

Data Basics

Colby Community College

Definition

A **data matrix** is common way to organize data, especially if collecting data in a spreadsheet.

Definition

A **data matrix** is common way to organize data, especially if collecting data in a spreadsheet.

Example 1

	loan_amount	interest_rate	term	grade	state	total_income	homeownership
1	7500	7.34	36	A	MD	70000	rent
2	25000	9.43	60	B	OH	254000	mortgage
3	14500	6.08	36	A	MO	80000	mortgage
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	3000	7.96	36	A	CA	34000	rent

Definition

A **data matrix** is common way to organize data, especially if collecting data in a spreadsheet.

Example 1

	loan_amount	interest_rate	term	grade	state	total_income	homeownership
1	7500	7.34	36	A	MD	70000	rent
2	25000	9.43	60	B	OH	254000	mortgage
3	14500	6.08	36	A	MO	80000	mortgage
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	3000	7.96	36	A	CA	34000	rent

Definition

Each row is called a **case** or **observational unit**.

Definition

A **data matrix** is common way to organize data, especially if collecting data in a spreadsheet.

Example 1

	loan_amount	interest_rate	term	grade	state	total_income	homeownership
1	7500	7.34	36	A	MD	70000	rent
2	25000	9.43	60	B	OH	254000	mortgage
3	14500	6.08	36	A	MO	80000	mortgage
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	3000	7.96	36	A	CA	34000	rent

Definition

Each row is called a **case** or **observational unit**.

Definition

Each column is called a **variable**.

Note

It is very important to provide descriptions of all the variables in a data matrix. Be sure to include the units of measurement.

Note

It is very important to provide descriptions of all the variables in a data matrix. Be sure to include the units of measurement.

Example 2

	loan_amount	interest_rate	term	grade	state	total_income	homeownership
1	7500	7.34	36	A	MD	70000	rent
2	25000	9.43	60	B	OH	254000	mortgage
3	14500	6.08	36	A	MO	80000	mortgage
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	3000	7.96	36	A	CA	34000	rent

variable	description
loan_amount	Amount of the load received, in US dollars.
interest_rate	Interest rate on the loan, in an annual percentage.
term	The length of the loan, which is always a whole number of months.
grade	Loan grade, which takes values A through G and represents the quality of the loan and its likelihood of being repaid.
state	US state where the borrower resides.
total_income	Borrower's total income, including any second income, in US dollars.
homeownership	Indicates whether the person owns, owns but has a mortgage, or rents.

Definition

Numerical data consisting of numbers representing counts or measurements. Sometimes referred to as **Quantitative** data.

Definition

Numerical data consisting of numbers representing counts or measurements. Sometimes referred to as **Quantitative** data.

Definition

Categorical data consisting of names or labels (not numbers). Sometimes referred to as **Qualitative** data.

Definition

Numerical data consisting of numbers representing counts or measurements. Sometimes referred to as **Quantitative** data.

Definition

Categorical data consisting of names or labels (not numbers). Sometimes referred to as **Qualitative** data.

Note

The names and labels in categorical data are sometimes coded with numbers. When a number is used as a name it is **not** quantitative.

Example 3

The ages (in years) of subjects enrolled in a clinical trial are?

Example 3

The ages (in years) of subjects enrolled in a clinical trial are?

Numerical data.

Example 3

The ages (in years) of subjects enrolled in a clinical trial are?

Numerical data.

Example 4

The genders (male / female / non-binary) of subjects enrolled in a clinical trial are?

Example 3

The ages (in years) of subjects enrolled in a clinical trial are?

Numerical data.

Example 4

The genders (male / female / non-binary) of subjects enrolled in a clinical trial are?

Categorical data.

Example 3

The ages (in years) of subjects enrolled in a clinical trial are?

Numerical data.

Example 4

The genders (male / female / non-binary) of subjects enrolled in a clinical trial are?

Categorical data.

Example 5

Identification numbers 1, 2, 3, ..., 25 are assigned randomly to 25 subjects in a clinical trial. The identification numbers are?

Example 3

The ages (in years) of subjects enrolled in a clinical trial are?

Numerical data.

Example 4

The genders (male / female / non-binary) of subjects enrolled in a clinical trial are?

