## STATISTICAL LABORATORY



Applied Mathematics for Economics and Management Ist Year/1st Semester 2025/2026

## CONTACT

**Professor**: Elisabete Fernandes

**E-mail**: efernandes@iseg.ulisboa.pt



https://doity.com.br/estatistica-aplicada-a-nutricao



https://basiccode.com.br/produto/informatica-basica/

## **PROGRAM**



I. Fundamental Concepts of Statistics



2. Exploratory Data Analysis



3. Organizing and Summarizing Data



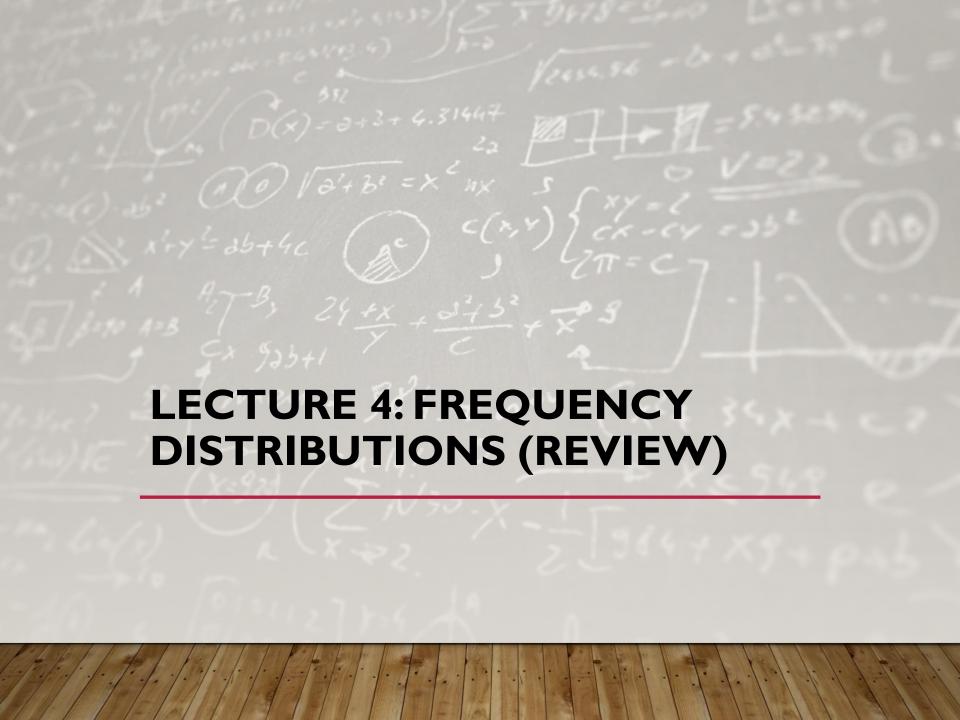
4. Association and Relationships Between Variables



5. Index Numbers



6.Time Series Analysis



# FREQUENCY TABLE FOR DISCRETE VARIABLE: EXAMPLE

#### 1. Discrete Variable Example

Suppose we collected n=15 observations of a discrete variable x:

$$x = (2, 3, 3, 4, 2, 5, 3, 4, 2, 5, 4, 3, 2, 4, 3)$$

Step 1: Identify distinct values

$$x_1'=2,\,x_2'=3,\,x_3'=4,\,x_4'=5$$

**Step 2: Count frequencies** 

#### **Absolute Frequencies**

$$n_1 = 4$$
,  $n_2 = 5$ ,  $n_3 = 4$ ,  $n_4 = 2$ 

Step 3: Compute relative frequencies

#### Relative Frequencies

$$f_j = \frac{n_j}{n} \Rightarrow f_1 = 0.267, f_2 = 0.333, f_3 = 0.267, f_4 = 0.133$$

Step 4: Determine range

$$x_{\min} = 2, \quad x_{\max} = 5, \quad \text{Range} \qquad = 3$$

# FREQUENCY TABLE FOR DISCRETE VARIABLE: EXAMPLE

Step 5: Table (Frequency Absolute Frequencies		Relative Frequencies
$x_j'$	$n_j$	$f_{j}$
2	4	0.267
3	5	0.333
4	4	0.267
5	2	0.133
Total	15	1.0

Sample Size n = 15

$$\sum_{j=1}^m n_j = n \qquad \sum_{j=1}^m f_j = 1$$

## FREQUENCY TABLE FOR CONTINUOUS VARIABLE: EXAMPLE

Suppose we collected n=12 observations of a continuous variable x:

$$x = (1.2, 2.3, 1.8, 2.0, 3.1, 2.7, 3.5, 1.5, 2.8, 3.0, 1.9, 2.5)$$

Step 1: Determine range and class width

$$x_{\min} = 1.2$$
,  $x_{\max} = 3.5$ , Range  $= 3.5 - 1.2 = 2.3$ 

Using **Sturges' rule**, m=4 classes:

$$ext{Class width} = rac{ ext{Range}}{m} = rac{2.3}{4} pprox 0.575 pprox 0.6$$

Step 2: Construct classes

#### Classes

Start at 1.2:

$$l_j = \left[l_{j-1}, l_j\right]$$
  $j = 1, 2, ..., m$   
 $m = \text{number of classes}$ 

$$I_j \cap I_k = \emptyset$$
 and

$$D \subset \bigcup_{j=1}^m I_j$$

Step 3: Find midpoints

$$D = [x_{min}, x_{max}]$$

$$x_1' = \frac{1.2 + 1.8}{2} = 1.5, \quad x_2' = \frac{1.8 + 2.4}{2} = 2.1, \quad x_3' = \frac{2.4 + 3.0}{2} = 2.7, \quad x_4' = \frac{3.0 + 3.6}{2} = 3.3$$

## FREQUENCY TABLE FOR CONTINUOUS VARIABLE: EXAMPLE

#### Step 4: Count frequencies

• Class [1.2, 1.8]: 1.2, 1.5, 1.8  $\rightarrow n_1 = 3$ 

#### **Absolute Frequencies**

• Class ]1.8, 2.4]: 1.9, 2.0, 2.3  $\rightarrow n_2 = 3$ 

• Class ]2.4, 3.0]: 2.5, 2.7, 3.0  $\rightarrow n_3 = 3$ 

• Class [3.0, 3.6]: 3.1, 3.5  $\rightarrow n_4 = 2$ 

Step 5: Compute relative frequencies

#### Relative Frequencies

$$f_1=rac{3}{12}=0.25, \quad f_2=0.25, \quad f_3=0.25, \quad f_4=rac{2}{12}pprox 0.167$$

#### Step 6: Table (Frequency Distribution with Totals)

Class $I_j$	$Midpoint x_j'$	$n_{j}$	$f_{j}$
[1.2, 1.8]	1.5	3	0.25
]1.8, 2.4]	2.1	3	0.25
]2.4, 3.0]	2.7	3	0.25
]3.0, 3.6]	3.3	2	0.167

