

$$H = \frac{12}{n(n+1)} \left[\sum_{i=1}^K R_i^2 / n_i - \frac{n(n+1)^2}{4} \right]$$

$$H = \frac{12}{n(n+1)} \sum_{i=1}^K R_i^2 / n_i - 3(n+1) - \textcircled{2}$$

- For every small values of K and n_i the test of the null hypothesis may be based on special tables but since the sampling distribution of H depends on the values of n_i it is impossible to tabulate it in compact form hence the test is usually based on the large sample theory that the distribution of H can be approximately closely with a χ^2 distribution with $K-1$ degrees of freedom.

(3)

1)

Example
A sampling of the acidity of rain on 40 randomly selected rainfall was recorded at 3 different locations i.e., 1, 2, 3.

1	2	3
4.45	4.60	4.55
4.02	4.27	4.31
4.13	4.31	4.84
3.51	3.88	4.67
3.81	4.22	4.95
4.18	4.54	4.72
3.95	4.76	4.63
4.07	4.36	4.36
4.29	4.21	4.67

Use the Kruskal-Wallis test at $\alpha = 0.05$ to test whether or not there is a difference in the average acidity of rain in the 3 locations.

$$k = 3$$

$$n = n_1 + n_2 + n_3 = 10 + 10 + 10$$

$$\Rightarrow n = 30$$

$$\sum_i^k n_i = n$$

for each location arrange the values in ascending order.

(4)

1	2	3
3.51	3.88	4.28
3.89	4.21	4.31
3.95	4.22	4.36
4.02	4.27	4.47
4.07	4.31	4.55
4.13	4.36	4.63
4.18	4.49	4.72
4.29	4.54	4.84
4.42	4.60	
4.45	4.76	4.95

- Rank the values jointly, we get

3.51, ①	3.88, ②	3.89, ③	3.95, ④	4.02, ⑤	4.07, ⑥	4.13, ⑦	4.18, ⑧
4.21, ⑨	4.22, ⑩	4.27, ⑪	4.28, ⑫	4.29, ⑬	4.31, ⑭	4.31, ⑮	4.36, ⑯
4.36, ⑰	4.42, ⑱	4.45, ⑲	4.47, ⑳	4.49, ㉑	4.54, ㉒	4.55, ㉓	4.60, ㉔
4.63, ㉕	4.67, ㉖	4.72, ㉗	4.76, ㉘	4.84, ㉙	4.95, ㉚		

Note: $4.31, 4.31$
 ⑯, ⑯ ⑰, ⑰ $14 + 15 = 29/2 = 14.5$

$4.36, 4.36$ $16 + 17 = 33/2 = 16.5$
 ⑮, ⑮ ⑯, ⑯

Ranks for London 1: 1 + 3 + 4 + 5 + 6 + 7 + 8 +
 13 + 18 + 19
 II : 34 + 11 + 14.5 + 2 + 21 + 10 + 22 + 28 + 16.5 + 9

⑤

$$\bar{W} : 13 + 14.5 + 29 + 26 + 12 + 30 + 27 + 25 + 16.5 + 20$$

$$R_1 = 84^1 \\ R_2 = 158 \\ R_3 = 223$$

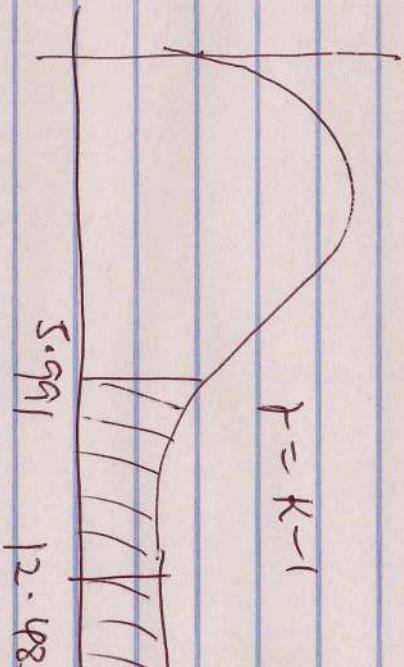
Therefore $H = 1^2 \sum_{i=1}^K \frac{R_i^2}{n_i} - 3(n+1)$

$$= \frac{1^2}{30(31)} \sum_{i=1}^3 \frac{R_i^2}{n_i} - 3(n+1)$$

$$= \frac{1^2}{30(31)} \left(\frac{84^2}{10} + \frac{158^2}{10} + \frac{223^2}{10} \right) - 3(31)$$

$$= \frac{1^2}{930} \left(\frac{84^2}{10} + \frac{158^2}{10} + \frac{223^2}{10} \right) - 93$$

$$\Rightarrow H = 12.4835$$



$$n = 3^0, K = 3, r = 3-1 = 2$$

$$\Rightarrow \chi^2_{0.05}(2) = \chi^2_{0.05}(2) = 5.991$$

$12.485 > 5.991$, we reject H_0 .