Categorical data.

Example 5

Identification numbers 1, 2, 3, ..., 25 are assigned randomly to 25 subjects in a clinical trial. The identification numbers are?

Categorical Data.

Example 3

The ages (in years) of subjects enrolled in a clinical trial are?

Numerical data.

Example 4

The genders (male / female / non-binary) of subjects enrolled in a clinical trial are?

Categorical data.

Example 5

Identification numbers 1, 2, 3, ..., 25 are assigned randomly to 25 subjects in a clinical trial. The identification numbers are?

Categorical Data.

Note

The numbers in Example 5 don't actually measure or count anything.

Definition

Discrete data result when the data values are numerical and the number of values is finite, or countable.

Definition

Discrete data result when the data values are numerical and the number of values is finite, or countable.

Example 6

Each of several physicians plans to count the number of physical examinations given during the next full week.

Definition

Discrete data result when the data values are numerical and the number of values is finite, or countable.

Example 6

Each of several physicians plans to count the number of physical examinations given during the next full week.

Example 7

Casino employees plan to roll a fair die until the number 5 is rolled, and they count the number of rolls required to get a 5. (It is possible, but unlikely, that a 5 could never be rolled.)

Definition

Continuous data result from infinitely many possible numerical values, where the collection of values is not countable.

Definition

Continuous data result from infinitely many possible numerical values, where the collection of values is not countable.

Example 8

When a typical patient has blood drawn as part of a routine examination, the volume of blood drawn is between 0 mL and 50 mL.

Definition

Continuous data result from infinitely many possible numerical values, where the collection of values is not countable.

Example 8

When a typical patient has blood drawn as part of a routine examination, the volume of blood drawn is between 0 mL and 50 mL.

Example 9

The unemployment rate for Los Angeles County is 4.69%.

Definition

Continuous data result from infinitely many possible numerical values, where the collection of values is not countable.

Example 8

When a typical patient has blood drawn as part of a routine examination, the volume of blood drawn is between 0 mL and 50 mL.

Example 9

The unemployment rate for Los Angeles County is 4.69%.

Example 10

A grade school teacher measures the heights of his students.

Definition

The **nominal level of measurement** is characterized by data that consist of names, labels, or categories only. The data cannot be arranged in some meaningful order, such as low to high.

Definition

The **nominal level of measurement** is characterized by data that consist of names, labels, or categories only. The data cannot be arranged in some meaningful order, such as low to high.

Example 11

A survey has responses of “yes”, “no”, and “undecided”

Definition

The **nominal level of measurement** is characterized by data that consist of names, labels, or categories only. The data cannot be arranged in some meaningful order, such as low to high.

Example 11

A survey has responses of “yes”, “no”, and “undecided”

Example 12

For an item on a survey, respondents are given a choice of possible answers, and they are coded as follows:

1 is coded as “I agree”

2 is coded as “I disagree”

3 is coded as “I don’t care”

4 is coded as “I refuse to answer”

The numbers 1,2,3, and 4 don’t count or measure anything.

Definition

Data are at the **ordinal level of measurement** if they can be arranged in some order, but differences between data values either cannot be determined or are meaningless.

Definition

Data are at the **ordinal level of measurement** if they can be arranged in some order, but differences between data values either cannot be determined or are meaningless.

Example 13

A college professor assigns grades of A, B, C, D, or F. These grades can be arranged in a meaningful order, but grades are very individualized so the difference between two students grades cannot be calculated.

Definition

Data are at the **ordinal level of measurement** if they can be arranged in some order, but differences between data values either cannot be determined or are meaningless.

Example 13

A college professor assigns grades of A, B, C, D, or F. These grades can be arranged in a meaningful order, but grades are very individualized so the difference between two students grades cannot be calculated.

Example 14

A survey asks people what they felt their blood pressure was. The possible answers are “Low”, “Normal”, “High.” These responses can be arranged in a meaningful order, but the differences between “Low” and “High” doesn’t make sense.