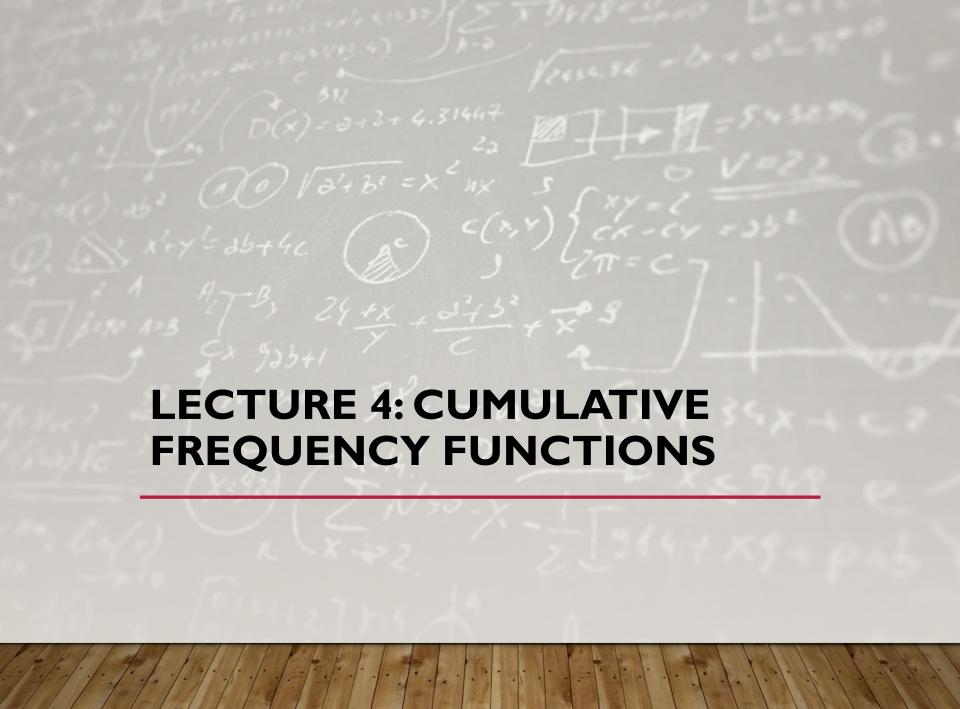
Total

12

1.0

Sample Size  $\sum_{j=1}^{m} n_j =$ 

 $\sum_{j=1}^{m} n_j = n \qquad \sum_{j=1}^{m} f_j = 1$ 



## **CUMULATIVE FREQUENCY FUNCTION** FOR DISCRETE VARIABLE

#### Cumulative Absolute Frequency Function

$$N(x) = \begin{cases} 0 & x < x_1' \\ n_1 & x_1' \le x < x_2' \\ n_1 + n_2 & x_2' \le x < x_3' \\ \dots & \dots \\ n_1 + n_2 + \dots + n_{m-1} & x_{m-1}' \le x < x_m' \\ n_1 + n_2 + \dots + n_m = n & x \ge x_m' \end{cases}$$
• It is a **non-decreasing function** and **right-continuous**.
• It is a **step function** where each jump corresponds to an absolute frequency. For exemple:  $N(x_3') - N(x_2') = (n_1 + n_2 + n_3) - (n_1 + n_2) = n_3$ 

#### • N(x) is a function whose domain is the real line, with $N(-\infty) = 0, N(+\infty) = n \text{ and } 0 \le N(x) \le n.$

#### Cumulative Relative Frequency Function

$$F^*(x) = \begin{cases} 0 & x < x_1' \\ f_1 & x_1' \le x < x_2' \\ f_1 + f_2 & x_2' \le x < x_3' \\ \dots & \dots \\ f_1 + f_2 + \dots + f_{m-1} & x_{m-1}' \le x < x_m' \\ f_1 + f_2 + \dots + f_m = 1 & x \ge x_m' \end{cases}$$

•  $F^*(x)$  has properties similar to N(x) , except with respect to the codomain:  $0 \le F_x^*(x) \le 1$ .

## CUMULATIVE FREQUENCY FUNCTION FOR DISCRETE VARIABLE: EXAMPLE

#### Sample information:

- Observations: x = (2, 3, 3, 4, 2, 5, 3, 4, 2, 5, 4, 3, 2, 4, 3)
- Number of observations: n=15
- Sample range: Sample range  $= x_{\rm max} x_{\rm min} = 5 2 = 3$

#### Frequency Distribution Table (with cumulative frequencies)

$x_j'$	$n_{j}$	$f_{j}$	$N(x_j')$	$F^*(x_j')$
2	4	0.267	4	0.267
3	5	0.333	9	0.600
4	4	0.267	13	0.867
5	2	0.133	15	1.0
Total	15	1.0	_	_

## Cumulative Absolute Frequency Function

$$N(x) = egin{cases} 0, & x < 2 \ 4, & 2 \leq x < 3 \ 9, & 3 \leq x < 4 \ 13, & 4 \leq x < 5 \ 15, & x \geq 5 \end{cases}$$

## Cumulative Relative Frequency Function

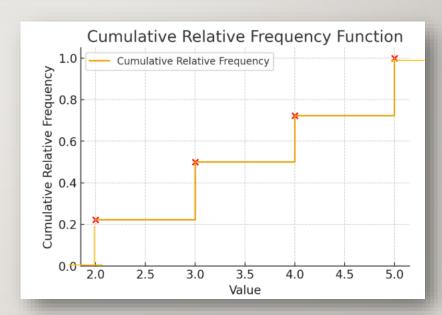
$$F^*(x) = egin{cases} 0, & x < 2 \ 0.267, & 2 \leq x < 3 \ 0.600, & 3 \leq x < 4 \ 0.867, & 4 \leq x < 5 \ 1, & x \geq 5 \end{cases}$$

## CUMULATIVE FREQUENCY GRAPH FOR A DISCRETE VARIABLE: EXAMPLE

## Cumulative Relative Frequency Function

$$F^*(x) = egin{cases} 0, & x < 2 \ 0.267, & 2 \leq x < 3 \ 0.600, & 3 \leq x < 4 \ 0.867, & 4 \leq x < 5 \ 1, & x \geq 5 \end{cases}$$





• Here is the graph of the **cumulative relative frequency function** (ogive) for your sample:

Values: 2, 3, 4, 5 Frequencies: 4, 5, 4, 5

- The red dots mark the cumulative proportions.
- This shows how the cumulative relative frequency increases step by step until it reaches 1 (100%).

## CUMULATIVE FREQUENCY FUNCTION FOR CONTINUOUS VARIABLE

## Cumulative Relative Frequency Function

$$F^*(l_0) = 0$$

$$F^*(l_1) = f_1$$

$$F^*(l_2) = f_1 + f_2$$
...
$$F^*(l_{m-1}) = f_1 + f_2 + \dots + f_{m-1}$$

$$F^*(l_m) = f_1 + f_2 + \dots + f_m = 1$$

It is assumed that the frequencies are uniformly distributed within each class and that the cumulative frequencies refer to the class boundaries.

