



STATISTICS

Adam Klepáč

October 10, 2023

WHAT EVEN IS STATISTICS?

STATISTICS

Statistics is a mathematical discipline concerned with predicting future state of a system based *solely* on its past behaviour.

WHAT EVEN IS STATISTICS?

STATISTICS

Statistics is a mathematical discipline concerned with predicting future state of a system based *solely* on its past behaviour.

The collective information about a system's past state is called **data**.

WHAT EVEN IS STATISTICS?

STATISTICS

Statistics is a mathematical discipline concerned with predicting future state of a system based *solely* on its past behaviour.

The collective information about a system's past state is called **data**.
It assigns **probabilities** to each possible future state of system based on data.

WHAT EVEN IS STATISTICS?

STATISTICS

Statistics is a mathematical discipline concerned with predicting future state of a system based *solely* on its past behaviour.

The collective information about a system's past state is called **data**.
It assigns **probabilities** to each possible future state of system based on data.
It also assigns probabilities to the **possibility of wrong prediction**.

EXAMPLE – BIASED COIN?

We throw a coin 10 times with the following outcome:

$$\{H, H, H, T, H, T, H, H, H, T\},$$

H for ‘heads’, *T* for ‘tails’.

EXAMPLE – BIASED COIN?

We throw a coin 10 times with the following outcome:

$$\{H, H, H, T, H, T, H, H, H, T\},$$

H for 'heads', T for 'tails'. We can ask two questions:

EXAMPLE – BIASED COIN?

We throw a coin 10 times with the following outcome:

$$\{H, H, H, T, H, T, H, H, H, T\},$$

H for ‘heads’, T for ‘tails’. We can ask two questions:

- What is the probability that the **next toss** will come out ‘heads’/‘tails’?

EXAMPLE – BIASED COIN?

We throw a coin 10 times with the following outcome:

$$\{H, H, H, T, H, T, H, H, H, T\},$$

H for ‘heads’, T for ‘tails’. We can ask two questions:

- What is the probability that the **next toss** will come out ‘heads’/‘tails’?
- Is this coin is **biased towards** ‘heads’/‘tails’ with *allowed probability of error* α ?

EXAMPLE – BIASED COIN?

We throw a coin 10 times with the following outcome:

$$\{H, H, H, T, H, T, H, H, H, T\},$$

H for ‘heads’, T for ‘tails’. We can ask two questions:

- What is the probability that the **next toss** will come out ‘heads’/‘tails’?
 - We got 7 heads out of 10 tosses, so the probability for the next toss being heads is $7/10$.
- Is this coin is **biased towards** ‘heads’/‘tails’ with *allowed probability of error* α ?

EXAMPLE – BIASED COIN?

We throw a coin 10 times with the following outcome:

$$\{H, H, H, T, H, T, H, H, H, T\},$$

H for ‘heads’, T for ‘tails’. We can ask two questions:

- What is the probability that the **next toss** will come out ‘heads’/‘tails’?
 - We got 7 heads out of 10 tosses, so the probability for the next toss being heads is $7/10$.
- Is this coin is **biased towards** ‘heads’/‘tails’ with *allowed probability of error* α ?
 - **No**, for $\alpha = 0.05$.
 - **Yes**, for $\alpha = 0.2$.

CONTENTS

Data

Types of Data

Visualizing Discrete Data

Mean – Median – Deviation – Correlation

The Mean

The Median

The Deviation

Correlation

Frequency Distribution

DATA

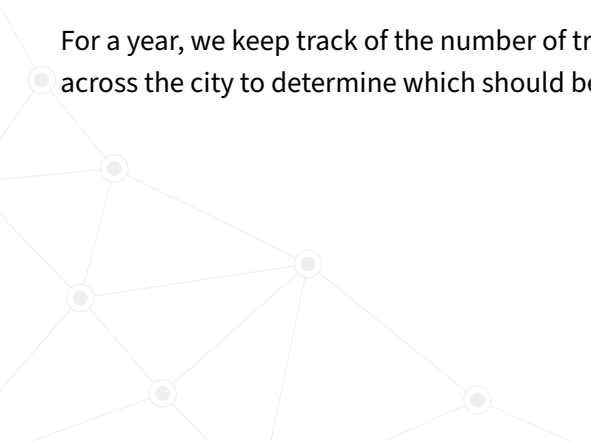
WHAT DO WE MEAN BY DATA?

DATA

Sets (called *inputs* and *outputs*) describing the studied system. There is typically only one set of inputs and possibly multiple sets of outputs.

EXAMPLE – JUNCTIONS

For a year, we keep track of the number of traffic accidents per day on road junctions across the city to determine which should be first replaced by roundabouts.



EXAMPLE – JUNCTIONS

For a year, we keep track of the number of traffic accidents per day on road junctions across the city to determine which should be first replaced by roundabouts. An **input** is a day in a year coupled with the location of the junction.

EXAMPLE – JUNCTIONS

For a year, we keep track of the number of traffic accidents per day on road junctions across the city to determine which should be first replaced by roundabouts.

An **input** is a day in a year coupled with the location of the junction.

An **output** is the number of traffic accidents in the given day on the given junction.

EXAMPLE – FIRST BABY

We study the age that women bear children for the first time across Europe.

EXAMPLE – FIRST BABY

We study the age that women bear children for the first time across Europe.

An **input** would be a name of a European country.

An **output** is the average age of a first-time mother in that country.

1

TYPES OF DATA



DISCRETE DATA VS. CONTINUOUS DATA

DISCRETE DATA

We call a data **discrete** if the set of *inputs* (and therefore also that of *outputs*) is **countable**.