(b)

Exercise

- 1) In an example to determine which of the 3 missiles systems is preferable, the propellant burning rate is measured. The data after cooling are given in the table below. Use the Kruskal-Wallis test at $\alpha = 0.05$ to test the hypothesis that the propellant burning rate are the same for the three missile systems.

propellant burning rate

1	2	3
14.0	13.2	18.4
16.7	19.8	19.1
12.8	18.1	17.3
19.8	17.6	19.7
18.9	20.2	18.9
	17.8	18.8
		19.3

- 2) The following data represent the operating time in hours for 3 types of scientific pocket calculators before a recharge is required

Calculator

A	B	C
4.9	5.5	6.4
6.1	5.4	6.8
4.3	6.2	5.6
4.6	5.8	6.5
5.3	5.5	6.5
	5.2	6.3
	4.8	6.6

Use Kruskal-Wallis test at $\alpha = 0.05$ to test the hypothesis that the operating times for all 3 calculators are the same

The run's test

- We look at a test for testing randomness of observed data on the basis of the order of which they were obtained.
- Definition: Run - A run is a succession of identical letters (or other kinds of symbols) which is preceded and followed by different letter or no letter at all.

- Consider the following arrangement of effective (d) and non-effective (n) pieces produced in the given order by a certain machine.

n n n n n d d d n n n n n
 ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗
 n n n n n n d d n n d d d
 ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗
 d n d d nn
 ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗
 ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗
 ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗ ↗

hence 9 runs.

- The total number of runs appearing in an arrangement of this kind is often a good indicator of a possible lack of randomness.
- If there are too few runs we might suspect a definite grouping or clustering or perhaps a trend.
- If there are too many runs we might suspect some sort of repeated alternating pattern.
- In our illustration there seems to be some sort of clustering (few runs).

- If we have n_1 letters of some kind and n_2 letters of another kind, there are:

(2)

$\binom{n_1+n_2}{n_1}$ possible arrangements of these letters

All arrangements are regarded to be equally likely.

Example

1) If we have 2 a's and b's then the no of possible arrangement of these 'letter is

$$\binom{2+2}{2} = \frac{4!}{2!2!} = \frac{4!}{2^2 \cdot 2!} = \frac{4 \times 3 \times 2 \times 1}{2^2 \times 2 \times 1} = 6$$

These arrangements are:

$$\begin{array}{c} \underline{aa} \\ \underline{bb} \\ \hline 2r \end{array} \quad \begin{array}{c} \underline{ab} \\ \underline{ab} \\ \hline 4r \end{array} \quad \begin{array}{c} \underline{abba} \\ \underline{bab} \\ \hline 3r \end{array} \quad \begin{array}{c} \underline{baba} \\ \underline{bab} \\ \hline 4r \end{array} \quad \begin{array}{c} \underline{bbaa} \\ \underline{bab} \\ \hline 2r \end{array}$$

baba

3r

- 2) We have five letters e's which are divided into 3 runs using vertical bar to represent the 5 letters into 3 runs. We find that there are 6 possibilities i.e.

e | e | eee e | e | eee ee | eee | ee

eee | e | e ee | e | eee e | eee | e

(Ans)

(3)

- Let there be n_1 letter of one kind and n_2 letters of another kind.
- Let k be the number of runs formed by these $n_1 + n_2$ letters
- If $u = 2k$ where k is a positive integer then we have k runs of n_1 letters and k runs of the n_2 letters.
- Number of ways of which n letters can form k runs is $\binom{n_1 - 1}{k - 1}$ and number of ways in which n_2 letters can form k runs is $\binom{n_2 - 1}{k - 1}$

$$\rho(n = 2k \text{ runs}) = F(u) = 2 \frac{\binom{n_1 - 1}{k - 1} \binom{n_2 - 1}{k - 1}}{\binom{n_1 + n_2}{n_1}}$$

if $n = 2k + 1$

$$\begin{aligned} F(u) &= \rho(u_1 = k+1, u_2 = k) + \rho(u_1 = k, u_2 = k+1) \\ &= \binom{n_1 - 1}{k} \binom{n_2 - 1}{k-1} + \binom{n_1 - 1}{k-1} \binom{n_2 - 1}{k} \\ &= \binom{n_1 - 1}{k} \binom{n_2 - 1}{k} + \binom{n_1 - 1}{k-1} \binom{n_2 - 1}{k} \\ &\quad \frac{\binom{n_1 + n_2}{n_1}}{\binom{n_1 + n_2}{n_1}} \end{aligned}$$

(4)

When n_1 and n_2 are small the test of randomness is based on u (the number of runs) are usually performed with the use of special tables. We reject the null hypothesis of randomness at level of significance if $u < u_{\alpha/2}$ or

$$u > u_{\alpha/2}$$

Where $u_{\alpha/2}$ is the largest value for which $P(u \leq u_{\alpha/2})$ does not exceed $\alpha/2$ and $u_{\alpha/2}$ is the smallest value for which $P(u > u_{\alpha/2})$ does not exceed $\alpha/2$.