## FREQUENCY CUMULATIVE FUNCTION FOR CONTINUOUS VARIABLE: EXAMPLE

#### Sample information:

- Observations: x = (1.2, 2.3, 1.8, 2.0, 3.1, 2.7, 3.5, 1.5, 2.8, 3.0, 1.9, 2.5)
- Number of observations: n=12
- ullet Sample range  $=x_{
  m max}-x_{
  m min}=3.5-1.2=2.3$
- Classes (Sturges, width ≈ 0.6): [1.2,1.8], ]1.8,2.4], ]2.4,3.0], ]3.0,3.6]
- Midpoints: 1.5, 2.1, 2.7, 3.3

#### Frequency Distribution Table (with cumulative frequencies)

Class	$l_j$	$n_{j}$	$f_{j}$	$N(x_j)$	$F^*(x_j)$
[1.2, 1.8]	1.2	3	0.25	3	0.25
]1.8, 2.4]	1.8	3	0.25	6	0.50
]2.4, 3.0]	2.4	3	0.25	9	0.75
]3.0, 3.6]	3.0	3	0.25	12	1.00
Total	_	12	1.00	_	_

Cumulative Absolute Frequency Function

$$N(x) = egin{cases} 0, & x < 1.2 \ 3, & 1.2 \leq x < 1.8 \ 6, & 1.8 \leq x < 2.4 \ 9, & 2.4 \leq x < 3.0 \ 12, & x \geq 3.0 \end{cases}$$

## Cumulative Relative Frequency Function

$$F^*(x) = egin{cases} 0, & x < 1.2 \ 0.25, & 1.2 \leq x < 1.8 \ 0.50, & 1.8 \leq x < 2.4 \ 0.75, & 2.4 \leq x < 3.0 \ 1.0, & x \geq 3.0 \end{cases}$$

## CUMULATIVE FREQUENCY GRAPH FOR A CONTINUOUS VARIABLE: EXAMPLE

Class	$l_j$	$n_{j}$	$f_{j}$	$N(x_j)$	$F^*(x_j)$
[1.2, 1.8]	1.2	3	0.25	3	0.25
]1.8, 2.4]	1.8	3	0.25	6	0.50
]2.4, 3.0]	2.4	3	0.25	9	0.75
]3.0, 3.6]	3.0	3	0.25	12	1.00
Total	_	12	1.00	_	_

Interval	
<b>Endpoints</b>	F*(x)
1,2	0
1,8	0,25
2,4	0,5
3	0,75
3,6	1

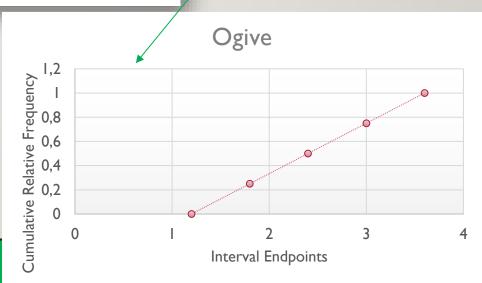
## Ogive (Cumulative Frequency Graph) Where does it start?

- It starts at the **lower limit of the first class**, with cumulative frequency = 0.
- The graph is built using the upper class limits on the x-axis and the cumulative frequencies (or percentages) on the y-axis.

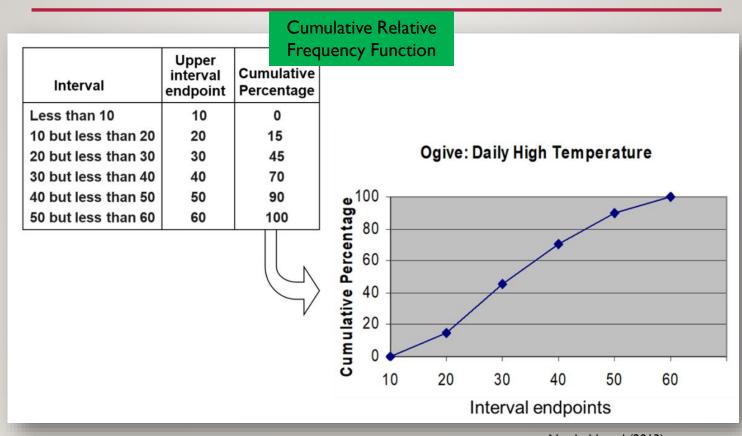
#### What is it used for?

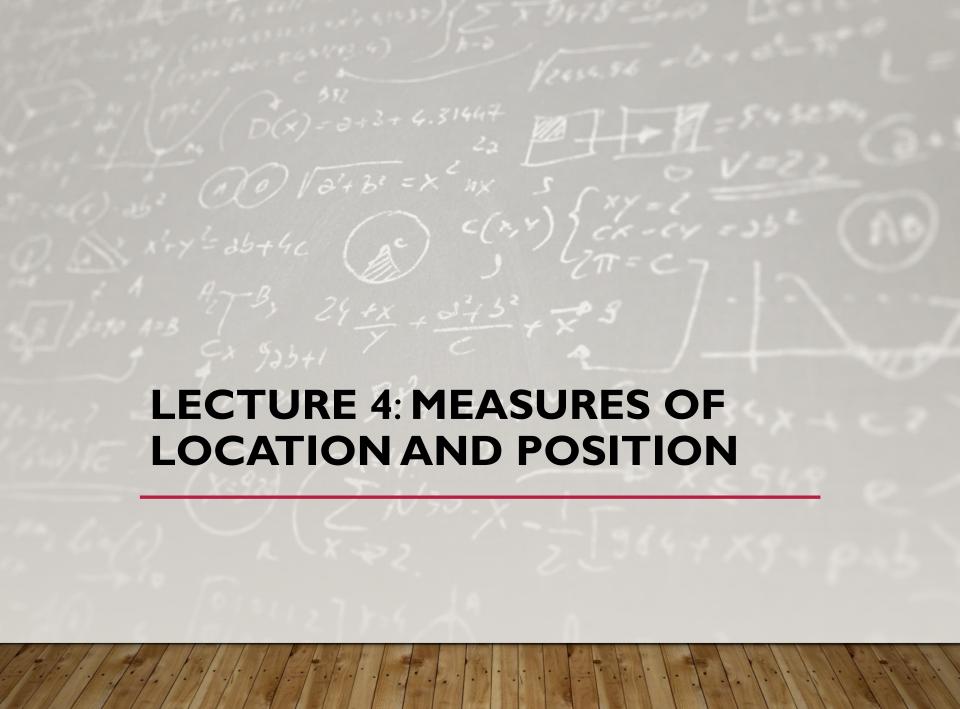
- To show how data accumulate across classes.
- To identify the median, quartiles, deciles, and percentiles.
- To compare cumulative distributions.

In this case, the ogive appears as a straight line rather than a curve, because the absolute and relative frequencies have the same values for all classes.



## THE OGIVE GRAPHING CUMULATIVE FREQUENCIES: OTHER EXAMPLE





# MEASURES OF LOCATION VS. MEASURES OF POSITION

#### Measures of Location

- Describe where the data are concentrated.
- Give an idea of the "center" or typical value of the data.
- Examples:
  - Mean sum of values divided by the number of observations.
  - Median value that separates the data into two equal halves.
  - Mode most frequently occurring value.
- Summary: indicate "on average, where the data are."

