

Descriptive Statistics

Formulary

I. Distribution Analysis

Absolute frequency:

$F(a_j) :=$ “Number of cases in which the attribute a_j occurs”, $j = 1, 2, \dots, k$.

Relative frequency:

$$f(a_j) := \frac{1}{n} \cdot F(a_j), \quad j = 1, 2, \dots, k.$$

n : Number of observations.

Empirical distribution function of an ungrouped data:

$$F(x) := \begin{cases} 0 & \text{for } x \leq a_1 \\ \sum_{i=1}^j f(a_i) & \text{for } a_j < x \leq a_{j+1} \\ 1 & \text{for } x > a_k \end{cases}$$

Properties of the empirical distribution function:

1.

$$F(x) = P(X < x) \quad (P: \text{proportion})$$

2.

$$0 \leq F(x) \leq 1$$

3.

$$\forall x_1, x_2 : x_1 < x_2 \Rightarrow F(x_1) \leq F(x_2)$$

4.

$$P(x_1 \leq X < x_2) = F(x_2) - F(x_1)$$

5.

$F(x)$ is at least left-sided continuous and has at most a finite number of jump discontinuities.

6.

$$x \rightarrow -\infty \Rightarrow F(x) \rightarrow 0$$

$$x \rightarrow +\infty \Rightarrow F(x) \rightarrow 1$$

Empirical distribution function of grouped data:

$$F(x) := \begin{cases} 0 & \text{for } x \leq m_1 \\ \sum_{i=1}^j f(a_i) & \text{for } m_j < x \leq m_{j+1} \\ 1 & \text{for } x > m_k \end{cases}$$

m_i : Midpoint of the class C_i , $i = 1, 2, \dots, p$.

II. Measures of Location

Arithmetic average (mean) of an ungrouped data:

$$\mu := \frac{1}{N} \sum_{i=1}^N x_i \quad (\text{Simple arithmetic average for population data})$$

$$= \frac{1}{N} \sum_{j=1}^k a_j \cdot F(a_j) \quad (\text{Weighted arithmetic average for population data})$$

N : Population size.

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{Simple arithmetic average for sample data})$$

$$= \frac{1}{n} \sum_{j=1}^k a_{ij} \cdot F(a_j) \quad (\text{Weighted arithmetic average for sample data})$$

n : sample size.

Arithmetic average (mean) of a grouped data:

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{Simple arithmetic average for sample data})$$

$$= \frac{1}{n} \sum_{j=1}^k a_j \cdot F(a_j) \quad (\text{Weighted arithmetic average for sample data})$$

n : Sample size.

$$\mu \approx \frac{1}{N} \cdot \sum_{i=1}^p \bar{x}_i \cdot F_i \quad (\text{Population mean})$$

$$\bar{x} \approx \frac{1}{n} \cdot \sum_{i=1}^p \bar{x}_i \cdot F_i \quad (\text{Sample mean})$$

(In case the informations \bar{x}_i are not available, they will be replaced by the interval *midpoints*:

$$m_i := \frac{b_i + B_i}{2}$$

$b_i, i = 1, 2, \dots, p$: lower bound of the class C_i

$B_i, i = 1, 2, \dots, p$: upper bound of the class C_i .)

Estimations:

$$\frac{1}{N} \cdot \sum_{i=1}^p b_i \cdot F_i \leq \mu \leq \frac{1}{N} \cdot \sum_{i=1}^p B_i \cdot F_i$$

$$\frac{1}{n} \cdot \sum_{i=1}^p b_i \cdot F_i \leq \bar{x} \leq \frac{1}{n} \cdot \sum_{i=1}^p B_i \cdot F_i$$

Median of an ungrouped data:

$$Me := \begin{cases} x_{\left[\frac{n+1}{2}\right]} & \text{if } n \text{ is odd} \\ \frac{1}{2} \left(x_{\left[\frac{n}{2}\right]} + x_{\left[\frac{n}{2}+1\right]} \right) & \text{if } n \text{ is even} \end{cases}$$

Median of a grouped data:

Let C_1, C_2, \dots, C_p be given classes with the relative frequencies f_1, f_2, \dots, f_p .

The *Median* of lies in the class C_i if

$$\sum_{j=1}^{i-1} f_j < 0.5 \quad \text{and} \quad \sum_{j=1}^i f_j \geq 0.5$$

and will be defined by

$$Me := b_i + \frac{0.5 - \sum_{j=1}^{i-1} f_j}{f_i} \cdot w_i$$

α -Quantile of an ungrouped data:

Given a *ranked* data set, the, α -Quantile ($0 < \alpha < 1$) is defined by:

$$\tilde{x}_\alpha := \begin{cases} x_{[k]} & \text{if } n \cdot \alpha \text{ is not integer} \\ & (k \text{ is then the following integer}) \\ \frac{1}{2} (x_{[k]} + x_{[k+1]}) & \text{if } n \cdot \alpha \text{ is integer} \\ & (k = n \cdot \alpha) \end{cases}$$

α -Quantile of a grouped data:

Let C_1, C_2, \dots, C_p be given classes with the relative frequencies f_1, f_2, \dots, f_p .

