

## Statistics Refresher I

18<sup>th</sup> JRC Annual training on Composite Indicators and Scoreboards

Elizabeth Casabianca

Joint Research Centre

## Notation

- An indicator: e.g. GDP  $\boldsymbol{\chi}$  *j*th value of an indicator: e.g. GDP for country *j*  $x_i$  *i*th indicator (in a group): e.g. group is GDP, life expectancy,  $\chi_i$ median income, ... • *j*th value, *i*th indicator: e.g. GDP for country *j*  $x_{i,i}$ N
  - European

e.g. number of countries

2021 JRC Week on Composite Indicators and Scoreboards

• Sample size

## Some data



Commission

### Rank



4 JRC-COIN © | 2021 JRC Week on Composite Indicators and Scoreboards

European Commission

#### Mean

$$\mu_i = E(x_i) = \frac{1}{N}(x_{i,1} + x_{i,2} + x_{i,3} + \dots + x_{i,N}) = \frac{1}{N}\sum_{j=1}^N x_{i,j}$$

 $\mu_i = \frac{1}{10} \left(9.6 + 10 + 9.3 + 7.1 + 5 + 8 + 7.5 + 8.2 + 6.4 + 1\right)$ 





## Median



The **middle value of the list of numbers**. Half of the values should be higher, half should be lower.

For an even number of values, it is the midway point between the two middle values.

The median is not equal to the mean unless it is a symmetric distribution.





### Percentiles

 $\chi_i$ 

5

6.4

Percentiles are a generalisation of the median.

Median: value such that 50% of values are below it

Commission

## xth percentile: value such that X% of values are below it



7 JRC-COIN © | 2021 JRC Week on Composite Indicators and Scoreboards

20th

percentile

## Variance

A measure of dispersion: how different are values from one another?

$$\operatorname{var}(x_i) = \sigma_i^2 = \frac{1}{N} \sum_{j=1}^N (\mu_i - x_{i,j})^2$$

The average squared distance of data points from the mean





## Variance and standard deviation



#### High variance

## $\sigma_i$ is the **standard deviation:** the square root of the variance



#### Low variance





**A histogram** shows how data is distributed into equally-spaced "bins".

Each bar gives an indication of how frequently it is that a value will fall inside that range.





A probability density function (pdf) shows the relative probability of different values occurring.





A probability density function (pdf) shows the relative probability of different values occurring.

The total area under the pdf (integral) is equal to 1 by definition.





A probability density function (pdf) shows the relative probability of different values occurring.

The total area under the pdf (integral) is equal to 1 by definition.

The area between two points shows the probability of a value falling in that range.





There are many theoretical types of probability distribution.

They aim to mimic real distributions observed in nature, economics, physics, etc.

**Discrete distributions** can only take a finite amount of values.

Continuous distributions can

take an infinite number of values (in a finite or infinite range).



## **Normal distribution**



The Normal (Gaussian) distribution is the simplest type of distribution!

- Values of a random variable fall into a (smooth) continuous curve with a **bell shaped**, **symmetric pattern**
- Main characteristics:
  - Mean=mode=median
  - Symmetric
  - Asymptotic



## **Skew and kurtosis**



**Skew** shows the extent to which values are clustered towards one end of the scale. Measure of lack of symmetry.

**Kurtosis** measures how heavy the tails of the distribution are, compared to a normal distribution.

**Be careful:** "kurtosis" may mean *true kurtosis*, or *excess kurtosis* (= true kurtosis – 3)



**Platykurtic** Kurtosis < 3 Excess kurtosis < 0



# Thank you



© European Union 2021

Unless otherwise noted the reuse of this presentation is authorised under the <u>CC BY 4.0</u> license. For any use or reproduction of elements that are not owned by the EU, permission may need to be sought directly from the respective right holders.