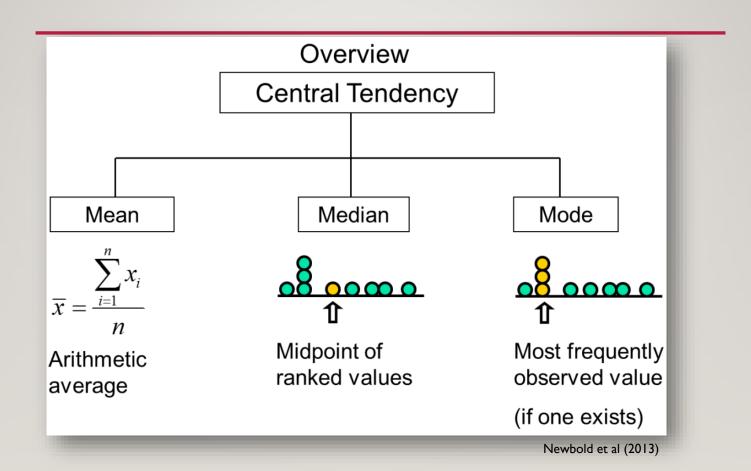
#### Measures of Position

- Indicate the relative position of a value within the dataset or the position of reference values.
- Useful for comparing values and identifying quantiles.
- Examples:
  - Quartiles, Deciles, Percentiles divide the data into equal parts and show relative positions.
  - **Z-score** shows how many standard deviations a value is above or below the mean.
- Summary: indicate "where a value stands in relation to the whole dataset."

#### **Key Difference:**

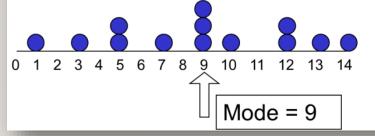
- Location: describes the center or typical value of the dataset.
- Position: describes the relative position of a value within the dataset.

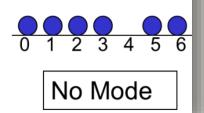
## **MEASURES OF LOCATION**



### MODE

- A measure of central tendency
- Value that occurs most often
- Not affected by extreme values
- Used for either numerical or categorical data
- There may be no mode
- There may be several modes





## **DISTRIBUTIONS BY MODE**

- 1 Amodal Distribution (No Mode)
- **Definition:** No value repeats → no mode
- Example: 2, 3, 5, 7, 11
- Unimodal Distribution (One Mode)
- Definition: One value appears most frequently
- Example: 1, 2, 2, 3, 4 → Mode = 2
- Bimodal Distribution (Two Modes)
- Definition: Two values appear with the same highest frequency
- Example: 1, 2, 2, 3, 3, 4 → Modes = 2, 3
- Multimodal Distribution (More than Two Modes)
- **Definition**: More than two values appear with the same highest frequency
- Example: 1, 1, 2, 2, 3, 3,  $4 \rightarrow$  Modes = 1, 2, 3

## MODE FOR GROUPED DATA

#### Definition:

The modal class is the class interval with the highest frequency in a grouped frequency distribution

King's Formula (Grouped Data):

$$\mathrm{Mo} = l + rac{f^{**}}{f^* + f^{**}} \cdot h$$

#### Legend / Notation:

- l = lower boundary of the modal class
- h = class width
- $f^*$  = relative frequency of the class before the modal class
- $f^{**}$  = relative frequency of the class after the modal class

Silvestre (2007)

## **MODE FOR GROUPED DATA: EXAMPLE**

#### Relative Frequencies

Class	Rel. Freq.	
0–10	0.10	
10–20	0.16	
20–30	0.24	
30–40	0.14	
40–50	0.06	

Modal class is the class interval with the highest absolute or relative frequency. Modal Class: 20 - 30

Step 1: Identify the modal class  $\rightarrow$  20–30

Step 2: Extract values for the formula:

• 
$$l = 20$$

• 
$$h = 10$$

• 
$$f^* = 0.16$$

• 
$$f^{**} = 0.14$$

Step 3: Apply King's formula:

$$\mathrm{Mo} = 20 + \frac{0.14}{0.16 + 0.14} \cdot 10 = 20 + \frac{0.14}{0.30} \cdot 10 = 20 + 4.67$$

Step 4: Result

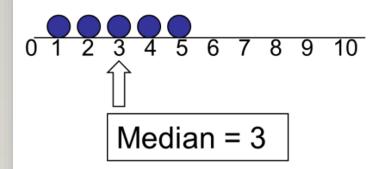
 $Mode \approx 24.67$ 

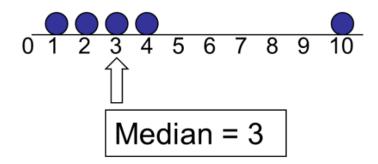


Mode ≈ 24.67

## **MEDIAN**

 In an ordered list, the median is the "middle" number (50% above, 50% below)





Not affected by extreme values

## FINDING THE MEDIAN

The location of the median:

Median position = 
$$\left(\frac{n+1}{2}\right)^{\text{th}}$$
 position in the ordered data

- If the number of values is odd, the median is the middle number
- If the number of values is even, the median is the average of the two middle numbers
- Note that  $\frac{n+1}{2}$  is not the value of the median, only the position of the median in the ranked data

# CALCULATING THE MEDIAN: EXAMPLES

#### Formula to Find Median Position

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

- n = number of observations
- If **Position is integer** → Median = value at that position
- If **Position is not integer** → Median = average of values at floor and ceil(Position)

#### **Examples**

#### Example 1 - Position is integer

• Data: 2, 4, 6, 8, 10 (n=5)



$$Position = \frac{5+1}{2} = 3$$
 (integer)

Median = 3rd value = 6

#### Example 2 - Position is not integer

Data: 3, 5, 8, 12, 15, 18 (n = 6)



$$Position = \frac{6+1}{2} = 3.5 \quad (not integer)$$

$$Median = \frac{3rd \ value + 4th \ value}{2} = \frac{8+12}{2} = 10$$

## MEDIAN FOR GROUPED DATA

#### Definition:

The **median class** is the class interval that contains the **middle value** (0.5 in cumulative relative frequency) of the distribution.

Formula (Grouped Data

$$\mathrm{Me} = l + rac{0.5 - F^*(l)}{F^*(L) - F^*(l)} \cdot h$$

#### Legend / Notation:

- ullet l = lower boundary of the median class
- L = upper boundary of the median class
- h = L l = class width
- ullet  $F^*(l)=$  cumulative relative frequency before the median class
- $F^*(L) =$  cumulative relative frequency at the upper boundary of the median class

Silvestre (2007)

## **MEDIAN FOR GROUPED DATA: EXAMPLE**

Relative Frequencies

**Cumulative Relative Frequencies** 

Median Class: 20 - 30

Class	Relative Frequency	Cumulative Rel. Freq.
0–10	0.10	0.10
10–20	0.16	0.26
20–30	0.24	0.50
30–40	0.14	0.64
40-50	0.06	0.70

The median class is the first class for which the cumulative relative frequency is equal to or greater than 0.5.

