# Statistical Inference

# Non-Parametric Tests II

# Contents

# Kruskal Wallis test

- The Kruskal Wallis test is considered a non-parametric alternative to one way analysis of variance (ANOVA).

- The Kruskal Wallis test is used to compare differences between more than two independent groups when the dependent variable is either ordinal or continuous, but not normally distributed.

- H0:  K samples come from the same population
  H1: Not H0.

# Kruskal Wallis test procedure

- Combine all the observations from k samples into a single sample of size n and arrange them in ascending order .

- Assign ranks to them from smallest to largest as 1 to n. if there is a tie at two or more places, each observation is given the mean of the ranks for which it is tied.

- The ranks assigned to observations in each of the k groups are added separately to give k rank sums.

- The test statistic is

$$H = \frac{12}{n(n+1)} \sum_{i=1}^{k} \frac{R_j^2}{n_i} - 3(n+1)$$

$n_j = number\ of\ observations\ in\ j^{th}\ sample$

$n = number\ of\ observations\ in\ the\ combined\ sample$

$R_j = sum\ of\ the\ ranks\ in\ the\ j^{th}\ sample.$

- H follows Chi Square Distribution with k-1 df

# Case Study - 1

## Background

The data consists of the aptitude scores of 3 groups of employees.

## Objective

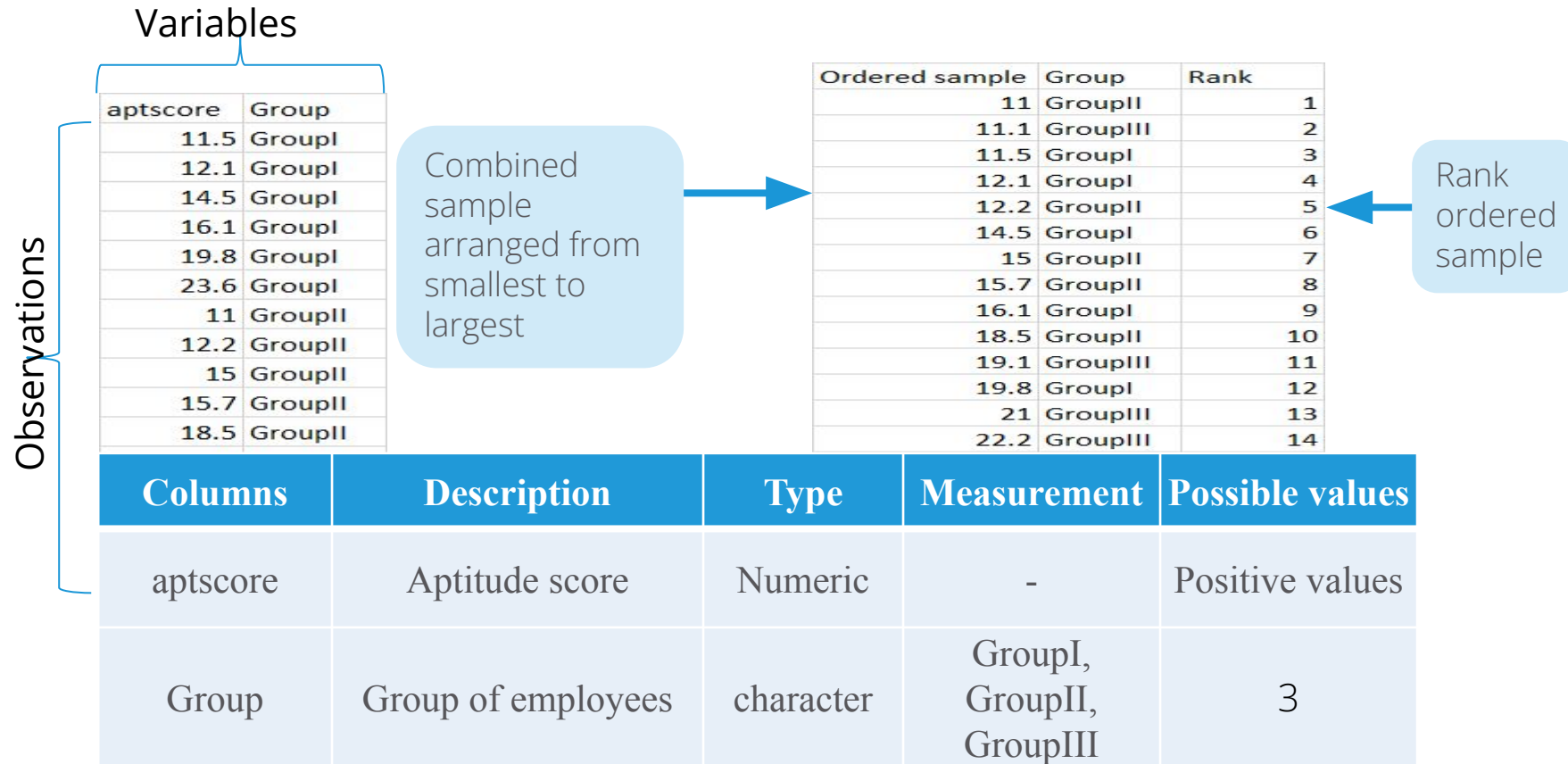**To check whether** there is difference in scores among the three groups.