Example

Checking on the palm tree that were planted many years ago along a country road a survey officer obtained the following arrangement of healthy (H) and diseased (S) trees i.e.

H H H H S S S H H H H H H
S S H H S S S H H H H H H

Test at $\alpha = 0.05$ whether this arrangement can be regarded to be random

Solution

We wish to test

H_0 : Arrangement is random

H_1 : Arrangement is not random

(5)

$$\text{No of } H_1's / n_1 = 13$$

$$u_{\alpha/2} = u_{0.05/2} \\ = 6$$

$$\text{No of } H_2's / n_2 = 9$$

We reject H_0 if $u \leq 6$ or $u \geq 7$
from the data $u = 6$, since $6 \leq 6$ we reject
 H_0 and conclude that the arrangement of healthy
and diseased tree is not random.

- If n_1 and n_2 are both greater or equal to
10, it is considered reasonable to assume that
the distribution of u can be closely approximated by
the normal.

- To perform the test on the basis of this assumption
we need the following result:

Theorem: Under the null hypothesis of randomness
the mean and variance of u are

$$E[u] = \frac{2n_1 n_2}{n_1 + n_2} + 1$$

$$\text{Var}(u) = \frac{2n_1 n_2 (2n_1 n_2 - n_1 - n_2)}{(n_1 + n_2)^2}$$

Example

The following series of boys and girls were
selected haphazardly in a co-education school

G B B G G B G B B B G G G G B G B G B B G G B G

Test for randomness at $\alpha=0.01$

solution

We want to test H_0 : Arrangement is random

H_1 : Arrangement is not random

⑥

If u is the number of runs, then

$$E[u] = 2n_1 n_2 - 1$$

$$\text{Var}(u) = \frac{2n_1 n_2 (2n_1 n_2 - n_1 - n_2)}{n_1}$$

We use the test statistic

$$Z = \frac{(u - E[u])}{\sqrt{\text{Var}(u)}}$$

We reject H_0 if $Z \leq -2.005 = -2.575$

or if $Z \geq 2.005 = 2.575$

$$\begin{aligned} \text{No of girls} &= n_1 = 20 \\ \text{No of boys} &= n_2 = 18 \end{aligned}$$

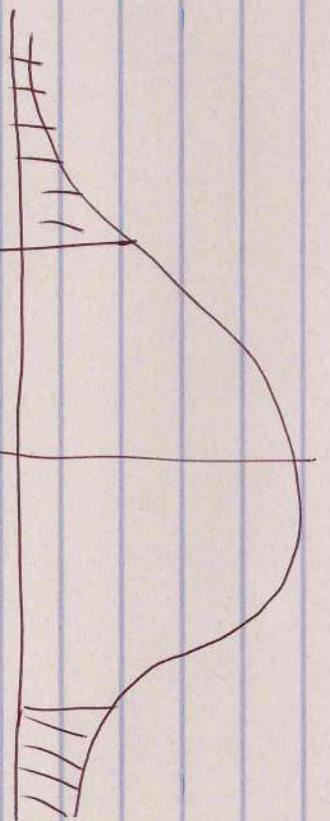
$$u = 25 (\text{runs})$$

$$E[u] = 2 \frac{(20)(18)}{20+18} = 1 = 19.987$$

$$\begin{aligned} \text{Var}(u) &= 2 \frac{(20)(18)}{(20+18)^2} (20+18-1) \\ &= 9.190 \end{aligned}$$

$$Z = \frac{25 - 19.987}{\sqrt{9.190}}$$

⑦



$$\text{Since } -2.57 < u < 2.57$$

$$-2.57 < 1.50 < 2.57$$

We accept H_0 and
We do not reject H_0 and say the arrangement
is random (1.50 is in the acceptable region)

Exercise

1) The following are the amounts of money in dollars
spent by 16 persons at a certain amusement bar

10.15	9.85	13.75	8.63	11.09	15.63
6.65	9.27	8.80	11.45	10.29	9.51
10.0	7.48	9.11			

- (a) Use the sign-test at $\alpha = 0.05$ to test the null hypothesis that on average a person spends 9 dollars at the bar.
- (b) Rewrite the preceding problem using the sign-rank test based on the normal approximation to the distribution of the test statistic.