Step 1: Identify the median class  $\rightarrow$  20–30

Step 2: Extract values for the formula:

- l = 20. L = 30. h = 10
- ullet  $F^*(l)=0.26$  (cumulative relative frequency before the median class)
- $F^*(L) = 0.50$  (cumulative relative frequency at the upper boundary of the median class)

Step 3: Apply formula:

$$\mathrm{Me} = 20 + rac{0.5 - 0.26}{0.50 - 0.26} \cdot 10 = 20 + rac{0.24}{0.24} \cdot 10 = 20 + 10$$

Step 4: Result

Median  $\approx 30$ 

Median ≈ 30

## **ARITHMETIC MEAN**

The arithmetic mean (mean) is the most common measure of central tendency

– For a population of N values:

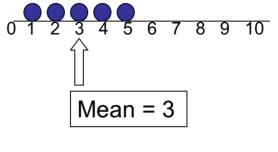
$$\mu = \frac{\sum_{i=1}^{N} x_i}{N} = \frac{x_1 + x_2 + \dots + x_N}{N}$$
Population values
Population size

– For a sample of size n:

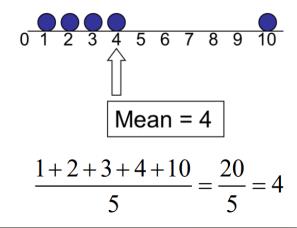
$$\overline{x} = \frac{\sum_{i=1}^{n} x_i}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}$$
 Observed values
Sample size

# ARITHMETIC MEAN: EXAMPLES

- The most common measure of central tendency
- Mean = sum of values divided by the number of values
- Affected by extreme values (outliers)



$$\frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$



## **TYPES OF MEANS**

#### Harmonic Mean

For a sample of n positive values  $x_1, x_2, \ldots, x_n$ :

$$H = rac{n}{\sum_{i=1}^n rac{1}{x_i}}$$

👉 It gives more weight to smaller values.

Typical use: averages of rates or speeds.

#### Geometric Mean

For a sample of n positive values  $x_1, x_2, \ldots, x_n$ :

$$G = \left(\prod_{i=1}^n x_i
ight)^{rac{1}{n}}$$

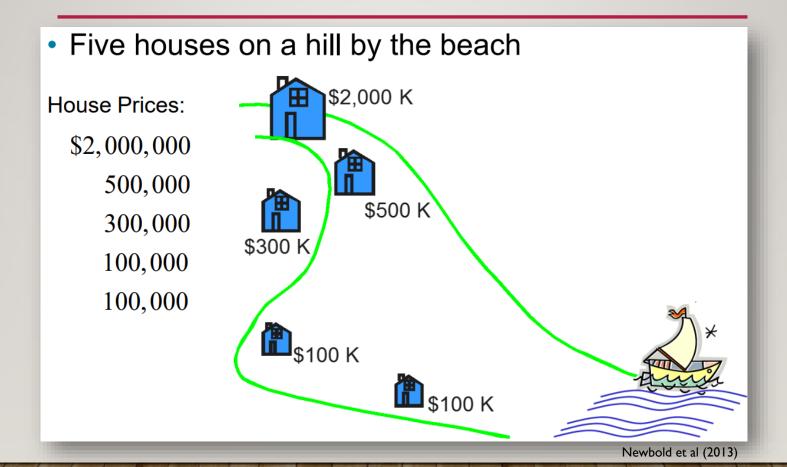
- 👉 Often used in population growth, compound interest, or growth rates.
- General comparison between means:

 $H \leq G \leq A$ 

# WHICH MEASURE OF LOCATION IS THE "BEST"?

- Mean is generally used, unless extreme values (outliers) exist ...
- Then median is often used, since the median is not sensitive to extreme values.
  - Example: Median home prices may be reported for a region – less sensitive to outliers

### **REVIEW EXAMPLE**



# REVIEW EXAMPLE: SUMMARY STATISTICS

#### House Prices:

\$2,000,000

500,000

300,000

100,000

100,000

Sum 3,000,000

• Mean:  $\left(\frac{\$3,000,000}{5}\right)$ 

= \$600,000

Median: middle value of ranked data

= \$300,000

Mode: most frequent value

= \$100,000

## **WEIGHTED MEAN**

#### Weighted Mean

The weighted mean of a set of data is

$$\overline{x} = \frac{\sum w_i x_i}{n}$$

where  $w_i$  = weight of the ith observation and  $n = \sum w_i$ .

### **WEIGHTED MEAN: EXAMPLE**

## Example 2.17 Stock Recommendation (Weighted Mean)

Zack's Investment Research is a leading investment research firm. Zack's will make one of the following recommendations with corresponding weights for a given stock: Strong Buy (1), Moderate Buy (2), Hold (3), Moderate Sell (4), or Strong Sell (5). Suppose that on a particular day, 10 analysts recommend Strong Buy, 3 analysts recommend Moderate Buy, and 6 analysts recommend Hold for a particular stock. Based on Zack's weights, find the mean recommendation.

Solution Table 2.8 shows the weights for each recommendation and the computation leading to a recommendation based on the following weighted mean recommendation conversion values: if the weighted mean is 1, Strong Buy; 1.1 through 2.0, Moderate Buy; 2.1 through 3.0, Hold; 3.1 through 4.0, Moderate Sell; 4.1 through 5, Strong Sell.

Table 2.8 Computation of Zack's Investment Research's Average Brokerage Recommendation

ACTION	Number of Analysts, $\boldsymbol{w}_i$	$V$ ALUE, $x_i$	$w_i x_i$
Strong Buy	10	1	10
Moderate Buy	3	2	6
Hold	6	3	18
Moderate Sell	0	4	0
Strong Sell	0	5	0

Newbold et al (2013)

$$\overline{x} = \frac{\sum_{i=1}^{n} w_i x_i}{n} = \frac{10 + 6 + 18 + 0 + 0}{19} = 1.79$$

The weighted mean of 1.79 yielded a Moderate Buy recommendation.

## **MEAN FOR GROUPED DATA**

### Definition:

The mean for grouped data is the weighted average of the class midpoints, weighted by the frequencies.

Formula (Grouped Data - Silvestre, 2007):

$$ar{x} = rac{\sum_{j=1}^m n_j \cdot x_j'}{n} = \sum_{j=1}^m f_j \cdot x_j'$$

## Legend / Notation:

- $j=1,2,...,m \rightarrow \text{class index}$
- ullet  $n_j=$  absolute frequency of class j
- ullet  $f_j=$  relative frequency of class j
- $x'_j = \text{midpoint of class } j$
- m = number of classes
- ullet  $n=\sum_{j=1}^m n_j o$  total number of observations

Silvestre (2007)

## **MEAN FOR GROUPED DATA: EXAMPLE**

Example Calculation:		Absolute Frequencies	Relative Frequencies
Class Interval	$x_j^\prime$ (Midpoint)	$n_{j}$	$f_{j}$
0–10	5	4	0.20
10–20	15	6	0.30
20–30	25	5	0.25
30–40	35	5	0.25

Step 1: Using absolute frequencies:

$$\bar{x} = \frac{(4\cdot 5) + (6\cdot 15) + (5\cdot 25) + (5\cdot 35)}{4+6+5+5} = \frac{410}{20} = 20.5$$

Step 2: Using relative frequencies:

$$ar{x} = (0.20 \cdot 5) + (0.30 \cdot 15) + (0.25 \cdot 25) + (0.25 \cdot 35) = 1 + 4.5 + 6.25 + 8.75 = 20.5$$

# PROPORTIONS: QUANTITIES AND FREQUENCIES

## Key Idea:

- A proportion represents the part of the total corresponding to a category.
- As a probability: A proportion can be interpreted as the likelihood of selecting an observation from that category, ranging from 0 to 1.
- As a percentage: A proportion can be expressed as a percentage of the total, ranging from 0% to 100%

#### Formulas:

1. Proportion / Probability:

## **Key Points:**

Proportions **sum to I** → total probability = I Percentages **sum to I00**%

$$ext{Percentage} = f_j \cdot 100$$

 $f_j = rac{n_j}{n}$ 

### Example:

Category	Quantity ( $n_j$ )	Proportion / Probability ( $f_j$ )	Percentage
А	10	0.25	25%
В	20	0.50	50%
С	10	0.25	25%

# FOUR PROPERTIES OF THE MEAN

Let m(x) or  $\bar{x}$  be the mean of variable x.

