Unit Two Descriptive Biostatistics

Dr. Hamza Aduraidi

Done by : farah mherat & Hamza abu mahmoud



Biostatistics

- What is the biostatistics? A branch of applied math. that deals with collecting, organizing and interpreting data using well-defined procedures.
- Types of Biostatistics:
 - Descriptive Statistics. It involves organizing, summarizing & displaying data to make them more understandable.
 - Inferential Statistics. It reports the degree of confidence of the sample statistic that predicts the value of the population parameter.



Descriptive Biostatistics

The best way to work with data is to summarize and organize them in a way that is meaningful and helpful for the next step

Numbers that have not been summarized and organized are called raw data.



Definition

- Data is any type of information
- Raw data is a data collected as they receive.
- Organize data is data organized either in ascending, descending or in a grouped data.

When we collect data from our samples, they are just meaningless data, and they can't take you anywhere and they don't provide you with any facts or info about the characteristic of your sample, so you need to organize them and then present them in a very meaningful way to describe them (and this is the descriptive measures).

Descriptive Measures

- Measures of Location (مقاییس التموضع)
 - Measures of central tendency:Mean; Median; Mode
 - Measures of non-central tendency Quantiles
 - Quartiles; Quintiles; Percentiles
- Measures of Dispersion (مقاييس التشنت)
 - Range
 - Interquartile range
 - Variance
 - Standard Deviation
 - Coefficient of Variation (CV)
- Measures of Shape
 - Mean > Median-positive or right Skewness
 - Mean = Median- symmetric or zero Skewness
 - Mean < Median-Negative of left Skewness</p>



Measures of Location notes

Measures of central tendency: measures the tendency of the sample's observations to be . located in the center (or around it) (على خط الاعداد) Measures of Non-central tendency: measures the tendency of the sample's observations to be located around interesting location other than the center or . the middle

Measures of Dispersion notes

- It do not measure the tendency of the sample's observations to be located in the center ,rather , it measures the tendency of the observations to be dispersed (distant) from the center .
- Standard Deviation is the most important one .

Measures of Shape notes

- If half of the observations is on the right of the center and the other half is on the left then this is a symmetrical distribution
- The observations may skewed to the right or left

Measures of Location

- It is a property of the data that they tend to be clustered about a center point.
- Measures of *central tendency* (i.e., central location) help find the approximate center of the dataset.
- Researchers usually do not use the term average, because there are three alternative types of average.
- These include the mean, the median, and the mode. In a perfect world, the mean, median & mode would be the same.
 - mean (generally not part of the data set)
 - median (may be part of the data set)
 - mode (always part of the data set)



Commonly Used Symbols

For a Sample

- x sample mean
- s² sample variance

Whenever you see symbols in English letter it's about a sample. Whenever you see symbols in Latin letter it's about population

- s sample standard deviation For a Population
 - μ population mean
 - population variance
 - 5 population standard deviation





Average

The sample mean is the sum of all the observations (ΣX_i) divided by the number of observations (n):

 $\overline{X} = \frac{\sum_{i=1}^{n} X_i}{n} \text{ where } \Sigma X_i = X_1 + X_2 + X_3 + X_4 + \dots + X_n$

Example. 1, 2, 2, 4, 5, 10. Calculate the mean. Note: n = 6 (six observations) Sample size

$$\Sigma X_i = 1 + 2 + 2 + 4 + 5 + 10 = 24$$

 $\overline{X} = 24 / 6 = 4.0$



General Formula--Population Mean



 $\mu = \text{population mean}$ $\Sigma = \text{summation sign}$ $x_i = \text{value of element i of the sample}$ N = population size

Notes on **Sample** Mean X

Formula



Note In sample size we use n In the population size we use N

Summation Sign

- In the formula to find the mean, we use the "summation sign" Σ
- This is just mathematical shorthand for "add up all of the observations"

$$\sum_{i=1}^n X_i = X_1 + X_2 + X_3 + \ldots + X_n$$



Notes on Sample Mean

- Also called *sample average* or *arithmetic mean*
- Mean for the sample = \overline{X} or M, Mean for population = mew (μ)
- Uniqueness: For a given set of data there is one and only one mean.
- Simplicity: The mean is easy to calculate.
- Sensitive to extreme values



Notes on Sample Mean

- The mean is used to fairly describe the observations (it gives the average)
- Despite it is a simple method to describe the sample's observations but it is sensitive to extreme values
- Example: if two students didn't study well for the exam and got 0 .. This will lower the mean of the marks
- So this will not be accurate for the of everybody's score.