Rank correlation coefficient

- Since the assumptions underlying the significance test for the correlation coefficient are rather stringent, it is sometimes preferable to use a non-parametric alternative.

- Most popular among the non-parametric measures of association is the rank correlation coefficient also called Spearman's rank correlation coefficient.
- For a given set of data $(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_n, y_n)$,
The rank correlation is obtained by ranking the x_i 's among themselves and the y_i 's both from low to high or high to low, the rank correlation coefficient is given by:

$$\rho_s = 1 - \frac{6 \sum_{i=1}^n \delta_i^2}{n(n^2-1)}$$

Where δ_i is the difference between the ranks assigned to x_i and y_i . Where there are ties in the ranks we assign the tied observations the mean of the ranks which they jointly occupy.

Example
Calculate the Spearman's coefficient rank correlation from the following:

x	73	98	95	107	95	81	79	75
y	147	125	132	137	130	140	139	145

Calculate r .

(2)

X	r_x	y	r_y	δ_i	δ_i^2
73	1	147	8	-7	49
118	8	125	1	7	49
115	7	132	3	4	16
107	6	137	2	2	4
95	5	130	4	3	9
81	4	140	6	-2	4
79	3	139	5	-2	4
75	2	145	7	-5	25

$$\sum \delta_i^2 = 160$$

$$r_s = 1 - \frac{6 \sum_{i=1}^n \delta_i^2}{n(n^2-1)}$$

$$= 1 - \frac{6(160)}{8(63)}$$

$$= 1 - \frac{120}{63}$$

$$= -0.9048$$

- For small values of n , i.e. $n \leq 10$ the test of H_0 of no correlation may be based on t -statistic tables determined from the exact sampling distribution of r_s .

- Most of the time though we use the fact that the distribution of random variables can be closely approximated to the normal distribution

③

Theorem: Under the null hypothesis of no correlation the mean and variance of random variable are:
 $E[r_s] = 0$ and $\text{Var}(r_s) = 1/n - 1$

Note: Strictly speaking the theorem applies when there are no ties, but the result can be used as an approximation unless the number of ties is large

Example

1) For the data in the previous example test at $\alpha = 0.01$ whether the rank correlation coefficient is significant

Solution

We want to test

H_0 : There is no correlation

H_1 : There is correlation

The test statistic is:

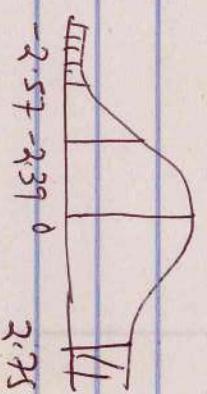
$$T = r_s - E[r_s] = \frac{r_s - 0}{\sqrt{\frac{1}{n-1}}}$$

We reject H_0 if $|z| > 2.75$ or $z > 2.75$

$$\text{i.e. } Z_{df/2} = Z_{0.01/2} = Z_{0.005} = \pm 2.75$$

$$r_s = -0.9048, n = 8$$

$$\Rightarrow T = -0.9048 \sqrt{8-1}$$



$$= -2.3938$$

$Z = -2.3938 > -2.57$ we reject H_0 and conclude that the correlation is not significant

(4)

Note: When there are no ties in ranks, r_s actually equals the correlation coefficient r , calculated for the ranks.

Let r_i and s_i be the ranks of x_i and y_i respectively.

We find that

$$\sum_{i=1}^n r_i = \sum_{i=1}^n s_i = \frac{n(n+1)}{2} (1 + 2 + 3 + \dots + n)$$

$$\sum_{i=1}^n r_i^2 = \sum_{i=1}^n s_i^2 = \frac{n(n+1)(2n+1)}{6} (1^2 + 2^2 + \dots + n^2)$$

Let $d_i = r_i - s_i$

$$d_i^2 = (r_i - s_i)^2 = r_i^2 - 2r_i s_i + s_i^2$$

$$\sum_{i=1}^n d_i^2 = \sum_{i=1}^n r_i^2 - 2 \sum_{i=1}^n r_i s_i + \sum_{i=1}^n s_i^2$$

$$\Rightarrow 2 \sum_{i=1}^n r_i s_i = \sum_{i=1}^n r_i^2 + \sum_{i=1}^n s_i^2 - \sum_{i=1}^n d_i^2$$

$$= \frac{2n(n+1)(2n+1)}{6} - \sum_{i=1}^n d_i^2$$

$$r = \frac{\sum r_i s_i - \frac{1}{n} (\sum r_i)(\sum s_i)}{\sqrt{\sum (r_i^2 - \bar{r}^2) \sum (s_i^2 - \bar{s}^2)}}$$

$$\sqrt{\sum (r_i^2 - \bar{r}^2) \sum (s_i^2 - \bar{s}^2)}$$

(5)