**Note:** The notation used here follows **Silvestre (2007)**.

- 1. Addition / Subtraction of a Constant:
- Adding a constant c to all values:

$$m(x+c) = m(x) + c$$

• Subtracting a constant c from all values:

$$m(x-c) = m(x) - c$$

- 2. Multiplication / Division by a Constant:
- Multiplying all values by a constant c:

$$m(c \cdot x) = c \cdot m(x)$$

• Dividing all values by a constant *c*:

$$m(x/c) = m(x)/c$$

Silvestre (2007)

# FOUR PROPERTIES OF THE MEAN

3. Mean of Deviations is Zero:

**Note:** The notation used here follows **Silvestre** (2007).

- $m(x-\bar{x})=m(x)-\bar{x}=0$
- The mean is the balance point of the data.
- 4. Mean of Grouped Values:
  - If the data are divided into groups  $G_1,G_2,...,G_k$  with group means  $\bar{x}_1,\bar{x}_2,...,\bar{x}_k$  and sizes  $n_1,n_2,...,n_k$ :

$$ar{x} = rac{n_1ar{x}_1 + n_2ar{x}_2 + ... + n_kar{x}_k}{n_1 + n_2 + ... + n_k}$$

Silvestre (2007)

# QUANTILES: DEFINITION (REVIEW)

- What are Quantiles?
- Quantiles are values that divide a dataset into equal parts.
- Special cases:
  - Quartiles → Q1, Q2, Q3, Q4 (divide data into 4 equal parts)
    - Median = Q2
  - Deciles → D1, D2, ..., D10 (divide data into 10 equal parts)
    - Median = D5
  - Percentiles → P1, P2, ..., P100 (divide data into 100 equal parts)
    - Median = P50

# QUANTILES FOR GROUPED DATA

$$q_p = l + rac{p - F^*(l)}{F^*(L) - F^*(l)} \cdot h$$

Silvestre (2007)

## Legend / Notes:

- $ullet q_p = extsf{p}$  -th quantile of grouped data
- ullet l= lower class boundary of the class containing  $q_p$
- ullet L= upper class boundary of that class
- $F^*(l) =$  cumulative relative frequency at l
- $F^*(L) =$  cumulative relative frequency at L
- p = quantile proportion (e.g., 0.25 for Q1)
- ullet h=L-l= class width of the class containing the quantile

# QUANTILE FOR GROUPED DATA: EXAMPLE

Suppose we have the following <b>frequency distribution</b> of exam scores:					
Score Interval	Frequency	Relative Frequencies	Cumulative Relative Frequencies		
0,10 7	2	0.10	0.10		
(10,20]	5	0.25	0.35		
(20,30]	8	0.40	0.75		
(30,40]	4	0.20	0.95		
(40,50]	1	0.05	1.00		
Question: Find the 1st quartile (Q1).					

## Step 1: Identify the class containing Q1

- Q1 corresponds to the **25th percentile**, so p=0.25.
- Look at the cumulative relative frequencies:

•

• (10,20]: 0.35 🔽

 $\rightarrow$  Q1 lies in the (10,20] class.

# QUANTILE FOR GROUPED DATA: EXAMPLE

Cupposouvo	have the	fallowing f		distribution	of over	
Suppose we	nave the	lollowing i	requency	distribution	oi exam so	Jores.

Question: Find the 1st quartile (Q1).

Score Interval	Frequency	Relative Frequency	Cumulative Relative Frequen
0,10 7	2	0.10	0.10
(10,20]	5	0.25	0.35
(20,30]	8	0.40	0.75
(30,40]	4	0.20	0.95
(40,50]	1	0.05	1.00

Step 2: Apply the formula

$$q_p = l + rac{p - F^*(l)}{F^*(L) - F^*(l)} \cdot h$$

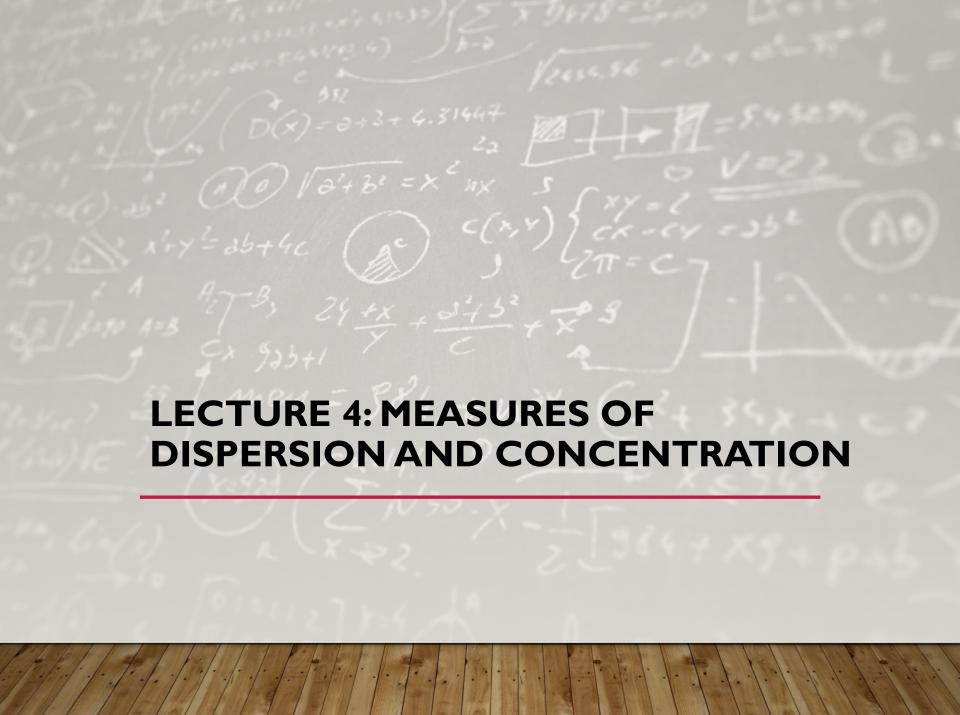
Where:

- $ullet \ l=10$  (lower boundary of the class)
- L=20 (upper boundary)
- $F^*(l) = 0.10$
- $F^*(L) = 0.35$
- p = 0.25
- h = L l = 10

$$Q1 = 10 + rac{0.25 - 0.10}{0.35 - 0.10} \cdot 10$$

Step 3: Calculate

$$Q1 = 10 + rac{0.15}{0.25} \cdot 10 = 10 + 0.6 \cdot 10 = 10 + 6 = 16$$



# MEASURES OF DISPERSION AND CONCENTRATION

- 1. Measures of Absolute Dispersion describe the spread of the data in absolute terms
  - Range (Amplitude de variação)
- Interquartile Range (IQR) (Amplitude inter-quartis)
- Mean Absolute Deviation (MAD) (Desvio médio absoluto)
- Variance and Standard Deviation
  - For ungrouped (raw) data
  - For grouped (classified) data
- 2. Measures of Relative Dispersion describe spread relative to a central value:
- Coefficient of Variation (CV)
- Quartile-based measure: (Q3 Q1) / Median
- 3. Measures of Concentration describe how data are focused or unequal:
- Gini Index