The Mean

Example.

For the data: 1, 1, 1, 1, 51. Calculate the mean. Note: n = 5 (five observations)

$$\Sigma X_i = 1 + 1 + 1 + 1 + 51 = 55$$

 $\overline{X} = 55 / 5 = 11.0$

The mean (11) here was influenced by the extreme value (51)

Here we see that the mean is affected by extreme values.

Extra note : the mean is not part of the data set. It's unique to the sample.



الوسيط الحسابي The Median

- □ The median is the middle value of the *ordered* data
- To get the median, we must first rearrange the data into an ordered array (in ascending or descending order). Generally, we order the data from the lowest value to the highest value.
- Therefore, the median is the data value such that half of the observations are larger and half are smaller. It is also the 50th percentile.
- If n is odd, the median is the middle observation of the ordered array. If n is even, it is midway between the *two* central observations.



The Median

- it's the middle value of the ordered data, it is less influenced by extreme values. Also it's not a thing that you can calculate (like the mean), you find the median by FIRSTLY rearranging the data into an ordered array (ascending or descending, generally we order the data from lowest value to the highest value) then you find the number in the middle.(may be part of the data set)
- Notice that half of the data are larger and half are smaller than the median .

The Median

Note:

n is the number of the ordered observations not their value.

Example:

0 2 3 5 20 99 100

Note: Data has been ordered from lowest to highest. Since n is odd (n=7), the median is the (n+1)/2 ordered observation, or the 4th observation. Answer: The median is 5.

The mean and the median are unique for a given set of data. There will be exactly one mean and one median.

Unlike the mean, the median is not affected by extreme values.

Q: What happens to the median if we change the 100 to 5,000? Not a thing, the median will still be 5. Five is still the middle value of the data set.



The Median

Example:



Note: Data has been ordered from lowest to highest. Since n is even (n=6), the median is the (n+1)/2ordered observation, or the 3.5^{th} observation, *i.e.*, the average of observation 3 and observation 4, The median is 35, the mean is 35. (Notice that the mean=median, that indicate that the mean is a true mean; it is exactly in the middle)

Answer: The median is 35.



Mean and Median in a set of Salary Data XYZ Corp



22

More explanation

Every dot in the picture indicates a person's salary in a year, the first 5 salaries are 10 thousands, the 6th and 7th are almost 20 thousands, the 10th person earns 120 thousands in a year. Notice that the last one will affect the mean and the distribution of the data. While the median is between the 5^{th} and 6^{th} person = (10+20)/2 = 15 thousands. But wait, the Mean is almost 32 thousands, which means that it's extremely influenced by the last person.

The Mode المنوال

- The mode is the value of the data that occurs with the greatest frequency.
- Unstable index: values of modes tend to fluctuate from one sample to another drawn from the same population
- **Example**. 1, 1, 1, 2, 3, 4, 5

Answer. The mode is 1 since it occurs three times. The other values each appear only once in the data set.

Example. 5, 5, 5, 6, 8, 10, 10, 10.

Answer. The mode is: 5, 10.

There are two modes. This is a *bi-modal* dataset.



The Mode

The mode is different from the mean and the median in that those measures always exist and are always unique. For any numeric data set there will be one mean and one median.

□ The mode may not exist.

- Data: 1, 2, 3, 4, 5, 6, 7, 8, 9, 0
- Here you have 10 observations and they are all different so it is called **non-modal sample**.

□ The mode may not be unique.

- Data: 0, 1, 1, 2, 2, 3, 3, 4, 4, 5, 5, 6, 6, 7
- Mode = 1, 2, 3, 4, 5, and 6. There are six modes.



Comparison of the Mode, the Median, and the Mean

When you compare between those three values and they are the same this is a good thing because the distribution of this sample is symmetric



Notice that Right and Left sides are mirror images)

Comparison of the Mode, the Median, and the Mean

- In a normal distribution, the mode, the median, and the mean have the same value.
- The mean is the widely reported index of central tendency for variables measured on an interval and ratio scale.
- The mean takes each and every score into account.
- It also the most stable index of central tendency and thus yields the most reliable estimate of the central tendency of the population.



Comparison of the Mode, the Median, and the Mean

- The mean is always pulled in the direction of the long tail, that is, in the direction of the extreme scores.
- For the variables that positively skewed (like income), the mean is higher than the mode or the median. For negatively skewed variables (like age at death) the mean is lower.
- When there are extreme values in the distribution (even if it is approximately normal), researchers sometimes report means that have been adjusted for outliers.
- To adjust means one must discard a fixed percentage (5%) of the extreme values from either end of the distribution.