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

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

$$= \frac{n(n+1)(2n+1)}{6} - \frac{1}{n} \left(\frac{n(n+1)}{2} \right)^2 - \frac{1}{2} \sum d_i^2$$

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

Note: $\sum x_i^2 - n\bar{x}^2$, $\bar{x} = \frac{1}{n} \sum x_i$

$$= \frac{1}{n} \left(\frac{n(n+1)}{2} \right)$$

$$n\bar{x}^2 = n \left(\frac{1}{n} \left(\frac{n(n+1)}{2} \right) \right)^2$$

$$= n \left(\frac{n+1}{2} \right)^2$$

$$= 1 - \frac{1}{2} \sum_{i=1}^n d_i^2$$

--- (x)

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

$$\frac{n(n+1)(2n+1)}{6} - \frac{1}{n} \left(\frac{n(n+1)}{2} \right)^2 = \frac{(n+1)}{12} \left[2n(2n+1) - 3n(n+1) \right]$$

(6)

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

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

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

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

Replacing in eqn (5) we find that

$$r = 1 - \frac{1}{12} \sum d_i^2$$

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

$$\sum d_i^2 \times \frac{1}{12} = 6 \sum d_i^2$$

$$n(n^2-1) \quad n(n^2-1)$$

$$r = 1 - 6 \sum d_i^2 = 1$$

$$\frac{n(n^2-1)}{n(n^2-1)}$$

check expansion of

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

$$\frac{1}{6} \left\{ n (2n^2 + 3n + 1) \right\} = \frac{1}{6} \left\{ 2n^3 + 3n^2 + n \right\}$$

⑦

$$\begin{aligned}\frac{1}{n} \left(\frac{n^2 + n}{2} \right)^2 &= \frac{1}{4n} (n^4 + 2n^3 + n^2) \\ &= \frac{1}{4n} (n^3 + 2n^2 - n)\end{aligned}$$

$$\Rightarrow \frac{3}{2n} + 3n^2 + n - \left(\frac{n^3 + 2n^2 + n}{4} \right)$$

$$= \frac{4n^3 + 6n^2 + 2n - 3n^3 - 6n^2 - 3n}{12}$$

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

(5)

(1)

Testing for goodness-of-fit

- Hypothesis testing procedures are designed for problems in which the probability distribution is known and hypothesis involve the parameters of the distribution.
 - Another kind of hypothesis is often encountered.
 - We do not know the underlying distribution of the population and we wish to test the hypothesis that particular distribution will be satisfactory as a population model.
 - We describe a formal goodness-of-fit test procedure based on the chi-square distribution.
 - The test procedure requires a random sample of size n from the population whose probability distribution is unknown.
 - These n observations are arranged in a frequency histogram having k - class intervals.
 - Let O_i be the observed frequency of the i th class.
 - From the hypothesized probability distribution we compute the expected frequency denoted by E_i .
 - The test statistic is $\chi^2_0 = \sum (O_i - E_i)^2 / E_i$
- which has an approximately χ^2_{k-p-1} distribution with $k-p-1$ degrees of freedom where p is the number of parameters of the hypothesized distribution estimated by sample statistics.
- This approximation improves as n increases.
 - We will reject the hypothesis that the distribution of the population is the hypothesized if the calculated value of the test statistic $\chi^2_0 > \chi^2_{\alpha, k-p-1}$

(2)

Example

The number of defects in a printed circuit board is hypothesized to follow a Poisson distribution. A random sample of $n = 60$ printed boards has been collected and the following number of defects observed

No of defects	Observed frequency (O_i)
0	32
1	15
2	9
3	4

The mean of the assumed Poisson distribution in this example is unknown and must be estimated from the sample data. The estimate of the mean is $(3x_1 + (15 \times 1) + (9 \times 2) + (4 \times 3)) / 60 = 0.75$

From the Poisson distribution with parameter 0.75 we may compute P_i , the theoretical hypothesized probability associated with the i th class interval. P_i are found as follows.

$$P_0 = P(X=0) = \frac{e^{-0.75}}{0!} = 0.472$$

$$P_1 = P(X=1) = \frac{e^{-0.75} 0.75^1}{1!} = 0.354$$

$$P_2 = P(X=2) = \frac{e^{-0.75} 0.75^2}{2!} = 0.133$$

$$P_3 = P(X=3) = \frac{e^{-0.75} 0.75^3}{3!} = 0.041$$

(3)

$$N_{\text{def}} \Rightarrow p_4 = \text{Pr}(X \geq 4) = 1 - (p_1 + p_2 + p_3)$$

The E_i are computed by multiplying the sample size $n=60$ times p_i i.e. $E_i = n p_i$

