

# Data analysis

Purpose: Analyze collected data to answer research questions

Steps:

- 1) Descriptive statistics
  - \* is data sensible?
  - \* unusual values
  - \* missing values, outliers & inconsistent values
- 2) Inferential Statistics
- 3) Interpret results

## Step 1: Descriptive statistics

- \* Frequencies
  - \* Absolute: Number of times a value appears in a dataset
  - \* Relative: Proportion or % a value appears respect dataset
- \* Measures of central tendency
  - \* Mode: Most common value
  - \* Median: Middle value
  - \* Mean: Arithmetic average
- \* Measures of distribution
  - \* Maximum/minimum: Self-descriptive
  - \* Interquartile range (IQR): Range between  $Q_1$  (25%) &  $Q_3$  (75%), spread of the middle 50% of the data.
  - \* Variance: Spread of a set of values around their mean. High value = data very spread respect mean
  - \* Std. Deviation:  $\sqrt{\text{variance}}$ . Average distance of data points respect the mean.  
Large value = variability among data points

## Data Validation

- \* Check consistency & credibility
- \* Ensuring correctness & completeness of collected data
- \* Identifying missing variables
- \* Identifying outliers, → How to deal?
  - Exclude
  - Report two analysis (w & w/o)
  - Transform variable
  - Gather more data
- \* Create box-plots
  - Analyse if there are systematic explanations for these values
  - Perform sensitive analysis (e.g. scatter plot)

## Dispersion example

We will compute the Standard deviation

$$\text{Dataset} = \{1, 2, 2, 2, 3, 14\}$$

From a population ( $\sigma$ )

From a sample ( $s$ )

### 1) Computation Std dev $\sigma$

#### 1) Compute mean of the dataset

$$\bar{x} = \frac{1+2+2+2+3+14}{6} = 4$$

#### 2) Compute $\sigma$

$$\sigma = \sqrt{\text{Variance}} \approx 4.151$$

#### 2) Compute difference between each point and the mean, and square it

$$(x_i - \bar{x})^2$$

$$(1-4)^2 = 9$$

$$(2-4)^2 = 4$$

$$(2-4)^2 = 4$$

$$(2-4)^2 = 4$$

$$(3-4)^2 = 1$$

$$(14-4)^2 = 100$$

#### 3) Compute the variance

$$\frac{\sum (x_i - \bar{x})^2}{N}$$

If we compute over population we use N

$$\text{Variance} = \frac{(9+4+4+4+1+100)}{6} = \frac{122}{6} = 20.33$$

### 2) Computation Std dev $s$

#### 1) Compute mean of the dataset

$$\bar{x} = 4$$

#### 2) Compute difference between each point and the mean, and square it

$$(x_i - \bar{x})^2: \{9, 4, 4, 4, 1, 100\}$$

#### 3) Compute the variance

$$\frac{\sum (x_i - \bar{x})^2}{N-1}$$

If we compute over sample we use N-1

$$\text{Variance} = \frac{(9+4+4+4+1+100)}{5} = 24.4$$

#### 4) Compute $s$

$$s = \sqrt{\text{Variance}} \approx 4.94$$

### 3) Computing IQR

Methods:

\* Tukey:

\* Moore & McCabe:

\* Mendenhall & Sincich:

\* Freund & Perles:

\* Minitab:

### Handling dependency

\* 2 variables

Relationship between X & Y →

\* linear regression:  $y = a + bx$

\* Correlation coefficient (Pearson)

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \text{degree of linear dependence}$$

between 2 variables

\* > 2 variables

\* Multivariate analysis

\* Principal component

...

## Step 2: Statistical analysis

### Hypothesis testing

#### \* Definition

Procedure for determining whether experimental results of a sample provide support for the validity of a thesis/hypothesis

#### \* Basic questions

1) Is the relationship between 2+ variables due to random chance?

2) If chance is excluded, what does it mean?

\* There might be a meaningful association between variables

#### \* Approach (Falsification principle)

Opposite-of-what-you-predict reasoning

Draw conclusions by

- rejecting the inverse hypothesis although it is true ( $< 0.05$ )

#### \* Procedure

1) Formulate alternative & null hypothesis

2) Select a statistical test

△ considering data distribution → \* Normal distribution → Parametric test

\* non-normal/ordinal/nominal distribution → Non-parametric test

3) Apply test

4) Use a significance level ( $\alpha$ )

It represents the probability of rejecting  $H_0$  when it is actually true → 0.01 or 0.05

→ Type I error

5) Perform power analysis ( $\beta$ )

Probability of failing to reject  $H_0$  when it is actually false → 0.2

→ Type II error

6) Interpretation

If  $H_0$  rejected:

There is an effect

↳ calculate real effect size

↳ Is effect in accordance with the theory?

↳ Modify theory?

If  $H_0$  rejected:

- It is not possible to conclude there is no effect
- There is no sufficient evidence to accept that there is an effect
- Discuss descriptive statistics
- Calculate real effect size → Determine sample size to reject  $H_0$
- Formulate new hypothesis

Example Simple test where 2 techniques are applied to students.  
on experiencing results are very similar. B seems better (but slightly),  
How do we make a decision?