DISCRETE DATA VS. CONTINUOUS DATA

DISCRETE DATA

We call a data **discrete** if the set of *inputs* (and therefore also that of *outputs*) is **countable**.

Both previous examples feature **discrete** data.

- There are only *finitely many* junctions in a city and days in a year.

DISCRETE DATA VS. CONTINUOUS DATA

DISCRETE DATA

We call a data **discrete** if the set of *inputs* (and therefore also that of *outputs*) is **countable**.

Both previous examples feature **discrete** data.

- There are only *finitely many* junctions in a city and days in a year.
- There are only *finitely many* countries on a continent.

DISCRETE DATA VS. CONTINUOUS DATA

CONTINUOUS DATA

We call a data **continuous** if the set of inputs is **uncountable**. In this case, the data is actually a **function**: set of inputs \rightarrow set of outputs.

DISCRETE DATA VS. CONTINUOUS DATA

CONTINUOUS DATA

We call a data **continuous** if the set of inputs is **uncountable**. In this case, the data is actually a **function**: set of inputs \rightarrow set of outputs.

More often than not, the inputs in a continuous data are **moments in time** or **coordinates in space**.

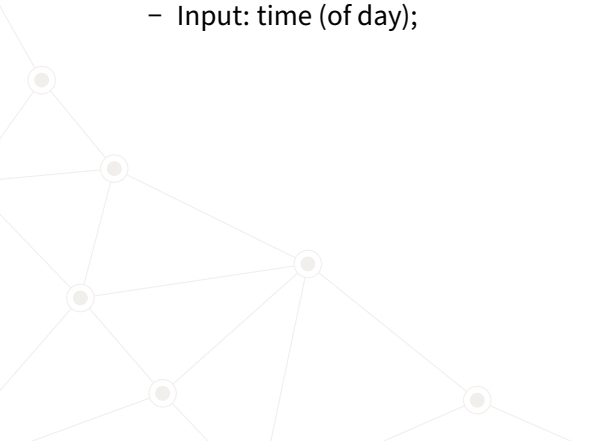
CONTINUOUS DATA – EXAMPLES

- We study the number of trains in a railway station at any given time.



CONTINUOUS DATA – EXAMPLES

- We study the number of trains in a railway station at any given time.
 - Input: time (of day);



CONTINUOUS DATA – EXAMPLES

- We study the number of trains in a railway station at any given time.
 - Input: time (of day);
 - Output: number of trains in the station.

CONTINUOUS DATA – EXAMPLES

- We study the number of trains in a railway station at any given time.
 - Input: time (of day);
 - Output: number of trains in the station.
 - The data is a function $f : [0, 24] \rightarrow \mathbb{N}$.

CONTINUOUS DATA – EXAMPLES

- We study the number of trains in a railway station at any given time.
 - Input: time (of day);
 - Output: number of trains in the station.
 - The data is a function $f : [0, 24] \rightarrow \mathbb{N}$.
- Another example is the density of air per cubic meter.

CONTINUOUS DATA – EXAMPLES

- We study the number of trains in a railway station at any given time.
 - Input: time (of day);
 - Output: number of trains in the station.
 - The data is a function $f : [0, 24] \rightarrow \mathbb{N}$.
- Another example is the density of air per cubic meter.
 - Input: Coordinates of a unit cube in space.

CONTINUOUS DATA – EXAMPLES

- We study the number of trains in a railway station at any given time.
 - Input: time (of day);
 - Output: number of trains in the station.
 - The data is a function $f : [0, 24] \rightarrow \mathbb{N}$.
- Another example is the density of air per cubic meter.
 - Input: Coordinates of a unit cube in space.
 - Output: The combined weight of air molecules.

CONTINUOUS DATA – EXAMPLES

- We study the number of trains in a railway station at any given time.
 - Input: time (of day);
 - Output: number of trains in the station.
 - The data is a function $f : [0, 24] \rightarrow \mathbb{N}$.
- Another example is the density of air per cubic meter.
 - Input: Coordinates of a unit cube in space.
 - Output: The combined weight of air molecules.
 - The data is a function $f : \mathbb{R}^3 \rightarrow [0, \infty)$.

VISUALIZING DISCRETE DATA

The background features abstract geometric shapes. A light teal triangle points downwards from the left edge. A dark blue triangle points upwards from the bottom right corner. These two triangles overlap in the center, creating a darker teal intersection. The top half of the image is a solid light gray.

TABLES



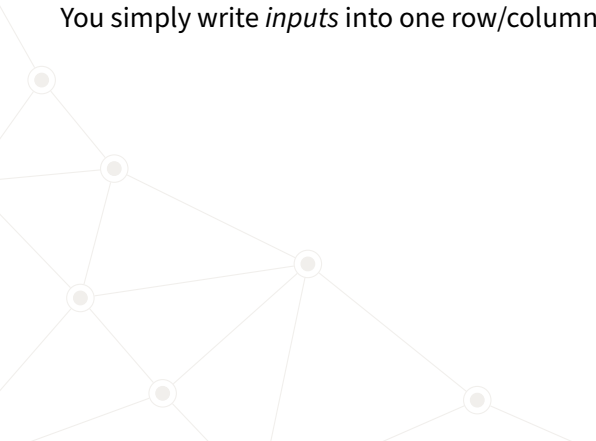
The simplest possible visualization.



TABLES

The simplest possible visualization.

You simply write *inputs* into one row/column and *outputs* into the other.



TABLES

The simplest possible visualization.

You simply write *inputs* into one row/column and *outputs* into the other.

