



STATISTICS

Adam Klepáč

September 21, 2023

MEAN – MEDIAN – DEVIATION – CORRELATION

USEFUL VALUES

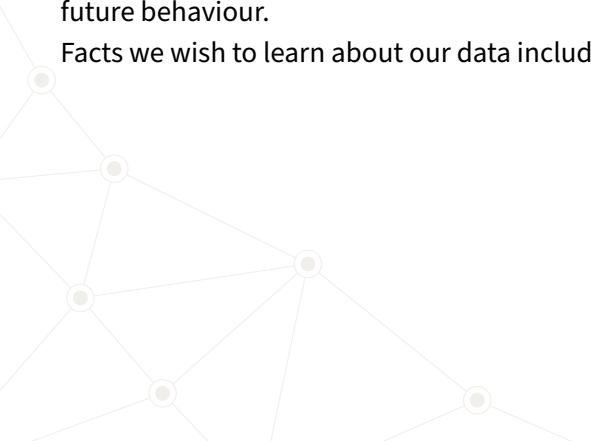
When studying data, there are **certain numerical values** which prove useful in predicting future behaviour.



USEFUL VALUES

When studying data, there are **certain numerical values** which prove useful in predicting future behaviour.

Facts we wish to learn about our data include:



USEFUL VALUES

When studying data, there are **certain numerical values** which prove useful in predicting future behaviour.

Facts we wish to learn about our data include:

- the expected value of the next experiment (the **mean**),

USEFUL VALUES

When studying data, there are **certain numerical values** which prove useful in predicting future behaviour.

Facts we wish to learn about our data include:

- the expected value of the next experiment (the **mean**),
- the 'middle' value of the outputs regardless of proportion (the **median**),

USEFUL VALUES

When studying data, there are **certain numerical values** which prove useful in predicting future behaviour.

Facts we wish to learn about our data include:

- the expected value of the next experiment (the **mean**),
- the 'middle' value of the outputs regardless of proportion (the **median**),
- the expected measure of difference of observed values from the mean (the **deviation**),

USEFUL VALUES

When studying data, there are **certain numerical values** which prove useful in predicting future behaviour.

Facts we wish to learn about our data include:

- the expected value of the next experiment (the **mean**),
- the 'middle' value of the outputs regardless of proportion (the **median**),
- the expected measure of difference of observed values from the mean (the **deviation**),
- dependence on any other data (the **correlation**).

1

THE MEAN



TYPES OF MEAN

The **mean** in some sense 'the most probable' next output based on the received data.

TYPES OF MEAN

The **mean** in some sense 'the most probable' next output based on the received data. What is 'the most probable' output however depends heavily on the type of experiment we are performing.

TYPES OF MEAN – ARITHMETIC MEAN

ARITHMETIC MEAN

The **arithmetic mean** is the sum of outputs divided by their number. If x_1, \dots, x_n are the outputs, their arithmetic mean (often denoted \bar{x}) is

$$\bar{x} := \frac{x_1 + x_2 + \dots + x_n}{n}.$$

TYPES OF MEAN – ARITHMETIC MEAN

ARITHMETIC MEAN

The **arithmetic mean** is the sum of outputs divided by their number. If x_1, \dots, x_n are the outputs, their arithmetic mean (often denoted \bar{x}) is

$$\bar{x} := \frac{x_1 + x_2 + \dots + x_n}{n}.$$

Most useful when dealing with data in **absolute** proportion.

TYPES OF MEAN – ARITHMETIC MEAN

ARITHMETIC MEAN

The **arithmetic mean** is the sum of outputs divided by their number. If x_1, \dots, x_n are the outputs, their arithmetic mean (often denoted \bar{x}) is

$$\bar{x} := \frac{x_1 + x_2 + \dots + x_n}{n}.$$

Most useful when dealing with data in **absolute** proportion.

Meaning that we're interested in '**how much**' is one output smaller/larger than another.

ARITHMETIC MEAN – EXAMPLE

The arithmetic mean is for most experiments the relevant one.



ARITHMETIC MEAN – EXAMPLE

The arithmetic mean is for most experiments the relevant one.

Consider for example an experiment tailored to determine the average height of a 15-year-old British male.



ARITHMETIC MEAN – EXAMPLE

The arithmetic mean is for most experiments the relevant one.

Consider for example an experiment tailored to determine the average height of a 15-year-old British male.