Definition

Data are at the **interval level of measurement** if they can be arranged in order, and differences between data values can be found and are meaningful. *Data at this level do not have a natural zero starting point at which none of the quantity is present.*

Definition

Data are at the **interval level of measurement** if they can be arranged in order, and differences between data values can be found and are meaningful. *Data at this level do not have a natural zero starting point at which none of the quantity is present.*

Example 15

Body temperatures of 98.2°F and 98.9°F are examples of data at this interval level of measurement. The values are ordered, and we can calculate their difference. There is no natural starting point. (The value 0°F is an arbitrary choice and doesn't represent the complete absence of heat.)

Definition

Data are at the **interval level of measurement** if they can be arranged in order, and differences between data values can be found and are meaningful. *Data at this level do not have a natural zero starting point at which none of the quantity is present.*

Example 15

Body temperatures of 98.2°F and 98.9°F are examples of data at this interval level of measurement. The values are ordered, and we can calculate their difference. There is no natural starting point. (The value 0°F is an arbitrary choice and doesn't represent the complete absence of heat.)

Example 16

The years 1492 and 1776 can be arranged in order, and the difference of 284 years is meaningful. But, time doesn't not have a natural starting point that represents "no time."

Definition

Data are at the **ratio level of measurement** if they can be arranged in order, differences can be found and are meaningful, and there is a natural starting point which indicates that none of the quantity is present.

Definition

Data are at the **ratio level of measurement** if they can be arranged in order, differences can be found and are meaningful, and there is a natural starting point which indicates that none of the quantity is present.

Example 17

Heights of 180cm and 90cm for a high school student and a preschool student are at the ratio level of measurement. (The starting point is 0cm and 180cm is twice as tall as 90cm.)

Definition

Data are at the **ratio level of measurement** if they can be arranged in order, differences can be found and are meaningful, and there is a natural starting point which indicates that none of the quantity is present.

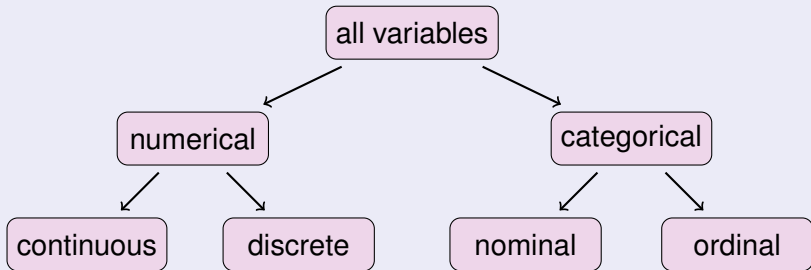
Example 17

Heights of 180cm and 90cm for a high school student and a preschool student are at the ratio level of measurement. (The starting point is 0cm and 180cm is twice as tall as 90cm.)

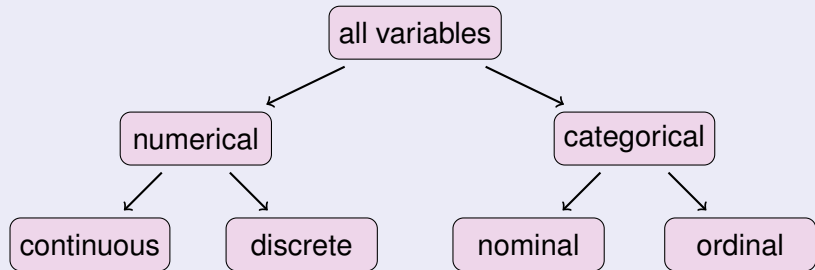
Example 18

The times of 50 minutes and 100 minutes for a math class are at the ratio level of measurement. (The starting point is 0 minutes and 100 minutes is twice as long as 50 minutes.)

Classification of Variables



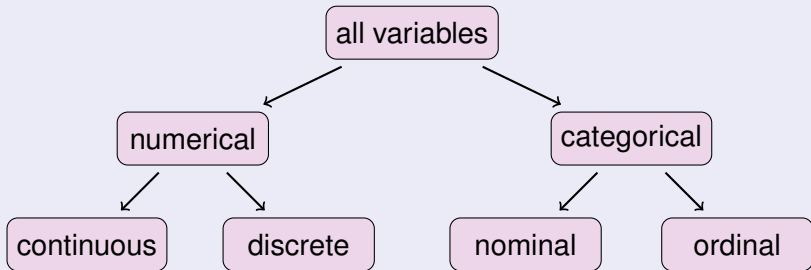
Classification of Variables



Note

For simplicities sake, we will often treat ordinal data as nominal data.