For example, suppose you measure the height of 10 random people. You can visualize your experiment like this:

The simplest possible visualization.

You simply write *inputs* into one row/column and *outputs* into the other.

For example, suppose you measure the height of 10 random people. You can visualize your experiment like this:

Input	1	2	3	4	5	6	7	8	9	10
Output	180	169	191	177	175	181	171	153	180	183

PIE CHART

Only usable if your outputs **total a predetermined number**, typically *percentages*.



PIE CHART

Only usable if your outputs **total a predetermined number**, typically *percentages*.

Suppose we have three inputs – I_1 , I_2 and I_3 – with three outputs – 30%, 10% and 60%.

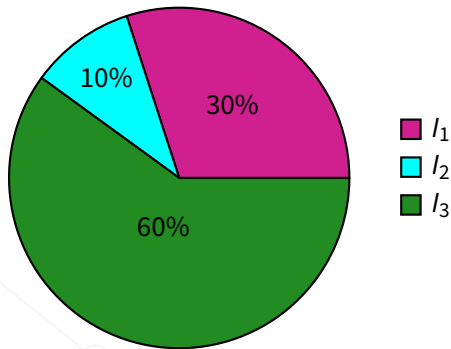


PIE CHART

Only usable if your outputs **total a predetermined number**, typically *percentages*.

Suppose we have three inputs – I_1 , I_2 and I_3 – with three outputs – 30%, 10% and 60%.

Pie chart of this data looks like this



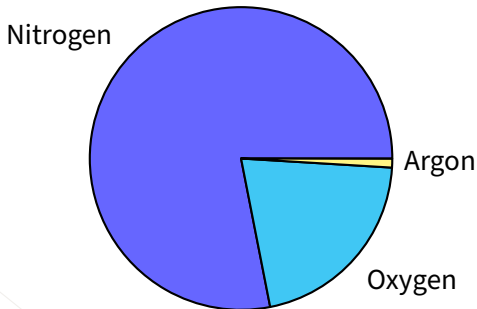
PIE CHART – EXAMPLES

Pie charts are frequently used to represent compositions of chemicals.



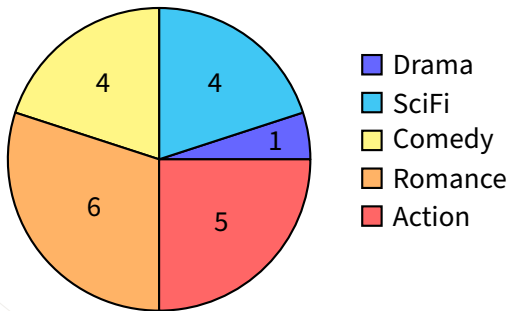
PIE CHART – EXAMPLES

Pie charts are frequently used to represent compositions of chemicals. For instance, here is a pie chart of the composition of *air*.



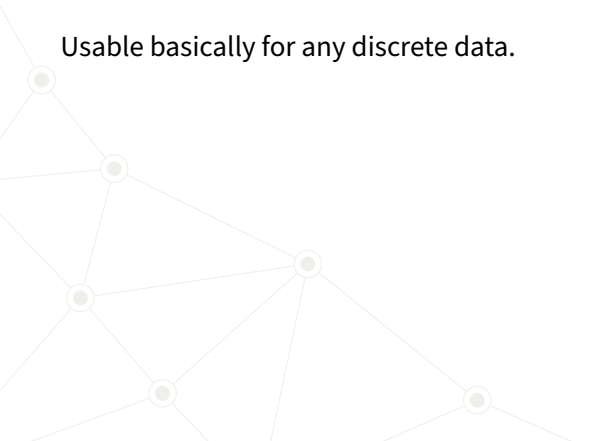
PIE CHART – EXAMPLES

Favourite type of movie as determined by a survey.



BAR CHART

Usable basically for any discrete data.



BAR CHART

Usable basically for any discrete data.

Especially useful when your inputs are ordered and when you expect the data to follow a certain trend – it can be easily approximated by a polygonal curve.

BAR CHART

Usable basically for any discrete data.

Especially useful when your inputs are ordered and when you expect the data to follow a certain trend – it can be easily approximated by a polygonal curve.

Also very good for comparing more outputs for the same inputs.

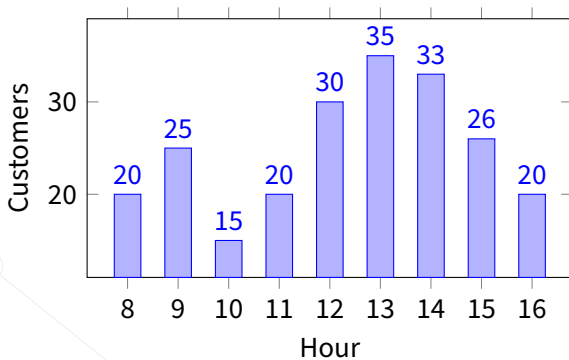
BAR CHART – EXAMPLE

Suppose we count the number of customers in our shop over each hour. If we're open from 8 AM to 5 PM, a bar chart of such an experiment can look like this:



BAR CHART – EXAMPLE

Suppose we count the number of customers in our shop over each hour. If we're open from 8 AM to 5 PM, a bar chart of such an experiment can look like this:



BAR CHART – EXAMPLE

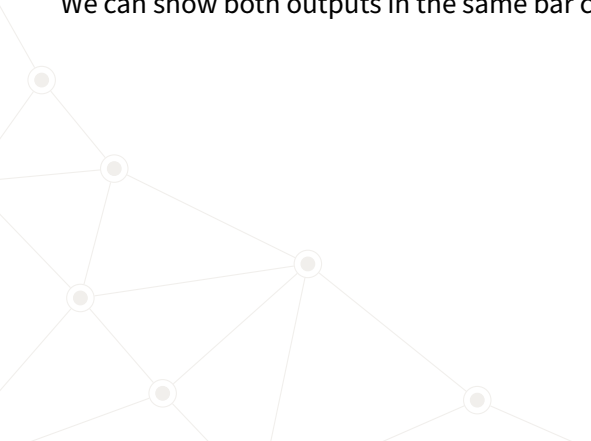
Let's say we open another shop and want to compare how the two shops are doing at each hour.