# MEASURES OF ABSOLUTE DISPERSION

1. Range (Amplitude de variação)

$$\mathrm{Range} = \max(x_i) - \min(x_i)$$

2. Interquartile Range (IQR) - Amplitude inter-quartis

$$IQR = Q3 - Q1$$

3. Mean Absolute Deviation (MAD) - Desvio médio absoluto (only for ungrouped data)

$$d = \frac{\sum_{i=1}^{n} |x_i - \bar{x}|}{n}$$

- 4. Variance Variância
- Ungrouped (raw) data:

$$s^2 = rac{\sum_{i=1}^n (x_i - ar{x})^2}{n} = rac{\sum_{i=1}^n x_i^2}{n} - ar{x}^2$$

### Notes:

- $x_i$  = individual data points (ungrouped)
- $x'_j$  = class midpoint (grouped)
- $\bar{x}$  = mean of the data
- $n_j$  = absolute frequency of the j-th class
- $f_j$  = relative frequency of the j-th class
- n = total number of observations

$$s=\sqrt{s^2}$$
 (standard deviation)

• Grouped (classified) data:

$$s^2 = rac{\sum_{j=1}^m n_j (x_j' - ar{x})^2}{n} = \sum_{j=1}^m f_j (x_j' - ar{x})^2 = rac{\sum_{j=1}^m n_j (x_j')^2}{n} - ar{x}^2$$

$$s = \sqrt{s^2}$$
 (standard deviation)

## ABSOLUTE DISPERSION MEASURES (UNGROUPED DATA): EXAMPLE

Data (scores):

$$12, 18, 20, 22, 25, 28, 30, 35, 38, 40$$

1. Range

$$\mathrm{Range} = \max(x_i) - \min(x_i) = 40 - 12 = 28$$

- 2. Interquartile Range (IQR) using position formula and interpolation
- n = 10
- Q1 position:  $(n+1) \cdot 0.25 = 2.75 \Rightarrow Q1 = 0.25 \cdot 18 + 0.75 \cdot 20 = 19.5$
- Q3 position:  $(n+1) \cdot 0.75 = 8.25 \Rightarrow Q3 = 0.75 \cdot 35 + 0.25 \cdot 38 = 35.75$

$$IQR = Q3 - Q1 = 16.25$$

3. Mean Absolute Deviation (MAD)

$$ar{x} = 26.8, \quad d = rac{\sum_{i=1}^{10} |x_i - ar{x}|}{10} pprox 7.3$$

4. Variance

$$s^2 = rac{(12-26.8)^2 + \cdots + (40-26.8)^2}{10} pprox 74.68$$

5. Standard Deviation

$$s=\sqrt{s^2}pprox 8.64$$

(standard deviation)

## ABSOLUTE DISPERSION MEASURES (GROUPED DATA): EXAMPLE

Data: Monthly Expenses of 20 Families (\$)				
Class Interval	$Midpoint\ x_j'$	Frequency $n_j$		
[0, 100]	50	2		
]100, 200]	150	5		
]200, 300]	250	8		
]300, 400]	350	4		
]400, 500]	450	1		

Step 1: Total observations

$$n = 2 + 5 + 8 + 4 + 1 = 20$$

Step 2: Range

 $Range = \max(upper\ limit) - \min(lower\ limit) = 500 - 0 = 500$ 

## ABSOLUTE DISPERSION MEASURES (GROUPED DATA): EXAMPLE

Data: Monthly Expenses of 20 Families (\$)		Absolute Frequen
Class Interval Midpoint $x_j^\prime$		Frequency $n_j$
[0, 100]	50	2
]100, 200]	150	5
]200, 300]	250	8
]300, 400]	350	4
]400, 500]	450	1

## Step 1: Total observations

$$n = 2 + 5 + 8 + 4 + 1 = 20$$

cies

Step 2: Mean

$$ar{x} = rac{\sum n_j x_j'}{n} = 235$$

Step 3: Variance

$$s^2 = rac{(50-235)^2 \cdot 2 + \dots + (450-235)^2 \cdot 1}{20} pprox 10750$$

Step 4: Standard Deviation

$$s=\sqrt{s^2}pprox 103.7$$

## PROPERTIES OF VARIANCE

#### 1. Variance of constants and shifts

$$v(c)=0,\quad v(x+c)=v(x),\quad v(x-c)=v(x)$$

• Adding or subtracting a constant does not change variance.

**Note:** The notation used here follows **Silvestre** (2007).

#### 2. Variance of scaled variables

$$v(cx)=c^2v(x), \quad v\left(rac{x}{c}
ight)=rac{v(x)}{c^2}, \quad c
eq 0$$

- Scaling by a factor c multiplies variance by  $c^2$ .
- 3. Variance decomposition for grouped data
- Suppose n observations divided into k groups  $g_1,g_2,\ldots,g_k$  with  $n_1,n_2,\ldots,n_k$  elements, such that  $n_1+n_2+\cdots+n_k=n$ .
- Let  $\bar{x}_1, \bar{x}_2, \ldots, \bar{x}_k$  be the group means and  $s_1^2, s_2^2, \ldots, s_k^2$  the group variances.

$$v_{ ext{total}} = \underbrace{\sum_{j=1}^k rac{n_j}{n} s_j^2}_{ ext{within-group}} + \underbrace{\sum_{j=1}^k rac{n_j}{n} (ar{x}_j - ar{x})^2}_{ ext{between-group}}$$

Note: This decomposition shows how total variability is explained by variability within groups and variability between groups.
Silvestre (2007)

## **RELATIVE DISPERSION MEASURES**

## 1. Coefficient of Variation (CV)

$$CV = rac{s}{ar{x}} imes 100\%$$

- Measures relative variability compared to the mean.
- Useful to compare variability across datasets with different units or scales.

## 2. Interquartile Range Ratio (IQR / Median)

$$\text{Relative IQR} = \frac{Q3 - Q1}{\text{Median}}$$

- Measures dispersion relative to the central value (median).
- Robust to extreme values (outliers).

### Notes:

- Both measures are dimensionless, allowing comparison between different datasets.
- CV is more sensitive to changes in the mean; Relative IQR is more robust.