- If the mode and the median is lower or greater than the mean then there is Asymmetry (skewness).
- It can be positive or negative and it depends on the direction of that skewness.

Distribution Characteristics

- Mode: Peak(s)
- Median: Equal areas point
- Mean: Balancing point

This is symmetrical distribution

Normal distribution is the most important example on symmetrical distribution



27 December 2021

Shapes of Distributions

- Symmetric (Right and left sides are mirror images)
 - Left tail looks like right tail
 - Mean = Median = Mode





Shapes of Distributions



Shapes of Distributions

Mear

Left skewed (negatively skewed)

- Long left tail
- Mean < Median</p>

Here the Mode>Median>Mean.

This distribution has long left tale and is influenced with extreme low values

The mean is skewed to the left because of an extreme value



mean

Mode

Median



- Measures of non-central location used to summarize a set of data
- Examples of commonly used quantiles:
 - Quartiles
 - Quintiles
 - Deciles
 - Percentiles



Quartiles

- Quartiles split a set of ordered data into four parts.
 - Imagine cutting a chocolate bar into four equal pieces... How many cuts would you make? (yes, 3!)

\Box Q₁ is the First Quartile

- 25% of the observations are smaller than Q₁ and 75% of the observations are larger
- Q₂ is the Second Quartile
 - 50% of the observations are smaller than Q₂ and 50% of the observations are larger. Same as the Median. It is also the 50th percentile.

\square Q₃ is the Third Quartile

75% of the observations are smaller than Q₃ and 25% of the observations are larger



All the observations between the upper limit and lower limit are divided into 4 equal parts

Descriptive Statistics I

Other Quantiles

- Similar to what we just learned about quartiles, where 3 quartiles split the data into 4 equal parts,
 - There are 9 *deciles* dividing the distribution into 10 equal portions (tenths). کل جزء یحتوي على عُشر العينة
 - On the left of d1 there are 10% of the sample , on the right there are 90%
 - D2: on the left 20% , on the right 80% of the sample
 - There are four *quintiles* dividing the population into 5 equal portions. کل جزء يحتوي على خُمس العينة
 - On the left of the first quintile there are 20% of the sample (1/5 of the sample), on the right 80% (4/5 of the sample)



Other Quantiles

- ... and 99 *percentiles*
- كل جزء يسمى المئين 🗧
- On the left of the first percentile there are 1% of the sample
- On the left of the P10 there 10% of the sample so it's equals the d1
- P20=d2=first quintile
- P50=median= quartile 2=d5
- P75=quartile 3
- P90=d9
- P80=d8=quintile 4
- In all these cases, the convention is the same. The point, be it a quartile, decile, or *percentile*, takes the value of one of the observations or it has a value halfway between two adjacent observations.
 It is never necessary to split the difference

In all these cases, the convention is the same. The point, be it a quartile, decile, or *percentile*, takes the value of one of the observations or it has a value halfway between two adjacent observations. It is never necessary to split the difference between two observations more finely.

It refers to how spread out the scores are.

- In other words, how similar or different participants are from one another on the variable. It is either homogeneous or heterogeneous sample.
- Why do we need to look at measures of dispersion?

Again, disperion : away from each other , away from the centre



Measures of Dispersion

- We will study these five measures of dispersion
 - Range
 - Interquartile Range
 - Standard Deviation
 - Variance
 - Coefficient of Variation
 - Relative Standing.





The range is important thing to look at the difference between different samples

- Is the simplest measure of variability, is the difference between the highest score and the lowest score in the distribution.
- In research, the range is often shown as the minimum and maximum value, without the abstracted difference score.
- It provides a quick summary of a distribution's variability.
- It also provides useful information about a distribution when there are extreme values.
- The range has two values, it is highly unstable.





Range = Largest Value – Smallest Value
 Example: 1, 2, 3, 4, 5, 8, 9, 21, 25, 30
 Answer: Range = 30 – 1 = 29.

Pros:

- Easy to calculate
- Cons:
 - Value of range is only determined by two values
 - The interpretation of the range is difficult.
 - One problem with the range is that it is influenced by extreme values at either end.

make the range wider



- The standard deviation, s, measures a kind of "average" deviation about the mean. It is not really the "average" deviation, even though we may think of it that way. مثلا اذا كنا بنحكي عن شعبة فيعني كم بعد كل شخص عن آفريج الشعبة بشكل معياري
- Why can't we simply compute the average deviation about the mean, if that's what we want?

$$\sum_{i=1}^{n} (X_i - \overline{X})$$

If you take a simple mean, and then add up the deviations about the mean, as above, this sum will be equal to 0. Therefore, a measure of "average deviation" will not work.