While comparing the heights of two people, we care about the **absolute** difference in centimetres.

ARITHMETIC MEAN – EXAMPLE

The arithmetic mean is for most experiments the relevant one.

Consider for example an experiment tailored to determine the average height of a 15-year-old British male.

While comparing the heights of two people, we care about the **absolute** difference in centimetres.

For example, if this is our data

Input	1	2	3	4	5
Output	165	161	164	172	168,

we conclude that the expected height of a randomly chosen 15-year-old British male is

$$\frac{165 + 161 + 164 + 172 + 168}{5} = 166.$$

TYPES OF MEAN – GEOMETRIC MEAN

GEOMETRIC MEAN

The **geometric mean** is the n -th root of the product of n outputs. That is, if x_1, \dots, x_n are the outputs, their geometric mean is

$$\bar{x} := \sqrt[n]{(x_1 \cdot x_2 \cdot \dots \cdot x_n)}.$$

TYPES OF MEAN – GEOMETRIC MEAN

GEOMETRIC MEAN

The **geometric mean** is the n -th root of the product of n outputs. That is, if x_1, \dots, x_n are the outputs, their geometric mean is

$$\bar{x} := \sqrt[n]{(x_1 \cdot x_2 \cdot \dots \cdot x_n)}.$$

Most useful when dealing with data in **relative** proportion.

TYPES OF MEAN – GEOMETRIC MEAN

GEOMETRIC MEAN

The **geometric mean** is the n -th root of the product of n outputs. That is, if x_1, \dots, x_n are the outputs, their geometric mean is

$$\bar{x} := \sqrt[n]{(x_1 \cdot x_2 \cdot \dots \cdot x_n)}.$$

Most useful when dealing with data in **relative** proportion.

Meaning that we're interested in '**how many times**' is one output smaller/larger than another.

GEOMETRIC MEAN – EXAMPLE

Suppose we're comparing the increase in population in Asian countries since the year 2000.



GEOMETRIC MEAN – EXAMPLE

Suppose we're comparing the increase in population in Asian countries since the year 2000.

When comparing two populations, we don't really care *by how much* they differ *how many times* is one larger than the other.

GEOMETRIC MEAN – EXAMPLE

Suppose we're comparing the increase in population in Asian countries since the year 2000.

When comparing two populations, we don't really care *by how much* they differ *how many times* is one larger than the other.

If this is our data

Input	India	China	Japan	South Korea	Mongolia	Taiwan
Output	1.328	1.118	0.991	1.100	1.366	1.078

This means that the expected increase in population in a randomly chosen Asian country is

$$\sqrt[6]{(1.328 \cdot 1.118 \cdot 0.991 \cdot 1.100 \cdot 1.366 \cdot 1.078)} = 1.156.$$

TYPES OF MEAN – HARMONIC MEAN

HARMONIC MEAN

The **harmonic mean** is the reciprocal of the sum of reciprocals divided by their number. If x_1, \dots, x_n are the outputs, their harmonic mean is

$$\bar{x} := \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}.$$

TYPES OF MEAN – HARMONIC MEAN

HARMONIC MEAN

The **harmonic mean** is the reciprocal of the sum of reciprocals divided by their number. If x_1, \dots, x_n are the outputs, their harmonic mean is

$$\bar{x} := \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}.$$

Most useful when dealing with **rates** and **ratios**.

TYPES OF MEAN – HARMONIC MEAN

HARMONIC MEAN

The **harmonic mean** is the reciprocal of the sum of reciprocals divided by their number. If x_1, \dots, x_n are the outputs, their harmonic mean is

$$\bar{x} := \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}.$$

Most useful when dealing with **rates** and **ratios**.

Meaning when comparing outputs which are actually ratios of two numbers.

HARMONIC MEAN – EXAMPLE

We study the speed of a train between individual stations.



HARMONIC MEAN – EXAMPLE

We study the speed of a train between individual stations.

If this is our data

Input	1 → 2	2 → 3	3 → 4	4 → 5	5 → 6	6 → 7
Output	65 km/h	52 km/h	71 km/h	60 km/h	62 km/h	53 km/h,

then the average speed of the train across the whole track is

$$\frac{6}{\frac{1}{65} + \frac{1}{52} + \frac{1}{71} + \frac{1}{60} + \frac{1}{62} + \frac{1}{53}} = 59.78 \text{ km/h.}$$

HARMONIC MEAN – EXAMPLE

We study the speed of a train between individual stations.

