

L3: Measures of Center

The median and the mean apply are only quantitative data, whereas the mode can be used with either quantitative or qualitative data.

4.1 The Mode

The **sample mode** of a qualitative or a discrete quantitative variable occurs with the greatest frequency in a data set.

Definition 4.1 (Mode). Obtain the frequency of each observed value of the variable in a data and note the greatest frequency.

1. If the greatest frequency is 1 (i.e. no value occurs more than once), then the variable has no mode.
2. If the greatest frequency is 2 or greater, then any value that occurs with that greatest frequency is called a sample mode of the variable.

To obtain the mode(s) of a variable, we first construct a frequency distribution for the data using classes based on single value. The mode(s) can then be determined easily from the frequency distribution.

Example 4.1. Let us consider the frequency table for blood types of 40 persons

BLOOD			
		Statistics	
BLOOD		Frequency	Percent
Valid	O	16	40.0
	A	18	45.0
	B	4	10.0
	AB	2	5.0
	Total	40	100.0

We can see from frequency table that the mode of blood types is A.

When we measure a continuous variable (or discrete variable having a lot of different values) such as height or weight of person, all the measurements may be different. In such a case there is no mode because every observed value has frequency 1. However, the data can be grouped into class intervals and the mode can then be defined in terms of class frequencies. With grouped quantitative variable, the mode class is the class interval with highest frequency.

Example 4.2. Let us consider the frequency table for prices of hot-dogs

Frequencies of prices of hotdogs (\$/oz.)				
		Frequency	Percent	Cumulative Percent
Valid	0.045-0.065	5	9.3	9.3
	0.065-0.085	15	27.8	37.0
	0.085-0.105	10	18.5	55.6
	0.105-0.125	9	16.7	72.2
	0.125-0.145	4	7.4	79.6
	0.145-0.165	2	3.7	83.3
	0.165-0.185	3	5.6	88.9
	0.185-0.205	3	5.6	94.4
	0.205-0.225	1	1.9	96.3
	0.225-0.245	1	1.9	98.1
	0.245-0.265	1	1.9	100.0
	Total	54	100.0	

4.2 The Median

The median is a "central" value – there are as many values greater than it as there are less than it. The value that divides the set of observed values in half, so that the observed values in one half are less than or equal to the median value and the other half are greater or equal to the median value.

To obtain the median of the variable, we arrange observed values in a data set in increasing order and then determine the middle value in the ordered list.

1. If the number of observation is odd, then the sample median is the observed value exactly in the middle of the ordered list.
2. If the number of observation is even, then the sample median is the number halfway between the two middle observed values in the ordered list.

-In both cases, if we let n denote the number of observations in a data set, then the sample median is at position $\frac{n+1}{2}$ *in the ordered list.* in the ordered list.

Example 7 participants in bike race had the following finishing times in minutes:

28,22,26,29,21,23,24.

What is the median?

Example 8 participants in bike race had the following finishing times

in minutes: 28,22,26,29,21,23,24,50.

What is the median?

The median in SPSS: Analyze -> Descriptive Statistics -> Frequencies

4.3 The Mean

The most commonly used measure of center for quantitative variable. It is the sum of observed values in a data divided by the number of observations.

The mean=

sum of value
Number of value

 or $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ or $\frac{\sum x_i}{n}$.

Example 7 participants in bike race had the following finishing times in minutes:
28,22,26,29,21,23,24.

What is the mean?

Example 8 participants in bike race had the following finishing times in minutes:
28,22,26,29,21,23,24,50.

What is the mean?

The mean in SPSS:

Analyze -> Descriptive Statistics -> Frequencies

Analyze -> Descriptive Statistics -> Descriptive

L4: Measures of variation

Measures of variation used mostly for quantitative variables.