Instead, we use:

$$s = \sqrt{\frac{\sum_{i=1}^{n} (X_i - \overline{X})^2}{n-1}}$$

- This is the "definitional formula" for standard deviation.
- The standard deviation has lots of nice properties, including:
 - By squaring the deviation, we eliminate the problem of the deviations summing to zero.
 - In addition, this sum is a minimum. No other value subtracted from X and squared will result in a smaller sum of the deviation squared. This is called the "least squares property."
- Note we divide by (n-1), not n. This will be referred to as a loss of one degree of freedom.



 The smaller the standard deviation, the better is the mean as the summary of a typical score.
 E.g. 10 people weighted 150 pounds, the SD would be zero, and the mean of 150 would communicate perfectly accurate information about all the participants wt. Another example would be a heterogeneous sample 5 people 100 pounds and another five people 200 pounds. The mean still 150, but the SD would be 52.7.



Example. Two data sets, X and Y. Which of the two data sets has greater variability? Calculate the standard deviation for each.

We note that both sets of data have the same mean:

	A i	• i :
$\overline{X} = 3$ This doesn't	1	0
$\overline{\mathbf{v}} = \mathbf{z}$ two samples are	2	0
I — J the same	3	0
	4	5
(continued)		10



$$S_{X} = \sqrt{\frac{10}{4}} = 1.58$$



$$S_{Y} = \sqrt{\frac{80}{4}} = 4.47$$

This mean that sample x is less diverse than sample Y

U	3	3-	9
0	3	3-	9
0	3	3-	9
5	3	2	4
10	3	7	49
		2=0	80

[Check these results with your calculator.]

Standard Deviation gives you a better image about the distribution of the sample



Variance

The variance, s^2 , is the standard deviation (*s*) squared. Conversely, $s = \sqrt{variance}$.

Definitional formula:
$$s^2 = \frac{\sum_{i=1}^{n} (X_i - \overline{X})^2}{n-1}$$

Computational formula: $s^2 = \frac{\sum_{i=1}^{n} (X_i)^2 - \frac{(\sum_{i=1}^{n} X_i)^2}{n}}{n-1}$

Variance is an important measure in the T test

This is what computer software (e.g., MS Excel or your calculator key) uses.



Coefficient of Variation (CV)

- The problem with s² and s is that they are both, like the mean, in the "original" units.
- This makes it difficult to compare the variability of two data sets that are in different units or where the magnitude of the numbers is very different in the two sets. For example,
 - Suppose you wish to compare two stocks and one is in dollars and the other is in yen; if you want to know which one is more volatile, you should use the coefficient of variation.
 - It is also not appropriate to compare two stocks of vastly different prices even if both are in the same units.
 - The standard deviation for a stock that sells for around \$300 is going to be very different from one with a price of around \$0.25.

The *coefficient of variation* will be a better measure of dispersion in these cases than the standard deviation (see example on the next slide). $CV = \frac{s}{\overline{v}}(100\%)$



Coefficient of Variation (CV)

$$CV = \frac{s}{\overline{X}}(100\%)$$

CV is in terms of a percent. What we are in effect calculating is what percent of the sample mean is the standard deviation. If CV is 100%, this indicates that the sample mean is equal to the sample standard deviation. This would demonstrate that there is a great deal of variability in the data set. 200% would obviously be even Huge amount of worse. variability

If CV= 50% means less variability which is good





a measure of dispersion that depends on the measure of location

- The Interquartile range (IQR) is the score at the 75th percentile or 3rd quartile (Q3) minus the score at the 25th percentile or first quartile (Q1). Are the most used to define outliers.
- It is not sensitive to extreme values.



How to find Q1 and Q3



1.Find the median

2. consider the part on the right side a new small sample, the same thing on the left part

3.Find the median of sample 1 so it's a Q1 The median of the sample 2 is Q3

Inter-Quartile Range (IQR)



Five Number Summary

When examining a distribution for shape, sometime the five number summary is useful: Smallest | Q1 | Median | Q3 | Largest



5-number summary: 2 | 8 | 10 | 16.5 | 63 This data is right-skewed.

In right-skewed distributions, the distance from Q_3 to $X_{largest}$ (16.5 to 63) is significantly greater than the distance from $X_{smallest}$ to $Q_1(2$ to 8).

They difference between the median and lower limit is less than the difference between the median and the upper limit so it's right-skewed