Silvestrel (2007)

# RELATIVE DISPERSION MEASURES: EXAMPLE

2000, 2200, 2500, 2700, 2800, 3000, 3200, 3500

Step 1: Mean and Standard Deviation

$$\bar{x} = 2712.5, \quad s \approx 467.3$$

Step 2: Coefficient of Variation (CV)

$$CV = rac{s}{ar{x}} imes 100\% pprox 17.2\%$$

Step 3: Interquartile Range (IQR) - using interpolation formula

- Q1 position:  $(n+1) \cdot 0.25 = 9 \cdot 0.25 = 2.25$ 
  - ullet r=2 , lpha=0.25

$$Q1 = (1 - 0.25) \cdot x_2 + 0.25 \cdot x_3 = 0.75 \cdot 2200 + 0.25 \cdot 2500 = 2275$$

- Q3 position:  $(n+1) \cdot 0.75 = 6.75$ 
  - $r = 6, \alpha = 0.75$

$$Q3 = (1 - 0.75) \cdot x_6 + 0.75 \cdot x_7 = 0.25 \cdot 3000 + 0.75 \cdot 3200 = 3150$$

$$IQR = Q3 - Q1 = 3150 - 2275 = 875$$

Step 4: Relative IQR

$$ext{Relative IQR} = rac{ ext{IQR}}{ ext{Median}} = rac{875}{2750} pprox 0.318$$

## K-TH MOMENT ABOUT THE ORIGIN

- 1. k-th Moment about the Origin (Positive Integer k = 1, 2, ...)
  - For ungrouped data (simple data):

$$m_k' = \frac{1}{n} \sum_{i=1}^n x_i^k$$

For grouped data:

$$m_k' = rac{1}{n}\sum_{j=1}^m n_j(x_j')^k$$

where  $x_j^\prime$  = class midpoint,  $n_j$  = frequency, m = number of classes.

## Note:

- Raw moments are computed about the origin.
- They are used to calculate variance, skewness, and kurtosis

2. First Moment (k=1) - Mean

Ungrouped: 
$$m_1' = \frac{1}{n} \sum_{i=1}^n x_i = ar{x}$$

$$ext{Grouped: } m_1' = rac{1}{n} \sum_{i=1}^m n_j x_j' = ar{x}_{grouped}$$

3. Second Moment (k=2) – Raw Second Moment

Ungrouped: 
$$m_2' = \frac{1}{n} \sum_{i=1}^n x_i^2$$

Grouped: 
$$m_2' = \frac{1}{n} \sum_{j=1}^m n_j (x_j')^2$$

## K-TH CENTRAL MOMENT ABOUT THE MEAN

k-th Central Moment (k = 1, 2, ...)

Ungrouped data:

$$m_k = rac{1}{n} \sum_{i=1}^n (x_i - ar{x})^k$$

Grouped data:

$$m_k = rac{1}{n}\sum_{j=1}^m n_j (x_j' - ar{x})^k$$

where  $x_j'$  = class midpoint,  $n_j$  = frequency, m = number of classes.

## Note:

 Central moments measure variability around the mean, with the second central moment equal to the variance.

## Special cases:

• k = 1 (First central moment):

$$m_1 = 0$$

• k = 2 (Second central moment – Variance):

$$m_2=s^2=rac{1}{n}\sum (x_i-ar{x})^2 \quad ext{or} \quad rac{1}{n}\sum n_j(x_j'-ar{x})^2$$

## MEASURE OF CONCENTRATION: GINI INDEX & LORENZ CURVE

### 1. Gini Index (G) - Discrete version

$$G = rac{\sum_{i=1}^{m-1} (p_i - q_i)}{\sum_{i=1}^{m-1} p_i} = 1 - rac{\sum_{i=1}^{m-1} q_i}{\sum_{i=1}^{m-1} p_i}$$

#### Notation:

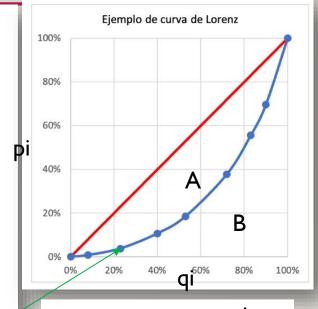
- $ullet p_i = rac{\sum_{j=1}^i n_j}{n} o$  cumulative proportion of **observations** up to class i
- $q_i = rac{\sum_{j=1}^i t_j}{\sum_{k=1}^m t_k}$  ightarrow cumulative proportion of the **variable** up to class i
- $t_j = n_j \cdot x_j' o$  total of the variable in class j ( $x_j'$  = class midpoint)
- m = total number of classes

### Conditions:

$$0 \le p_i \le 1$$
,  $0 \le q_i \le 1$ ,  $p_i \ge q_i$ ,  $i = 1, \ldots, m$ 

## Interpretation:

- G=0  $\rightarrow$  perfect equality
- ullet G=1 ightarrow maximum inequality



$$Gini\ Index = \frac{A}{A+B}$$

## 2. Lorenz Curve

- ullet Graph of  $q_i$  (cumulative share of variable) vs.  $p_i$  (cumulative population)
- ullet The area between the Lorenz curve and the line of equality is used to calculate G

# GINI INDEX FOR DATA GROUPED: EXAMPLE

Data			
Class (k€)	$n_{j}$	$x_j'$	$t_j=n_jx_j'$
0–20	10	10	100
20–40	15	30	450
40-60	20	50	1000
60-80	5	70	350

•  $n_i$ : class frequency

•  $x_i'$ : class midpoint

•  $t_j = n_j x_j'$  (class total income)

ullet  $N=\sum_{j}n_{j}$  (total population)

•  $T = \sum_{j} t_{j}$  (total income)

 $n_j/N$  (class population share)

 $t_{i}/T$  (class income share)

•  $P_i$  = (cumulative population share up to class i)

•  $Q_i$ : (cumulative income share up to class i)

## Totals:

$$N = 10 + 15 + 20 + 5 = 50$$

$$T = 100 + 450 + 1000 + 350 = 1900$$

## Shares and cumulative shares

Class	$n_j/N$	$t_j/T$	$P_i$ (cum.)	$Q_i$ (cum.)
0–20	0.20000	0.0526316	0.20000	0.0526316
20–40	0.30000	0.2368421	0.50000	0.2894737
40–60	0.40000	0.5263158	0.90000	0.8157895
60–80	0.10000	0.1842105	1.00000	1.0000000

(Values shown to 5–7 significant digits for clarity)

# GINI INDEX FOR DATA GROUPED: EXAMPLE

## Apply Silvestre's formula

Silvestre's version sums the cumulative shares **up to the penultimate class** (i.e. exclude the last row where  $P_m=Q_m=1$ ):

$$G = 1 - \frac{\sum_{i=1}^{m-1} Q_i}{\sum_{i=1}^{m-1} P_i}.$$

Compute the sums (exclude last class):

- $\sum_{i=1}^{m-1} P_i = 0.20000 + 0.50000 + 0.90000 = 1.60000$
- $\sum_{i=1}^{m-1} Q_i = 0.0526316 + 0.2894737 + 0.8157895 = 1.1578948$

Now

$$G = 1 - \frac{1.1578948}{1.60000} = 1 - 0.72368425 = 0.27631575 \approx \mathbf{0.27632}.$$

#### • G=0 $\rightarrow$ perfect equality (everyone has exactly the same income). • G=1 $\rightarrow$ maximum inequality (one person has all the income) ★ In practice: Researchers often classify Gini values into approximate ranges: Gini Value Typical Interpretation 0.00 - 0.20 Very low inequality (almost perfect equality, rarely observed in large societies). 0.20 - 0.30Low inequality. 0.30 - 0.40Moderate inequality 0.40 - 0.50High inequality. > 0.50 Very high (extreme) inequality.

# THANKS!

**Questions?**