1-Range

The sample range of the variable is the difference between its maximum and minimum values in a data set:

$$\text{Range} = \text{Max} - \text{Min}.$$

There are several characteristics of range:

- 1- The sample range of the variable is quite easy to compute.
- 2- In using the range, a great deal of information ignored, that is, only the largest and smallest values of the variable are considered; the other observed values are disregarded.
- 3- The range cannot ever decrease, but can increase, when additional observations are included in the data set and that in sense the range is overly sensitive to the sample size.

Example 5.3. Prices of hotdogs (\$/oz.):

0.11, 0.17, 0.11, 0.15, 0.10, 0.11, 0.21, 0.20, 0.14, 0.14, 0.23, 0.25, 0.07,
0.09, 0.10, 0.10, 0.19, 0.11, 0.19, 0.17, 0.12, 0.12, 0.12, 0.10, 0.11, 0.13,
0.10, 0.09, 0.11, 0.15, 0.13, 0.10, 0.18, 0.09, 0.07, 0.08, 0.06, 0.08, 0.05,
0.07, 0.08, 0.08, 0.07, 0.09, 0.06, 0.07, 0.08, 0.07, 0.07, 0.07, 0.08, 0.06,
0.07, 0.06

The range in SPSS:

Analyze -> Descriptive Statistics -> Frequencies,

Analyze -> Descriptive Statistics -> Descriptives

Range of the prices of hotdogs

	N	Range	Minimum	Maximum
Price (\$/oz)	54	.20	.05	.25
Valid N (listwise)	54			

Example 5.2. 7 participants in bike race had the following finishing times in minutes:
28, 22, 26, 29, 21, 23, 24.

What is the range?

Example 5.3. 8 participants in bike race had the following finishing times in minutes:
28, 22, 26, 29, 21, 23, 24, 50.

What is the range?

2- Interquartile range

Before we can define the sample interquartile range, we have to first define the percentiles, the deciles and the quartiles of the variable in a data set.

1- The percentiles of the variable divide observed values into **hundredths**, or 100 equal parts.

The first percentile, P1, is the number that divides the bottom 1% of the observed values from the top 99%; **second percentile, P2**, is the number that divides the bottom 2% of the observed values from the top 98%; and so forth.

The median is the 50th percentile.

2- The deciles of the variable divide the observed values into **tenths**, or 10 equal parts. The variable has nine deciles, denoted by D1, D2, . . . , D9. **The first decile D1 is 10th percentile**, the second decile D2 is the 20th percentile, and so forth.

3- The most commonly used percentiles are quartiles. The quartiles of the variable divide the observed values **into quarters**, or **4 equal parts**.

The variable has three quartiles, denoted by Q1, Q2 and Q3: Arrange the observed values of variable in a data in **increasing order**.

1. The first quartile Q1 is at position $n+1/4$

2. The second quartile Q2 (the median) is at position $n+1/2$

3. The third quartile Q_3 is at position $\frac{3(n+1)}{4}$ in the ordered list.

4- The sample interquartile range of the variable(IQR) :is the difference between the first and third quartiles of the variable, that is:

$$IQR = Q_3 - Q_1$$

IQR gives the range of the middle 50% of the observed values.

Example 5.4. 7 participants in bike race had the following finishing times in minutes:

28,22,26,29,21,23,24.

What is the interquartile range?

Example 5.5. 8 participants in bike race had the following finishing times in minutes:

28,22,26,29,21,23,24,50.

What is the interquartile range?