BAR CHART – EXAMPLE

Let's say we open another shop and want to compare how the two shops are doing at each hour.

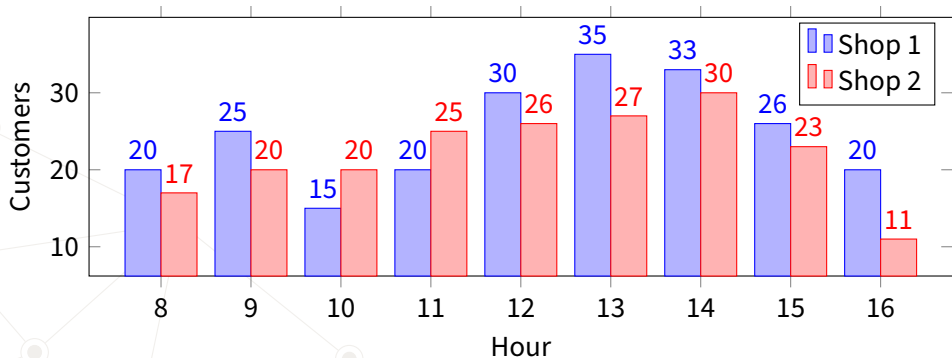
We can show both outputs in the same bar chart:



BAR CHART – EXAMPLE

Let's say we open another shop and want to compare how the two shops are doing at each hour.

We can show both outputs in the same bar chart:



BAR CHART – EXAMPLE

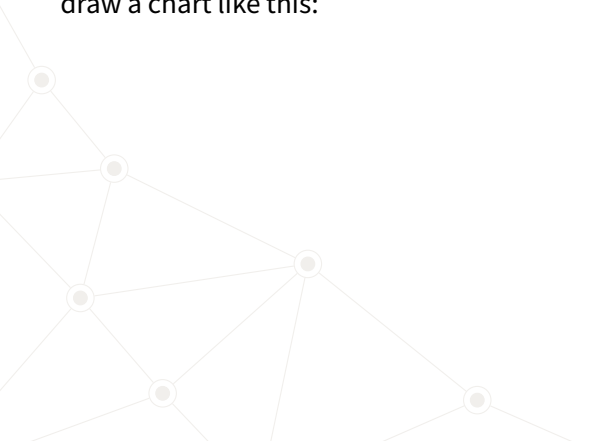
You can also use bar chart to stack outputs on top of each other.



BAR CHART – EXAMPLE

You can also use bar chart to stack outputs on top of each other.

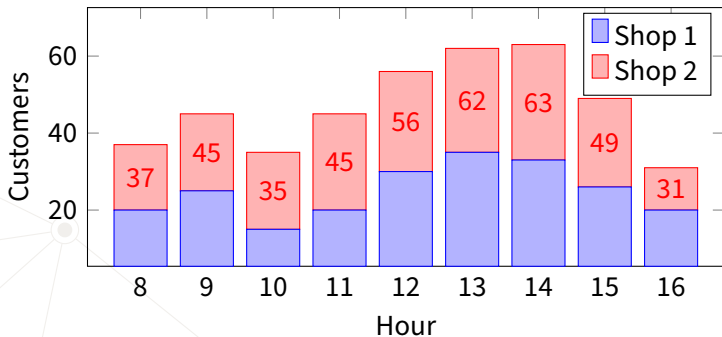
For example, if I wanted to know the **total** number of customers in both my shops, I could draw a chart like this:



BAR CHART – EXAMPLE

You can also use bar chart to stack outputs on top of each other.

For example, if I wanted to know the **total** number of customers in both my shops, I could draw a chart like this:



SCATTER PLOT

Scatter plots are useful when studying 'random' data.



SCATTER PLOT

Scatter plots are useful when studying 'random' data.

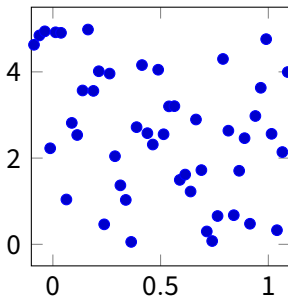
Something like the position of an air molecule in a box over time.



SCATTER PLOT

Scatter plots are useful when studying 'random' data.

Something like the position of an air molecule in a box over time.



Here, the x-axis represents time (0 to 1s) and the y-axis represents one coordinate of the molecule (say the box is a 5x5x5 cube).

SCATTER PLOT – EXAMPLE

Of course, you can also display multiple outputs with the same inputs in a scatter plot.



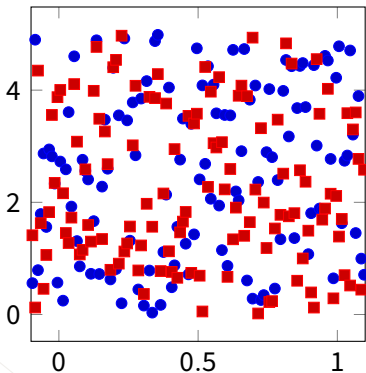
SCATTER PLOT – EXAMPLE

Of course, you can also display multiple outputs with the same inputs in a scatter plot.
Let's add another air molecule.