No of defects	p_i	E_i
0	0.472	28.32
1	0.354	21.24
2	0.133	7.98
3 (or more)	0.041	2.46

Since the expected frequency in the last cell is less than 3, we combine the last 2 cells

No of defects	E_i
0	28.32
1	21.24
2 (or more)	10.44

The degree of freedom are $k-p-1 = 3-1-1 = 1$

H_0 : The form of the dist of defect is Poisson

H_1 : The form of the dist of defect is not a Poisson

$$\chi^2 = 0.05$$

$$\chi^2_0 = \sum \frac{(O_i - E_i)^2}{E_i}$$

Reject H_0 if $\chi^2_0 > 0.05, 1 = 3.84$

(F)

$$\chi^2 = \frac{(37 - 28.32)^2}{28.32} + \frac{(15 - 21.24)^2}{21.24} + \frac{(13 - 10.44)^2}{10.44}$$

$$= 2.94$$

Since $\chi^2 = 2.94 < \chi^2_{0.05} = 3.84$ we are unable to reject H_0 that the distribution of defects printed on the circuit board is Poisson.

Contingency table test

- If is used to test whether independence of the categorical variables classified in a table.
- A contingency table is a two-way classification of data i.e. a company has to choose among 3 pension plans. The management wishes to know whether the preference of plans is independent of job classification using $\alpha = 0.05$

- The opinion of a random sample of 500 employees are shown below

Job classification	Pension plan			Total
	1	2	3	
Salary workers	160	140	40	340
Hourly workers	40	60	60	160
Total	200	200	100	500

H_0 : preference is independent of salaried workers versus hourly job classification

H_A : preference is not independent of salaried versus hourly job classification.

⑤

$$\chi^2_{\text{calculated}} = \sum_{i=1}^m \sum_{j=1}^n \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

m = rows

n = cols

Where O_{ij} is the observed frequency in row i column j

E_{ij} is the expected frequency in row i column j

and is given by

$$E_{ij} = \frac{\text{ith row total} \times \text{jth column total}}{\text{Sample size}}$$

$$= \frac{R_i - C_j}{n}$$

The critical value is given by $\chi^2_{\text{crit}} =$

$\chi^2_{(m-1) \times (n-1)}$. Reject the null if $\chi^2_{\text{calc}} > \chi^2_{\text{crit}}$

$$E_{ij} \Rightarrow E_{11} = \frac{340 \times 200}{500} = 136$$

$$\text{then } E_{12} = \frac{340 \times 200}{500} = 136$$

$$E_{13} = \frac{340 \times 100}{500} = 68$$

(6)

$$E_{21} = \frac{160 \times 200}{500} = 64$$

$$E_{22} = \frac{160 \times 200}{500} = 64$$

$$E_{23} = \frac{160 \times 100}{500} = 32$$

Job Classification	Pentin	Tidit	Total
Salaried Workers	160 (136)	140 (36)	340
Handy Workers	40 (64)	60 (64)	160

$$\chi^2_{\text{calc}} = \frac{(160 - 136)^2}{136} + \frac{(40 - 136)^2}{136} + \frac{(40 - 68)^2}{68} +$$

$$\frac{(40 - 64)^2}{64} + \frac{(60 - 64)^2}{64} + \frac{(60 - 32)^2}{32}$$

$$= 49.63$$

$$\chi^2_{\text{crit}} = \chi^2_{(n-1)(p-1)} = \chi^2_{2(0.05)} = 5.991$$

Conclusion: Since $\chi^2 = 49.63 > \chi^2_{(0.05)} = 5.991$
 We reject Ho and conclude preference for pentin
 plans is not dependent on job classification

Exercise

- Grades in statistics and DE courses taken simultaneously were as follows for a group of students.

(T)

Statistics grade	A	B	C	Others	Total
A	25	6	7	13	51
B	17	16	15	6	54
C	18	14	8	10	50
Others	10	8	11	20	49
Total	70	34	51	49	204

Are the grades in statistics and OR related? $\alpha = 0.04$
in reaching your conclusion.