Analyze -> Descriptive Statistics -> Explore

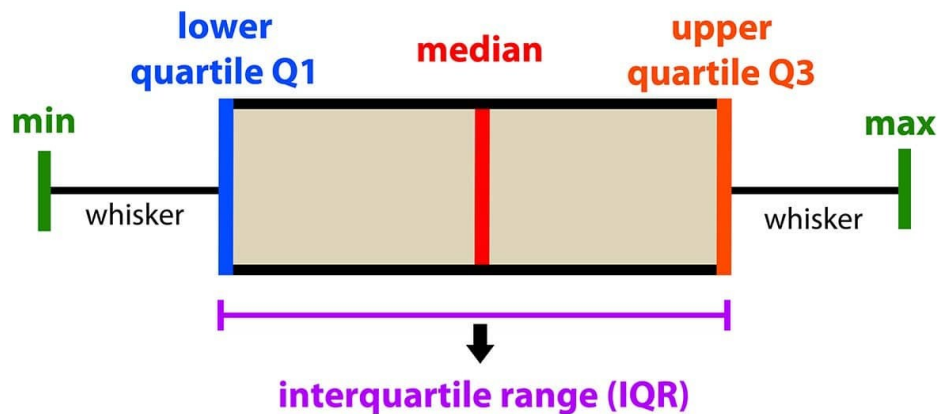
Noticeable: Five-number summary and boxplots

The five-number summary of the variable consists of minimum, maximum, and quartiles written in increasing order:

Min, Q_1 , Q_2 , Q_3 , Max.

A boxplot is based on the five-number summary and can be used to provide a graphical display of the center and variation of the observed values of variable in a data set.

introduction to data analysis: Box Plot



Standard deviation –SD

The sample standard deviation is the most frequently used measure of variability, although it is not as easily understood as ranges. It can be considered as a kind of average of the absolute deviations of observed values from the mean of the variable in question.

Standard deviation is a statistic that measures the dispersion of a dataset relative to its mean and is calculated as the square root of the variance.

$$s_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

$$SD = \sqrt{\frac{\sum |x - \mu|^2}{N}}$$

where:

\sum means "sum of"

x is a value in the data set

μ is the mean of the data set

N is the number of data points in the population.

Since the standard deviation is defined using the sample mean μ of the variable x .

Note :

1- the standard deviation is always positive number, i.e., $SD \geq 0$.

In a formula of the standard deviation, the sum of the squared deviations from the mean.

2- The formula above is for finding the standard deviation of a population. If you're dealing with a sample, you'll want to use a slightly different formula (below), which uses ***n-1*** instead of *N*.

$$SD_{\text{sample}} = \sqrt{\frac{\sum |x - \bar{x}|^2}{n - 1}}$$

Step 1: Find the mean. **Step 1: Finding μ in** $\sqrt{\frac{\sum |x - \mu|^2}{N}}$

Step 2: For each data point, find the square of its distance to the mean.

Step 2: Finding $|x - \mu|^2$

In this step, we find the distance from each data point to the mean (i.e., the deviations) and square each of those distances.

Step 3: Sum the values from Step 2. **Step 3: Finding $\sum |x - \mu|^2$**

Step 4: Divide by the number of data points. **Step 4: Finding $\frac{\sum |x - \mu|^2}{N}$**

Step 5: Take the square root.

Step 5: Finding the standard deviation

$$\sqrt{\frac{\sum |x - \mu|^2}{N}}$$

EX: find the SD of these values (6, 3, 2, 1)

Step 1: Find the mean μ .

$$\mu = \frac{6 + 2 + 3 + 1}{4} = \frac{12}{4} = 3$$

Step 2: Find the square of the distance from each data point to the mean

$$|x - \mu|^2.$$

x	$ x - \mu ^2$
6	$ 6 - 3 ^2 = 3^2 = 9$
2	$ 2 - 3 ^2 = 1^2 = 1$
3	$ 3 - 3 ^2 = 0^2 = 0$
1	$ 1 - 3 ^2 = 2^2 = 4$

Steps 3, 4, and 5:

$$\begin{aligned} \text{SD} &= \sqrt{\frac{\sum |x - \mu|^2}{N}} \\ &= \sqrt{\frac{9 + 1 + 0 + 4}{4}} \end{aligned}$$

$$= \sqrt{\frac{14}{4}} \quad \text{Sum the squares of the distances (Step 3).}$$

$$= \sqrt{3.5} \quad \text{Divide by the number of data points (Step 4).}$$

$$\approx 1.87 \quad \text{Take the square root (Step 5).}$$

Q- Find the standard deviation of the data set.

1, 4, 7, 2, 6