SCATTER PLOT – EXAMPLE

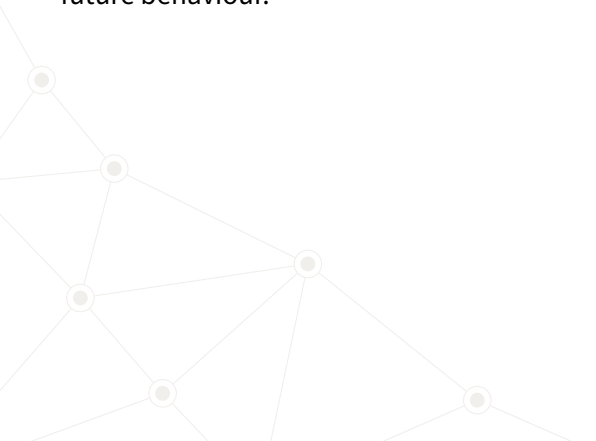
Of course, you can also display multiple outputs with the same inputs in a scatter plot.
Let's add another air molecule.



MEAN – MEDIAN – DEVIATION – CORRELATION

USEFUL VALUES

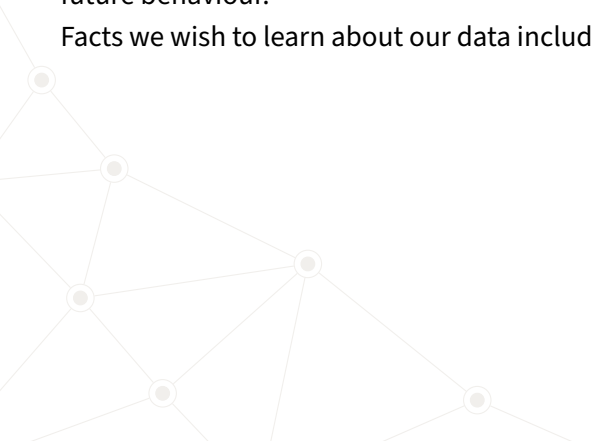
When studying data, there are **certain numerical values** which prove useful in predicting future behaviour.



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/observed 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/observed 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 is 'the most probable' next output based on the received data.

TYPES OF MEAN

The **mean** in some sense is '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, but *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, but *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

then 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 over 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 over 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 over 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!**

2

THE MEDIAN



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 **ordered from least to greatest**, 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 dealing with data manifesting radical extremes.



THE MEDIAN – EXAMPLE

The median is useful when dealing with data manifesting radical extremes.
It is very often used in relationship to **location**.

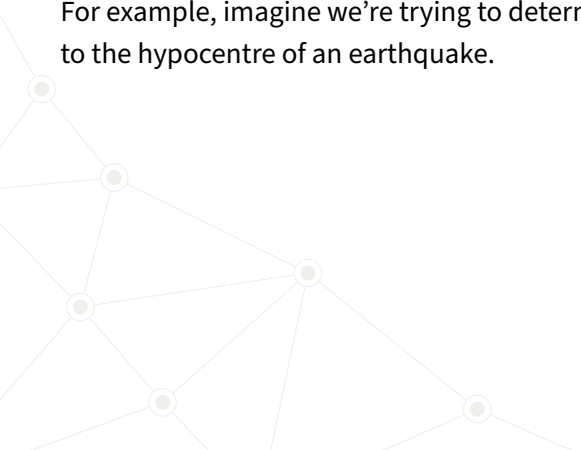


THE MEDIAN – EXAMPLE

The median is useful when dealing with data manifesting 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 dealing with data manifesting 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 dealing with data manifesting 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 dealing with data manifesting 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.

3

THE DEVIATION



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.

4

CORRELATION



CORRELATION

CORRELATION

Correlation or **dependence** is a relationship between two outputs of a dataset. A 'correlation' in many cases means some type of association.

CORRELATION

CORRELATION

Correlation or **dependence** is a relationship between two outputs of a dataset. A 'correlation' in many cases means some type of association.

Correlation is an number between -1 and 1 .

CORRELATION

CORRELATION

Correlation or **dependence** is a relationship between two outputs of a dataset. A 'correlation' in many cases means some type of association.

Correlation is an number between -1 and 1 . Intuitively, a

- negative correlation means that the two series of outputs **contradict** each other;

CORRELATION

CORRELATION

Correlation or **dependence** is a relationship between two outputs of a dataset. A 'correlation' in many cases means some type of association.

Correlation is an number between -1 and 1 . Intuitively, a

- negative correlation means that the two series of outputs **contradict** each other;
- zero correlation means that the two series of outputs are **unrelated**;

CORRELATION

CORRELATION

Correlation or **dependence** is a relationship between two outputs of a dataset. A 'correlation' in many cases means some type of association.

Correlation is an number between -1 and 1 . Intuitively, a

- negative correlation means that the two series of outputs **contradict** each other;
- zero correlation means that the two series of outputs are **unrelated**;
- positive correlation means that the two series of outputs **influence** each other.