(H)

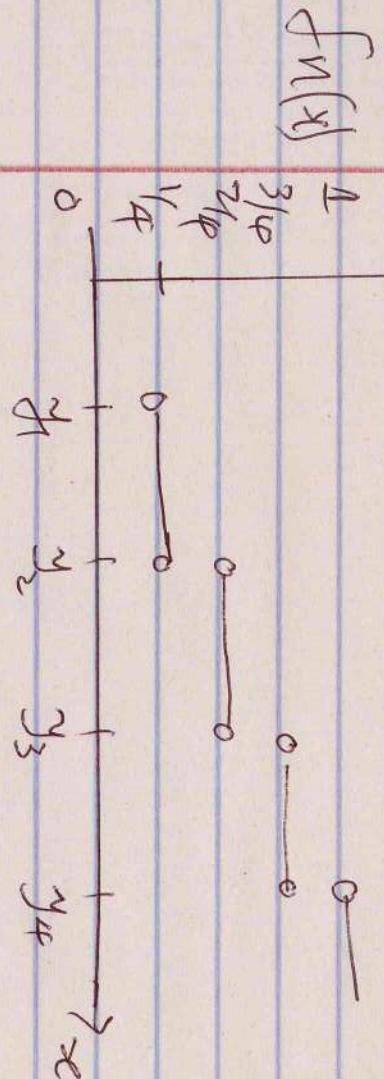
(I)

KD INDIGOOR-SMIRNOV TEST

- Suppose the random variables x_1, x_2, \dots, x_n is from a continuous random sample and let x_1, x_2, \dots, x_n be the observed values of X_1, X_2, \dots, X_n .
- Since the observations come from a continuous distribution there is a probability of zero (0) that any of the observed values x_1, x_2, \dots, x_n will be equal. Thus we shall assume for simplicity that all the n -values are different.
- We consider a function $f_n(x)$ which follows is constant from the values x_1, x_2, \dots, x_n as follows.
 - For each number x ($-\infty < x < \infty$) the value of $f_n(x)$ is defined as the proportion of observed values in the sample which are less than or equal to x .
 - In other words if exactly k of the observed values in the sample are less than or equal to x , then
$$f_n(x) = \frac{k}{n}.$$
- The function $f_n(x)$ defined this way is called the sample distribution function (cdf).
- Sometimes $f_n(x)$ is called the empirical distribution function. If we let $y_1 < y_2 < \dots < y_n$ denote the value of the order statistics of the sample, then
$$f_n(x) = 0 \text{ for } x < y_1$$
$$f_n(x) \text{ jumps to the value } \frac{1}{n} \text{ at } x = y_1 \text{ and remains}$$
 at $\frac{1}{n}$ for $y_1 \leq x \leq y_2$.
- $f_n(x)$ jumps to $\frac{2}{n}$ at $x = y_2$ and remains at $\frac{2}{n}$ for $y_2 \leq x \leq y_3$ and so on.

(2)

Graphically



- Now let $f(x)$ denote the distribution function of the distribution from which the sample was drawn.
- For a given number x ($-\infty < x < \infty$) the probability that any particular x_i is less than or equal to x is $F(x)$.

Therefore it follows that from the law of large numbers (\bar{X}_n) that as $n \rightarrow \infty$ the proportion of $f_n(x)$ of the observations in the sample where one less than or equal to x will converge to $f(x)$ in probability.

$$P\left(\lim_{n \rightarrow \infty} f_n(x)\right) = f(x) \quad \text{--- --- ---} \quad (4)$$

for $-\infty < x < \infty$

where n is the sample size of x_1, x_2, \dots, x_n .
 $f_n(x)$ is the sample distribution function.
 $f(x)$ is the distribution function

- The relation (1) expresses the fact that at each point x the sample distribution function $f_n(x)$ will converge to the actual $f(x)$ of the distribution from which the random sample was drawn.

(3)

$$\text{Let } D_n = \sup_{-\infty < x < \infty} |F_n(x) - F(x)|$$

- Suppose we now wish to test the simple hypothesis H_0 : the unknown distribution function $f(x)$ is actually a particular continuous distribution $f^*(x)$ vs the general alternative that the actual distribution function is not $f^*(x)$ i.e

$$H_0: f_n(x) = f^*(x) \quad -\infty < x < \infty \quad (2)$$

H_1 : the hypothesis H_0 is not true

- This is a non parameter problem because the unknown distribution from which the random sample was taken might be any continuous distribution.
 - Let $F_N(x)$ denote the sample distribution function (sdf) and let
- $$D_n = \sup_{-\infty < x < \infty} |F_n(x) - F^*(x)|$$
- In other words D_n is the maximum difference between the sdf $F_n(x)$ and the hypothesized $F^*(x)$
 - When H_0 in eqn (2) is true the probability distribution of D_n will be a certain distribution which is the same for any possible continuous distribution $F^*(x)$ does not depend on a particular distribution function being studied in a specific problem.
 - Tables of this distribution for various values of n (sample size) have been developed and are tabulated in collection of statistical tables. The tables gives the values of D_n such that:

(4)

$$\begin{aligned} P(\hat{\theta}_n \leq \theta_n - \alpha) &= 1 - \alpha \\ \Rightarrow P(\hat{\theta}_n > \theta_n - \alpha) &= \alpha \end{aligned}$$