If this is our data

Input	1 → 2	2 → 3	3 → 4	4 → 5	5 → 6	6 → 7
Output	65 km/h	52 km/h	71 km/h	60 km/h	62 km/h	53 km/h,

then the average speed of the train across the whole track is

$$\frac{6}{\frac{1}{65} + \frac{1}{52} + \frac{1}{71} + \frac{1}{60} + \frac{1}{62} + \frac{1}{53}} = 59.78 \text{ km/h.}$$

Actually, here the arithmetic mean is 60.5 km/h which is not just an *inadequate* estimate, it's simply **the wrong answer!**

HARMONIC MEAN – EXAMPLE

We study the speed of a train between individual stations.

If this is our data

Input	1 → 2	2 → 3	3 → 4	4 → 5	5 → 6	6 → 7
Output	65 km/h	52 km/h	71 km/h	60 km/h	62 km/h	53 km/h,

then the average speed of the train across the whole track is

$$\frac{6}{\frac{1}{65} + \frac{1}{52} + \frac{1}{71} + \frac{1}{60} + \frac{1}{62} + \frac{1}{53}} = 59.78 \text{ km/h.}$$

Actually, here the arithmetic mean is 60.5 km/h which is not just an *inadequate* estimate, it's simply **the wrong answer!**

If you summed up all the distances between stations and divided them by the total time, you would get the **harmonic mean!**

THE MEDIAN

MEDIAN

The **median** is the value which lies exactly in the middle of a dataset. It is essentially the value separating the lower and upper half of outputs. If x_1, \dots, x_n are the outputs, the median is

$$\text{median}(x) := \begin{cases} x_{(n+1)/2} & \text{if } n \text{ is odd,} \\ \frac{x_{n/2} + x_{n/2+1}}{2} & \text{if } n \text{ is even.} \end{cases}$$

THE MEDIAN – EXAMPLE

The median is useful when we deal with data featuring radical extremes.



THE MEDIAN – EXAMPLE

The median is useful when we deal with data featuring radical extremes.
It is very often used in relationship to **location**.

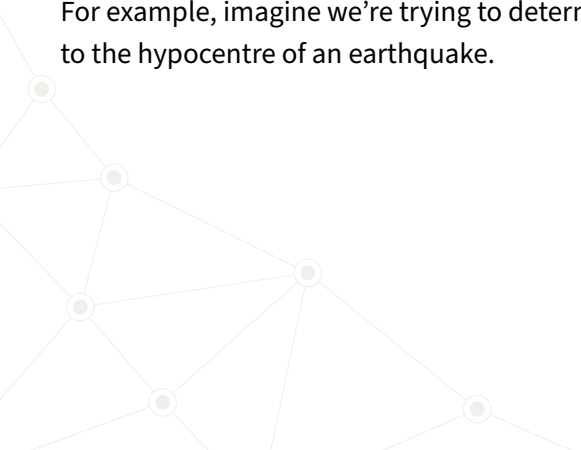


THE MEDIAN – EXAMPLE

The median is useful when we deal with data featuring radical extremes.

It is very often used in relationship to **location**.

For example, imagine we're trying to determine the distance from our measuring station to the hypocentre of an earthquake.



THE MEDIAN – EXAMPLE

The median is useful when we deal with data featuring radical extremes.

It is very often used in relationship to **location**.

For example, imagine we're trying to determine the distance from our measuring station to the hypocentre of an earthquake.

We can detect where the quake is strongest, giving us this data:

Input	1	2	3	4	5	6
Output	1 km	2 km	2 km	2 km	3 km	14 km

THE MEDIAN – EXAMPLE

The median is useful when we deal with data featuring radical extremes.

It is very often used in relationship to **location**.

For example, imagine we're trying to determine the distance from our measuring station to the hypocentre of an earthquake.

We can detect where the quake is strongest, giving us this data:

Input	1	2	3	4	5	6
Output	1 km	2 km	2 km	2 km	3 km	14 km

The median of this dataset is 2 km which is a much better estimate of a 'centre' than for example the arithmetic mean, being equal to 4, is.

THE MEDIAN – EXAMPLE

The median is useful when we deal with data featuring radical extremes.

It is very often used in relationship to **location**.

For example, imagine we're trying to determine the distance from our measuring station to the hypocentre of an earthquake.