## Sample Size

Sample size: 20
Variables: aptscore, Group

# Data Snapshot

**Kruskal Wallis Test**

Variables

Observations

| aptscore | Group |
|---|---|
| 11.5 | GroupI |
| 12.1 | GroupI |
| 14.5 | GroupI |
| 16.1 | GroupI |
| 19.8 | GroupI |
| 23.6 | GroupI |
| 11 | GroupII |
| 12.2 | GroupII |
| 15 | GroupII |
| 15.7 | GroupII |
| 18.5 | GroupII |

Combined sample arranged from smallest to largest

| Ordered sample | Group | Rank |
|---|---|---|
| 11 | GroupII | 1 |
| 11.1 | GroupIII | 2 |
| 11.5 | GroupI | 3 |
| 12.1 | GroupI | 4 |
| 12.2 | GroupII | 5 |
| 14.5 | GroupI | 6 |
| 15 | GroupII | 7 |
| 15.7 | GroupII | 8 |
| 16.1 | GroupI | 9 |
| 18.5 | GroupII | 10 |
| 19.1 | GroupIII | 11 |
| 19.8 | GroupI | 12 |
| 21 | GroupIII | 13 |
| 22.2 | GroupIII | 14 |

Rank ordered sample

| Columns | Description | Type | Measurement | Possible values |
|---|---|---|---|---|
| aptscore | Aptitude score | Numeric | - | Positive values |
| Group | Group of employees | character | GroupI, GroupII, GroupIII | 3 |

# Kruskal Wallis test

Testing distribution of more than two samples

| Objective | To test the **null hypothesis** that all the samples came from same population |
|---|---|

Null Hypothesis (H$_0$):  The three samples are from the same population
Alternate Hypothesis (H$_1$): The three samples do not come from the same population

| Test Statistic | $H = \dfrac{12}{n(n+1)} \displaystyle\sum_{j=1}^{k} \dfrac{R_j^2}{n_j} - 3(n+1)$ | $n_j = $ number of observations in $j^{th}$ sample <br> $n = $ number of observations in the combined sample <br> $R_j = $ sum of the ranks in the $j^{th}$ sample. |
|---|---|---|
| Decision Criteria | Reject the null hypothesis if the p-value < 0.05 | |

# Kruskal Wallis test example

Calculations :

| | Value |
|---|---|
| Sample size | $n_1 = 6$ <br> $n_2 = 7$ <br> $n_3 = 7$ |
| $R_1$ | 50 |
| $R_2$ | 68 |
| $R_3$ | 92 |
| H | 2.2309 |
| p-value | 0.3278 |

# Kruskal Wallis test in R

```
# Import the CSV file

data<-read.csv("Kruskal Wallis Test.csv",header=TRUE)


# Kruskal walis test

kruskal.test(formula=aptscore~Group,data=data)
```

- ❏ *kruskal.test performs the Kruskal waliss test on the data.*
- ❏ *aptscore is the analysis variable.*
- ❏ *Group is the factor variable.*

# Kruskal Walis test in R

```
# Output:
```

```
        Kruskal-Wallis rank sum test

data:   aptscore by Group
Kruskal-Wallis chi-squared = 2.2309, df = 2, p-value = 0.3278
```

*Interpretation :*
- *Since the p-value is >0.05, do not reject H0.  Aptitude score is the same for all three groups of employees.*

# Chi-square test of Association

- The chi-square test for independence, also called as Pearson's chi-square test or the chi-square test of association, is used to test if there is a relationship between two categorical variables.

- The two categorical variables can be nominal or ordinal.

- H0:  Two attributes are independent (not associated)
  H1: Not H0.

# Chi-square test procedure

- Assume that there are 'r' categories of attribute A and 'c' categories of attribute B. Therefore, we have a cross table of r*c (r rows and c columns).

- Let Ri be the total of the ith row and Cj be the total of the jth column.

- Observed frequencies are calculated from the data. Oij: Observed frequency in ith row and jth column.

- Expected frequencies are given by Eij = (Ri * Cj)/ n where n is total sample size. Expected frequencies are computed under the null hypothesis.

- Test statistic

$$\chi^2 = \sum_{i=1}^{r} \sum_{j=1}^{c} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where Oij are the observed frequencies in the ith row and jth column. Eij are the expected frequencies in the ith row and jth column.

- $\chi^2$ follows a Chi-Square Distribution with **(r-1)(c-1)** degrees of freedom.

# Case Study - 2

**Background**

The data consists of information regarding the Performance & Recruitment Source of employees.

**Objective**

To check whether Performance & Source of Recruitment are associated.