1) Compute difference between means

$$\text{diff} = (\bar{x}_b - \bar{x}_a)$$

2) locate diff in the histogram

3) calculate area that falls in the right side of the histogram

That area is the probability that we could obtain a result  $\geq (\bar{x}_b - \bar{x}_a)$   $\rightarrow p\text{-value}$

4) Analyze

If  $p\text{-value} < \alpha$

Techniques are not alike  $\rightarrow$  significant result

## Parametric VS non-parametric

\* T-test  $\rightarrow$  parametric

\* Main diff  $\rightarrow$  Assumption of the distribution of the sample

Non parametric: does not make any assumption

## Parametric testing

### 1) T-Test (independent sample)

#### 1.1) One sample t-test

Compare  $\bar{x}$  response of a group against a specific value

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

$\bar{x}$  = mean of group  
 $\mu_0$  = specified value  
 $n$  = # subjects/group  
 $s$  = std dev group

Look in t-distribution table to obtain p-value

#### ⚠ Requirements

- 1) Samples must be independent & identically distributed
  - 2) Mean estimator should be normally distributed
  - 3) Response variables measured on ratio scales
- If any does not hold ↓  
 ↗ non-parametric test
- ⚠ ordinal metrics can't be used

#### 1.2) Two sample t-test

• Checks → statistical significance of the difference between the mean responses of two levels of a factor

→  $H_0$  of 2 subpopulations where the mean is the same

$$H_0: \mu_1 = \mu_2$$

$$\sigma_1 = \sigma_2$$

• Pre-requisites:

- \* The two sample sizes are equal
- \* The two distributions have the same variance

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_{x_1, x_2} \cdot \sqrt{\frac{1}{n_1}}}$$

$$s_{x_1, x_2} = \sqrt{\frac{1}{2} (s_{x_1}^2 + s_{x_2}^2)}$$

•  $\bar{x}$  = mean

•  $s$  = std dev

•  $n$  = # subjects (equal in both groups)      •  $s^2$  = unbiased estimators of the variance

### 1.3) Special cases

#### 1.3.1) Unequal sample size

⚠ Equal variance assumed

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_{x_1, x_2} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$s_{x_1, x_2} = \sqrt{\frac{(n_1-1)s_{x_1}^2 + (n_2-1)s_{x_2}^2}{n_1+n_2-2}}$$

#### 1.3.2) Equal/Unequal sample sizes & unequal variances

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_{\bar{x}_1 - \bar{x}_2}}$$

$$s_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

## 2) Paired T-test (Dependent sample)

Repeated measures or matched pairs

⚠ Diff between all pairs must be calculated

$$t = \frac{\bar{X}_D - \mu_0}{S_D / \sqrt{n}} \quad \bar{X}_D = \text{mean of diff between pairs}$$

## 3) One factor ANOVA

Checks the statistical significance of the difference between the mean responses of one factor with several levels

Steps

### 1) Identify mathematical model

$$Y_{ij} = \mu + X_j + e_{ij}$$

•  $\mu$ : overall mean response  
•  $X_j$ : level of factor  $j$   
•  $e_{ij}$ : error

\* Factor: independent variable selected to determine its effect on the response variable

\* Level of factor: Different categories/distinct values that the factor can take on

$$\text{Where } x_{ij} = \bar{Y}_{\cdot j} - \mu = \bar{Y}_{\cdot j} - \bar{Y}_{..}$$

### 2) Validation of the basic model that relates the experiment variables

Confirm assumptions of ANOVA are met

### 3) Calculate the factor-induced variation in the response variable.

ANOVA decomposes total variation in the response variable into different components:

- 1) Variation between groups (induced by the factor)
  - 2) Variation within groups (random error)
- How much of the variability in the response variable is attributable to differences between the levels of factor

### 4) Calculate the statistical significance of the factor-induced variable

- Asses whether the differences in mean responses among the levels of factor are significant

• F-test → if f-ratio large & f-ratio > f-distribution diff between means among levels of the factor

### 5) Establish recommendations on the alternative that provides the best response variable values

After determining the statistical significance of the factor-induced variation we draw conclusions.

If  $H_0$  rejected:

≥ 1 level of factor differs from the others

## Non-parametric testing

### 1) MANN-WHITNEY-U Test (Independent sample)

- For independent groups
- It can be used instead of t-test if data is not normally distributed
- Robust against outliers

#### Prerequisites

- Responses are on an ordinal scale
- Distributions of both groups are equal under the  $H_0$

#### 1.1) Method 1

Small samples.

Choose smaller sample (for easy computation)

For each observation in sample 1  
count smaller rank obs. in sample 2      } sum result = U

#### 1.2) Method 2

Large samples

$$U_1 = n_1 \cdot n_2 \cdot \frac{n_1(n_1+1)}{2} - R_1 \quad \text{sum of ranks of group 1}$$

Reject  $H_0$  if:

$$U_2 = n_1 \cdot n_2 \cdot \frac{n_2(n_2+1)}{2} - R_2 \quad \min(U_1, U_2) \leq \text{critical value}$$

$\alpha$

### 2) Wilcoxon signed rank test

Alternative to paired t-test

#### Pre-requisite

It must be possible to determine which value is larger and to rank the differences

### 3) Sign-test

- Alternative to paired t-test
- Use if:

Not possible to rank the differences

△ Must be on ordinal scale at least

• Formula:  $p = \frac{1}{2^N} \sum_{i=0}^n (N)_i$