COMPUTING CORRELATION

CORRELATION FORMULA

If x_1, \dots, x_n and y_1, \dots, y_n are two series of outputs for the same inputs with means \bar{x} and \bar{y} , their correlation is

$$\text{cor}(x, y) := \frac{(x_1 - \bar{x})(x_2 - \bar{x}) \cdots (x_n - \bar{x})(y_1 - \bar{y})(y_2 - \bar{y}) \cdots (y_n - \bar{y})}{\sqrt{(x_1 - \bar{x})^2(x_2 - \bar{x})^2 \cdots (x_n - \bar{x})^2(y_1 - \bar{y})^2(y_2 - \bar{y})^2 \cdots (y_n - \bar{y})^2}}.$$

INTERPRETING CORRELATION – TABLE

A crude interpretation of correlation is given in the following table:

Coefficient	Strength	Type
-0.7 to -1	Very strong	Negative
-0.5 to -0.7	Strong	Negative
-0.3 to -0.5	Moderate	Negative
0 to -0.3	Weak	Negative
0 to 0.3	Weak	Positive
0.3 to 0.5	Moderate	Positive
0.5 to 0.7	Strong	Positive
0.7 to 1	Very strong	Positive

INTERPRETING CORRELATION – CHART

You can use a scatter plot to visualize two sets of outputs with the same inputs.



INTERPRETING CORRELATION – CHART

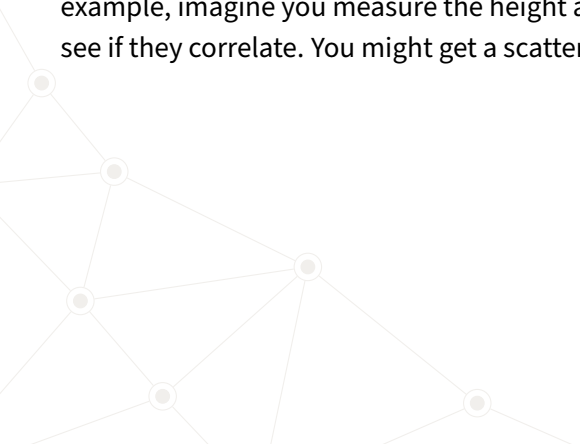
You can use a scatter plot to visualize two sets of outputs with the same inputs.
Correlation tells you **how well you can approximate** this scatter plot by a straight line.



INTERPRETING CORRELATION – CHART

You can use a scatter plot to visualize two sets of outputs with the same inputs.

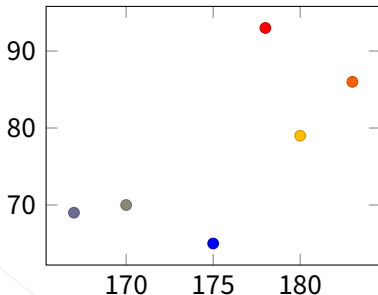
Correlation tells you **how well you can approximate** this scatter plot by a straight line. For example, imagine you measure the height and weight of a sample of people and want to see if they correlate. You might get a scatter plot like this:



INTERPRETING CORRELATION – CHART

You can use a scatter plot to visualize two sets of outputs with the same inputs.

Correlation tells you **how well you can approximate** this scatter plot by a straight line. For example, imagine you measure the height and weight of a sample of people and want to see if they correlate. You might get a scatter plot like this:



FREQUENCY DISTRIBUTION

The background features two large, overlapping geometric shapes. On the left, a light teal triangle points downwards towards the center. On the right, a dark teal triangle points upwards towards the center. These shapes create a V-shape in the middle, framing the title text.

WHAT IS FREQUENCY DISTRIBUTION?

FREQUENCY DISTRIBUTION

A **frequency** of a value is the number of times it occurs in a dataset. A **frequency distribution** is the number of times each variable occurs in a dataset.

TYPES OF FREQUENCY DISTRIBUTIONS

- **Ungrouped frequency distribution:** the number of observations of each output. It's usable for *categorical data*.



TYPES OF FREQUENCY DISTRIBUTIONS

- **Ungrouped frequency distribution:** the number of observations of each output. It's usable for *categorical data*.
- **Grouped frequency distribution:** the number of observations of each **class interval** of a variable. Useful for *quantitative data*.

TYPES OF FREQUENCY DISTRIBUTIONS

- **Ungrouped frequency distribution:** the number of observations of each output. It's usable for *categorical data*.
- **Grouped frequency distribution:** the number of observations of each **class interval** of a variable. Useful for *quantitative data*.
- **Relative frequency distribution:** the proportion of each value or class interval of a variable. Useful for any type of data **if we care about comparing frequencies** rather than amounts.

TYPES OF FREQUENCY DISTRIBUTIONS

- **Ungrouped frequency distribution:** the number of observations of each output. It's usable for *categorical data*.
- **Grouped frequency distribution:** the number of observations of each **class interval** of a variable. Useful for *quantitative data*.
- **Relative frequency distribution:** the proportion of each value or class interval of a variable. Useful for any type of data **if we care about comparing frequencies** rather than amounts.
- **Cumulative frequency distribution:** the sum of frequencies less than or equal to each value or class interval of a variable. Useful when we want to understand how often observations fall below certain values.

CATEGORICAL VS QUANTITATIVE DATA

Quantitative Data represent real amounts that can be added, subtracted etc.



CATEGORICAL VS QUANTITATIVE DATA

Quantitative Data represent real amounts that can be added, subtracted etc.

- Can be discrete or continuous:

CATEGORICAL VS QUANTITATIVE DATA

Quantitative Data represent real amounts that can be added, subtracted etc.

- Can be discrete or continuous:
 - **Discrete data** represents counts of individual items like number of students in a class.

CATEGORICAL VS QUANTITATIVE DATA

Quantitative Data represent real amounts that can be added, subtracted etc.

- Can be discrete or continuous:
 - **Discrete data** represents counts of individual items like number of students in a class.
 - **Continuous data** represents measurements of uncountable values like density, volume or time.