**Sample Size**

Sample size: 870
Variables: sn, performance, source

# Data Snapshot

**Variables**

**chi square test of association**

**Observations**

| sn | performan | source |
|---|---|---|
| 1 | Excellent | Internal |
| 2 | Excellent | Internal |
| 3 | Excellent | Internal |
| 4 | Excellent | Internal |
| 101 | Excellent | Campus |
| 102 | Excellent | Campus |
| 251 | Excellent | Jobportal |
| 252 | Excellent | Jobportal |
| 253 | Excellent | Jobportal |
| 254 | Excellent | Jobportal |
| 291 | Good | Internal |
| 292 | Good | Internal |
| 293 | Good | Internal |
| 491 | Good | Jobportal |
| 492 | Good | Jobportal |
| 493 | Good | Jobportal |
| 591 | Poor | Internal |

| Columns | Description | Type | Measurement | Possible values |
|---|---|---|---|---|
| sn | Serial number | Numeric | - | - |
| performance | Employee performance | character | Excellent, Good,Poor | 3 |
| source | Source of recruitment | Character | Campus, Internal, Jobportal | 3 |

Get the observed frequency (count) table from this data.

# Chi-square test of Association

Testing association between two categorical variables

| Objective | To test the **null hypothesis** that two categorical variables are independent |
|---|---|

Null Hypothesis (H$_0$):  performance and source are not associated

Alternate Hypothesis (H$_1$): performance and source are associated

| Test Statistic | $$\chi^2 = \sum_{i=1}^{r} \sum_{j=1}^{c} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$ | Oij = observed frequencies in the ith and jth column. Eij = expected frequencies in the ith row and jth column. |
|---|---|---|
| Decision Criteria | Reject the null hypothesis if the p-value < 0.05 | |

# Chi-square test example

## Observed Frequency table

| Performance | Recruitment Source | | | |
| --- | --- | --- | --- | --- |
| | Campus | Internal | Jobportal | Total |
| Excellent | 150 | 100 | 40 | 290 |
| Good | 100 | 100 | 100 | 300 |
| Poor | 80 | 50 | 150 | 280 |
| Total | 330 | 250 | 290 | 870 |

## Expected Frequency table

| Performance | Recruitment Source | | | |
| --- | --- | --- | --- | --- |
| | Campus | Internal | Jobportal | Total |
| Excellent | =(330*290)/870 | 83 | 97 | 290 |
| Good | 114 | =(250*300)/870 | 100 | 300 |
| Poor | 106 | 80 | =(290*280)/870 | 280 |
| Total | 330 | 250 | 290 | 870 |

| | Value |
| --- | --- |
| r | 3 |
| c | 3 |
| $\chi^2$ | 107.3786 |

# Chi-Square test in R

```
# Import the CSV file
data<-read.csv("chi square test of association.csv", header=TRUE)

# Install and use the package "gmodels"
install.packages("gmodels")
library(gmodels)
```

*"gmodels" is needed for the contingency table. The table displays frequencies, relative frequencies of two categorical variables.*

```
# Chi-square test of association
CrossTable(data$performance, data$source, chisq=TRUE)
```

*CrossTable function performs Chi-square test of association when chisq=TRUE.*

# Chi-Square test in R

# Output:

```
   Cell Contents
|-------------------------|
|                      N  |
| Chi-square contribution |
|          N / Row Total  |
|          N / Col Total  |
|        N / Table Total  |
|-------------------------|


Total Observations in Table:  870
```

*Interpretation :*
- *Since the p-value is <0.05, reject H0. Recruitment source and employee performance are associated.*

| data$performance | data$source Campus | Internal | Jobportal | Row Total |
|---|---|---|---|---|
| Excellent | 150 | 100 | 40 | 290 |
| | 14.545 | 3.333 | 33.218 | |
| | 0.517 | 0.345 | 0.138 | 0.333 |
| | 0.455 | 0.400 | 0.138 | |
| | 0.172 | 0.115 | 0.046 | |
| Good | 100 | 100 | 100 | 300 |
| | 1.672 | 2.207 | 0.000 | |
| | 0.333 | 0.333 | 0.333 | 0.345 |
| | 0.303 | 0.400 | 0.345 | |
| | 0.115 | 0.115 | 0.115 | |
| Poor | 80 | 50 | 150 | 280 |
| | 6.467 | 11.531 | 34.405 | |
| | 0.286 | 0.179 | 0.536 | 0.322 |
| | 0.242 | 0.200 | 0.517 | |
| | 0.092 | 0.057 | 0.172 | |
| Column Total | 330 | 250 | 290 | 870 |
| | 0.379 | 0.287 | 0.333 | |

```
Statistics for All Table Factors


Pearson's Chi-squared test
------------------------------------------------------------
Chi^2 =  107.3786      d.f. =  4      p =  2.635987e-22
```

# Quick Recap

| Kruskal Wallis test | • Nonparametric alternative to one way ANOVA. |
|---|---|
| Chi-Square test | • Also called Pearson's chi-square test or the chi-square test of association. It is used to test if there is a relationship between two categorical variables (nominal or ordinal). |