- Below are some critical values for the Kolmogorov-Smirnov test

n	α	0.1 0.05 0.01
5	0.45	0.51
10	0.32	0.37
15	0.27	0.26
20	0.23	0.24
25	0.21	0.22
30	0.19	0.20
35	0.18	0.19
40	0.17	0.18
45	0.16	0.17
50	0.15	0.17
750	1.07	1.22

n	\sqrt{n}	$\frac{1}{\sqrt{n}}$	$\frac{1.36}{\sqrt{n}}$	$\frac{1.03}{\sqrt{n}}$
750	25	0.04	0.0544	0.0412

Example

Test the hypothesis by the Kolmogorov-Smirnov test that the following sample come from a standard normal distribution

-1.23	1.64	-0.24	0.70	1.40	0.44	-0.07	-0.02
-0.15	1.76	1.62	0.40	-2.11	-0.99	-0.42	0.81
1.47	-2.46	0.88	1.39	0.42	0.37	-0.39	-0.10
1.07							

(5)

Solution

We want to test

$H_0: f(x) = f^*(x)$ for $-d < x < d$

$H_1:$ the hypothesis is not true.

where $f^*(x) = \Phi(x)$ the cdf of the standard normal distribution at $x = 0.05$

$$n = 25, b_{25}^{0.05} = 0.27 \text{ (from tables)}$$

We reject H_0 if $b_n \geq b_{0.05} = 0.27$

where $b_{25} = \sup$

$$\sup_{-d < x < d} |f_{25}(x) - f^*(x)|$$

Arrange the sample values in ascending order and for each value determine $f_n(x)$ and $f^*(x)$.

i	x	$F_N(x)$	$f^*(x)$	$ f_n(x) - f^*(x) $
1	-2.46	$\frac{1}{25} = 0.04$	0.007	0.033
2	-2.11	$\frac{2}{25} = 0.08$	0.017	0.063
3	-1.23	$\frac{3}{25} = 0.12$	0.109	0.011
4	-0.99	$\frac{4}{25} = 0.16$	0.161	0.001
5	-0.42	0.20	0.337	0.137 $\rightarrow \sup$
6	-0.39	0.24	0.348	0.108
7	-0.21	0.28	0.417	0.137 $\rightarrow \sup$
8	-0.15	0.32	0.440	0.120
9	-0.10	0.36	0.460	0.10
10	-0.07	0.40	0.472	0.072
11	-0.02	0.44	0.492	0.052
12	0.27	0.48	0.606	0.126
13	0.40	0.52	0.655	0.135
14	0.42	0.56	0.663	0.103
15	0.44	0.60	0.670	0.070
16	0.70	0.64	0.758	0.114
17	0.81	0.68	0.791	0.111

(6)

i.	x	$F_N(x)$	$F^*(x)$	$F_N(x) - F^*(x)$
18	0.88	0.72	0.81	0.091
19	1.07	0.76	0.858	0.093
20	1.31	0.80	0.918	0.118
21	1.40	0.84	0.919	0.079
22	1.47	0.88	0.929	0.049
23	1.62	0.92	0.947	0.027
24	1.64	0.96	0.950	0.010
25	1.76	1.0	0.961	0.039

Getting $f^*(x)$

$$f^*(x) = 2 \cdot 46 = \rho \left(x \leq -\frac{x-\mu}{\sqrt{\sigma^2}} \right)$$

$$= -3.46 - \mu$$

$$\sqrt{\sigma^2}$$

$$\Rightarrow \rho \left(z \leq -\frac{3.46 - \mu}{\sqrt{\sigma^2}} \right)$$

for standard normal
 $Z \sim N(0, 1)$

$$\mu = 0, \sigma^2 = 1$$

$$\Rightarrow \rho(z \leq -3.46) = 1 - \rho(z \geq -3.46)$$

$$= 1 - \phi(-3.46)$$

$$= 1 - 0.9931$$

$$\approx 0.0069$$

$$\lambda_{25} = 0.137 \text{ and } \lambda_{25}^{0.05} = 0.27$$

$$\lambda_{25} = 0.137 < \lambda_{25}^{0.05} = 0.27 \text{ hence we do not}$$

Reject H_0 and conclude that the sample came from a standard normal distribution.

Exercise

Test the hypothesis by Kolmogorov-Smirnov test that the following sample values came from a normal distribution with mean μ and variance σ^2 at $\alpha = 0.01$

2.72	3.84	0.88	5.72	5.48	3.12	0.10	2.48
1.76	0.52	2.64	3.64	3.40	1.80	-0.52	-0.12
2.30	3.10	1.04	1.02				