The α -Quantile ($0 < \alpha < 1$) of the data set lies in the class C_i if

$$\sum_{j=1}^{i-1} f_j < \alpha \quad \text{and} \quad \sum_{j=1}^i f_j \geq \alpha$$

and will be defined by

$$\bar{x}_\alpha := b_i + \frac{\alpha - \sum_{j=1}^{i-1} f_j}{f_i} \cdot w_i$$

Mode of an ungrouped data:

The value of the item which occurs most frequently is called the *mode* of it is unique.

It will be denoted by *Mo*.

Mode of a grouped data:

The *mode* of a grouped data set lies in the class with the highest frequency. It will be defined by

$$Mo \approx b_i + \frac{f_i - f_{i-1}}{2f_i - f_{i-1} - f_{i+1}} \cdot w_i$$

Geometric mean of an ungrouped data

$$\bar{x}_g := \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n} = \sqrt[n]{\prod_{i=1}^n x_i}, \quad x_i > 0, \quad i = 1, 2, \dots, n$$

III. Measures of Distribution

Range:

1. *Ungrouped*

$$R := x_{\max} - x_{\min}$$

2. *Grouped*

$$R \approx B_p - b_1$$

Interquartile range:

$$R_Q := x_{0.75} - x_{0.25}$$

Variance (ungrouped):

$$\sigma^2 := \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 \quad (\text{population variance})$$

$$s^2 := \frac{1}{n-1} \sum_{i=1}^n \left(x_i - \bar{x} \right)^2 \quad (\text{sample variance})$$

Standard deviation (ungrouped):

$$\sigma := \sqrt{\sigma^2}, \quad \sigma > 0 \quad (\text{population standard deviation})$$

$$s := \sqrt{s^2}, \quad s > 0 \quad (\text{population standard deviation})$$

Variance (grouped):

$$\sigma^2 := \frac{1}{N} \sum_{i=1}^p \left(\bar{x}_i - \mu \right)^2 \cdot F_i \quad (\text{population})$$

$$s^2 := \frac{1}{n-1} \sum_{i=1}^p \left(\bar{x}_i - \bar{x} \right)^2 \cdot F_i \quad (\text{sample})$$

In case the information \bar{x}_i is not available, it will be replaced by m_i .

Standard deviation (grouped):

$$\sigma := \sqrt{\sigma^2}, \quad \sigma > 0 \quad (\text{population})$$

$$s := \sqrt{s^2}, \quad s > 0 \quad (\text{sample})$$

Chebyshev Theorem

For any number $k > 1$, at least $(1 - 1/k^2)$ of the data values lie within k standard deviations of the mean.

At least 50% of the observations in a data set lie in the interval

$$[\mu - \sigma, \quad \mu + \sigma], \quad (\text{population})$$

$$\left[\bar{x} - s, \quad \bar{x} + s \right] \quad (\text{sample})$$

Coefficient of Variation:

$$v := \frac{\sigma}{\mu} \quad (\text{population})$$

$$v := \frac{s}{\bar{x}} \quad (\text{sample})$$

A Rule of Thumb:

The arithmetic average of a data set is to be considered as representative only if its coefficient of variation is less than 0.5 (or 50%).

IV. Correlation and Regression

Spearman's Rank Coefficient of Correlation)

$$\rho := 1 - \frac{6 \cdot \sum_{i=1}^n (R_i - R'_i)^2}{(n-1) \cdot n \cdot (n+1)}$$

Here are:

$R_i, i = 1, 2, \dots, n$: ranks of the characteristic X

$R'_i, i = 1, 2, \dots, n$: ranks of the characteristic Y ,

Karl Pearson's Coefficient of Correlation:

Let X and Y be two variates.

$$r := \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}, \quad -1 \leq r \leq +1.$$

Simple Linear Regression:

The coefficients of a simple linear regression function

$$y^* = a_0 + a_1 x$$

can be obtained by solving the following system of linear equations:

$$\begin{aligned} n \cdot a_0 + a_1 \cdot \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ a_0 \cdot \sum_{i=1}^n x_i + a_1 \cdot \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i \cdot y_i \end{aligned}$$

Coefficients of Correlation and Determination:

1. The *coefficient of correlation of a simple linear regression function* is defined by

$$r := \frac{n \cdot \sum_{i=1}^n x_i \cdot y_i - \sum_{i=1}^n x_i \cdot \sum_{i=1}^n y_i}{\sqrt{\left(n \cdot \sum_{i=1}^n x_i^2 - \sum_{i=1}^n x_i \cdot \sum_{i=1}^n x_i \right) \left(n \cdot \sum_{i=1}^n y_i^2 - \sum_{i=1}^n y_i \cdot \sum_{i=1}^n y_i \right)}},$$
$$-1 \leq r \leq +1$$

2. r^2 ($0 \leq r^2 \leq 1$) is called *coefficient of determination*.

Coefficient of Determination

The *coefficient of determination* in general is defined as:

$$r^2 := \frac{SSR}{SST}$$

where:

$$SSR = \sum_{i=1}^n (y_i^* - \bar{y})^2 : \text{sum of square due to regression}$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 : \text{total sum of squares}$$