CATEGORICAL VS QUANTITATIVE DATA

Categorical Data represents groupings. They can be recorded as numbers but the numbers represent categories and not actual amounts.

CATEGORICAL VS QUANTITATIVE DATA

Categorical Data represents groupings. They can be recorded as numbers but the numbers represent categories and not actual amounts.

- Can be binary, nominal or ordinal:

CATEGORICAL VS QUANTITATIVE DATA

Categorical Data represents groupings. They can be recorded as numbers but the numbers represent categories and not actual amounts.

- Can be binary, nominal or ordinal:
 - **Binary data** represents yes or no outcomes like coin flips or win/loss situations.

CATEGORICAL VS QUANTITATIVE DATA

Categorical Data represents groupings. They can be recorded as numbers but the numbers represent categories and not actual amounts.

- Can be binary, nominal or ordinal:
 - **Binary data** represents yes or no outcomes like coin flips or win/loss situations.
 - **Nominal data** represents groups without rank or order between them – like the names of species or colours.

CATEGORICAL VS QUANTITATIVE DATA

Categorical Data represents groupings. They can be recorded as numbers but the numbers represent categories and not actual amounts.

- Can be binary, nominal or ordinal:
 - **Binary data** represents yes or no outcomes like coin flips or win/loss situations.
 - **Nominal data** represents groups without rank or order between them – like the names of species or colours.
 - **Ordinal data** represents groups that are ranked – like finishing place in a race.

UNGROUPED FREQUENCY DATA – TABLE

1. Create a table with one column for inputs and as many columns as there are output with a row for each input.

UNGROUPED FREQUENCY DATA – TABLE

1. Create a table with one column for inputs and as many columns as there are output with a row for each input.
 - For **ordinal** variables, the values should be ordered from smallest to largest.

UNGROUPED FREQUENCY DATA – TABLE

1. Create a table with one column for inputs and as many columns as there are output with a row for each input.
 - For **ordinal** variables, the values should be ordered from smallest to largest.
 - For **nominal** variables, the rows can be ordered arbitrarily.
2. Count the **frequencies**.

UNGROUPED FREQUENCY DATA – EXAMPLE 1

A gardener sets up a bird feeder in his backyard. He wishes to know which type of bird species visit the feeder the most.



UNGROUPED FREQUENCY DATA – EXAMPLE 1

A gardener sets up a bird feeder in his backyard. He wishes to know which type of bird species visit the feeder the most. His observations are in the following table:

Species	Frequency
Chickadee	3
Dove	1
Finch	4
Grackle	2
Sparrow	4
Starling	2

UNGROUPED FREQUENCY DATA – EXAMPLE 2

We observe how many times a specific type of tram (based on age) stops at a chosen station each day.



UNGROUPED FREQUENCY DATA – EXAMPLE 2

We observe how many times a specific type of tram (based on age) stops at a chosen station each day.

This experiment may yield a table like this:

Type	Frequency
1990	6
1996	11
2005	3
2017	5

GROUPED FREQUENCY DATA – TABLE

1. Divide the variables into **class intervals**. There's no 'best choice' for the width of a class interval. Different choices convey different meanings.



GROUPED FREQUENCY DATA – TABLE

1. Divide the variables into **class intervals**. There's no 'best choice' for the width of a class interval. Different choices convey different meanings.
 - Calculate the **range** = highest value – lowest value.

GROUPED FREQUENCY DATA – TABLE

1. Divide the variables into **class intervals**. There's no 'best choice' for the width of a class interval. Different choices convey different meanings.
 - Calculate the **range** = highest value – lowest value.
 - Decide on the **class interval width**. **ALWAYS THINK** about that the best width should be! But, if you can't decide, a rule of thumb is the width

$$\text{width} = \frac{\text{range}}{\sqrt{\text{number of inputs}}}.$$

It is typically beneficial to round this value to an integer.

GROUPED FREQUENCY DATA – TABLE

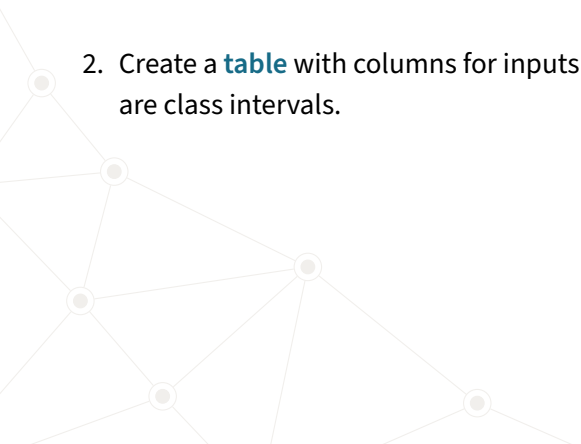
1. Divide the variables into **class intervals**. There's no 'best choice' for the width of a class interval. Different choices convey different meanings.
 - Calculate the **range** = highest value – lowest value.
 - Decide on the **class interval width**. **ALWAYS THINK** about that the best width should be! But, if you can't decide, a rule of thumb is the width

$$\text{width} = \frac{\text{range}}{\sqrt{\text{number of inputs}}}.$$

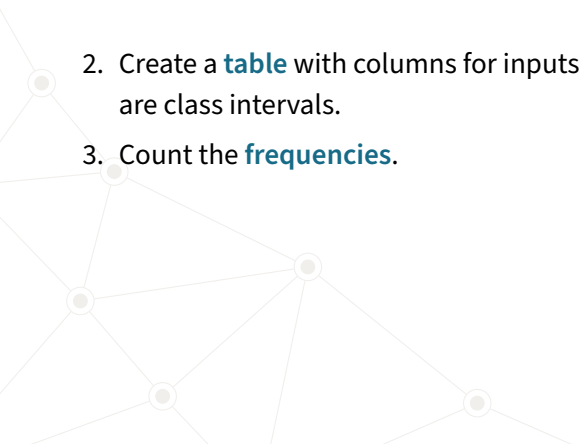
It is typically beneficial to round this value to an integer.