We can detect where the quake is strongest, giving us this data:

Input	1	2	3	4	5	6
Output	1 km	2 km	2 km	2 km	3 km	14 km

The median of this dataset is 2 km which is a much better estimate of a 'centre' than for example the arithmetic mean, being equal to 4, is.

Also, the mean and the median cannot be 'too far' apart and the median requires at most two values to calculate, making it a very resource efficient approximation of the mean.

DEVIATION

DEVIATION

Deviation is a measure of the difference between the observed outputs and the computed mean.

DEVIATION

DEVIATION

Deviation is a measure of the difference between the observed outputs and the computed mean.

There are many types of deviations, we'll focus on two of them:

DEVIATION

DEVIATION

Deviation is a measure of the difference between the observed outputs and the computed mean.

There are many types of deviations, we'll focus on two of them:
the **standard deviation** (a measure of 'dispersion'),

DEVIATION

DEVIATION

Deviation is a measure of the difference between the observed outputs and the computed mean.

There are many types of deviations, we'll focus on two of them:

- the **standard deviation** (a measure of 'dispersion'),
- the **average absolute deviation** (a measure of actual 'difference').

DEVIATION

DEVIATION

Deviation is a measure of the difference between the observed outputs and the computed mean.

There are many types of deviations, we'll focus on two of them:

- the **standard deviation** (a measure of 'dispersion'),
- the **average absolute deviation** (a measure of actual 'difference').

A **very important distinction** is that the *standard deviation* concerns **future** measurements while the *average absolute deviation* concerns **past** measurements.

TYPES OF DEVIATION – STANDARD DEVIATION

STANDARD DEVIATION

The **standard deviation** measures the dispersion of a set of values. Basically, it measures how likely the data is to concentrate around the mean. If x_1, \dots, x_n are the outputs and \bar{x} is their **arithmetic** mean, then their standard deviation is

$$\sigma := \sqrt{\frac{1}{n}((x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2)}.$$

STANDARD DEVIATION – EXAMPLE

Let us repeat the height experiment. We measured the heights of 5 15-year-old British males to try to determine the national average. This is the data:

Input	1	2	3	4	5
Output	165	161	164	172	168

STANDARD DEVIATION – EXAMPLE

Let us repeat the height experiment. We measured the heights of 5 15-year-old British males to try to determine the national average. This is the data:

Input	1	2	3	4	5
Output	165	161	164	172	168

We computed the arithmetic mean to be 166. This means that the standard deviation of this data is

$$\sigma = \sqrt{\frac{1}{5}((165 - 166)^2 + (161 - 166)^2 + (164 - 166)^2 + (172 - 166)^2 + (168 - 166)^2)}$$

$$= 3.742,$$

meaning we can expect most new values to concentrate 3.742 cm around 166 cm.

TYPES OF DEVIATION – AVERAGE ABSOLUTE DEVIATION



AVERAGE ABSOLUTE DEVIATION

The **average absolute deviation** is the average of the absolute deviations from a chosen central point (typically the mean). If x_1, \dots, x_n are the outputs and \bar{x} is the chosen central point, then the average absolute deviation of this dataset is

$$\frac{|x_1 - \bar{x}| + |x_2 - \bar{x}| + \dots + |x_n - \bar{x}|}{n}.$$

AVERAGE ABSOLUTE DEVIATION – EXAMPLE

If we return to the height experiment yet again, we can calculate that the average absolute deviation of the data (with the central point being the arithmetic mean)

Input	1	2	3	4	5
Output	165	161	164	172	168

AVERAGE ABSOLUTE DEVIATION – EXAMPLE

If we return to the height experiment yet again, we can calculate that the average absolute deviation of the data (with the central point being the arithmetic mean)

Input	1	2	3	4	5
Output	165	161	164	172	168

$$\frac{|165 - 166| + |161 - 166| + |164 - 166| + |172 - 166| + |168 - 166|}{5} = 3.2,$$

AVERAGE ABSOLUTE DEVIATION – EXAMPLE

If we return to the height experiment yet again, we can calculate that the average absolute deviation of the data (with the central point being the arithmetic mean)

Input	1	2	3	4	5
Output	165	161	164	172	168

is

$$\frac{|165 - 166| + |161 - 166| + |164 - 166| + |172 - 166| + |168 - 166|}{5} = 3.2,$$

meaning that the measured heights differ on average by 3.2 cm from the calculated arithmetic mean.