- Calculate the **class intervals**. Each interval is of the form [lower limit, lower limit + width). Simply divide the outputs into these intervals.

GROUPED FREQUENCY DATA – TABLE

- 
- A decorative geometric pattern in the bottom-left corner consisting of several interconnected triangles and lines, with small circles at the vertices.
2. Create a **table** with columns for inputs and each output and as many rows as there are class intervals.

GROUPED FREQUENCY DATA – TABLE

- 
- A decorative geometric pattern in the bottom-left corner consisting of several interconnected triangles and lines, with small circles at the vertices.
2. Create a **table** with columns for inputs and each output and as many rows as there are class intervals.
 3. Count the **frequencies**.

GROUPED FREQUENCY DATA – EXAMPLE

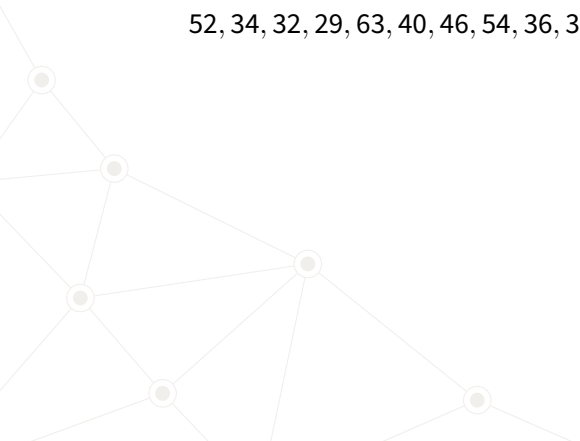
In sociological surveys, you typically want to find the distribution of respondents by age.



GROUPED FREQUENCY DATA – EXAMPLE

In sociological surveys, you typically want to find the distribution of respondents by age. Let's say a survey had 20 respondents. Their ages go like this:

52, 34, 32, 29, 63, 40, 46, 54, 36, 36, 24, 19, 45, 20, 28, 29, 38, 33, 49, 37.



GROUPED FREQUENCY DATA – EXAMPLE

In sociological surveys, you typically want to find the distribution of respondents by age. Let's say a survey had 20 respondents. Their ages go like this:

52, 34, 32, 29, 63, 40, 46, 54, 36, 36, 24, 19, 45, 20, 28, 29, 38, 33, 49, 37.

We calculate the range as highest – lowest = $63 - 19 = 44$.

We calculate the interval width as

$$\text{width} = \frac{\text{range}}{\sqrt{\text{sample size}}} = \frac{44}{\sqrt{20}} = 9.84,$$

and round it up to 10.

GROUPED FREQUENCY DATA – EXAMPLE

In sociological surveys, you typically want to find the distribution of respondents by age. Let's say a survey had 20 respondents. Their ages go like this:

52, 34, 32, 29, 63, 40, 46, 54, 36, 36, 24, 19, 45, 20, 28, 29, 38, 33, 49, 37.

We calculate the range as highest – lowest = $63 - 19 = 44$.

We calculate the interval width as

$$\text{width} = \frac{\text{range}}{\sqrt{\text{sample size}}} = \frac{44}{\sqrt{20}} = 9.84,$$

and round it up to 10.

Therefore, we have the following intervals

$[19, 29)$, $[29, 39)$, $[39, 49)$, $[49, 59)$, $[59, 69)$.

GROUPED FREQUENCY DATA – EXAMPLE

Counting the numbers of outputs falling into each of those intervals gives the table:

Age	Frequency
19 – 28	4
29 – 38	9
39 – 48	3
49 – 58	3
59 – 68	1

RELATIVE FREQUENCY DATA – TABLE

1. Simply create a grouped or ungrouped frequency table.

RELATIVE FREQUENCY DATA – TABLE

1. Simply create a grouped or ungrouped frequency table.
2. To each output add another column to represent **relative frequencies**.

RELATIVE FREQUENCY DATA – EXAMPLE

In our gardener example, the relative frequency table would look like this:

Species	Frequency	Relative Frequency
Chickadee	3	$\frac{3}{3+1+4+2+4+2} = 0.19$
Dove	1	0.06
Finch	4	0.25
Grackle	2	0.13
Sparrow	4	0.25
Starling	2	0.13

CUMULATIVE FREQUENCY DATA – TABLE

1. Create an ungrouped or grouped frequency table for **an ordinal or quantitative variable**. Cumulative frequencies make no sense for nominal variables because they're not ordered.

CUMULATIVE FREQUENCY DATA – TABLE

1. Create an ungrouped or grouped frequency table for **an ordinal or quantitative variable**. Cumulative frequencies make no sense for nominal variables because they're not ordered.
2. Add another column for each output with **cumulative frequency**. The cumulative frequency is the number of observations less than or equal to a certain value or class interval.

CUMULATIVE FREQUENCY DATA – EXAMPLE

Going back to our example of a sociological survey. The cumulative frequency table of the age of survey participants would look like this:

Age	Frequency	Cumulative Frequency
19 – 28	4	4
29 – 38	9	$9 + 4 = 13$
39 – 48	3	$9 + 4 + 3 = 16$
49 – 58	3	19
59 – 68	1	20