Classification of Variables



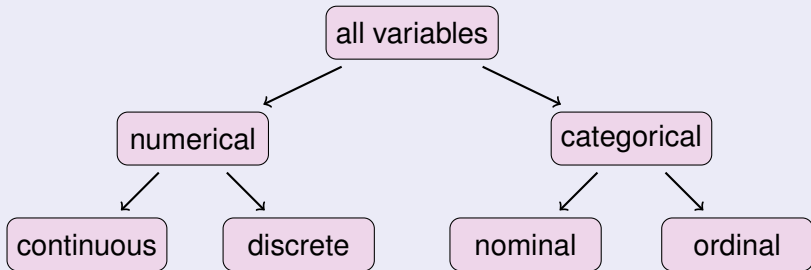
Note

For simplicities sake, we will often treat ordinal data as nominal data.

Note

If your data consists of real numbers, then it is generally continuous.

Classification of Variables



Note

For simplicities sake, we will often treat ordinal data as nominal data.

Note

If your data consists of real numbers, then it is generally continuous.

Note

If your data consists of only integers, then it is generally discrete.

Relationships Between Variables

Often, researchers will want to study the relationship between variables.

Relationships Between Variables

Often, researchers will want to study the relationship between variables.

- If homeownership is lower than the national average in one county, will the percent of multi-unit structures in that county tend to be above or below the national average?

Relationships Between Variables

Often, researchers will want to study the relationship between variables.

- If homeownership is lower than the national average in one county, will the percent of multi-unit structures in that county tend to be above or below the national average?
- Does a higher than average increase in county population tend to correspond to counties with higher or lower median household incomes?

Relationships Between Variables

Often, researchers will want to study the relationship between variables.

- If homeownership is lower than the national average in one county, will the percent of multi-unit structures in that county tend to be above or below the national average?
- Does a higher than average increase in county population tend to correspond to counties with higher or lower median household incomes?
- How useful a predictor is median education level for the median household income for US counties?

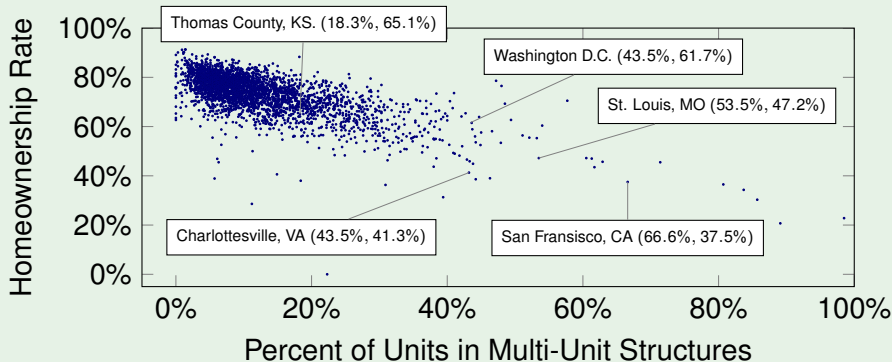
Definition

A **scatterplot** is a plot of paired (x, y) numerical data with a horizontal x -axis and a vertical y -axis. The horizontal axis is used for the first variable (x), and the vertical axis for the second variable (y).

Definition

A **scatterplot** is a plot of paired (x, y) numerical data with a horizontal x-axis and a vertical y-axis. The horizontal axis is used for the first variable (x), and the vertical axis for the second variable (y).

Example 19



Definition

When two variables show some connection with one another, they are called **associated** variables. Sometimes referred to as **dependent**.

Definition

When two variables show some connection with one another, they are called **associated** variables. Sometimes referred to as **dependent**.

Definition

If there is a downward trend in the scatter plot, the variables are said to be **negatively associated**.

Definition

When two variables show some connection with one another, they are called **associated** variables. Sometimes referred to as **dependent**.

Definition

If there is a downward trend in the scatter plot, the variables are said to be **negatively associated**.

Definition

If there is an upward trend in the scatter plot, the variables are said to be **positive associated**.

Definition

When two variables show some connection with one another, they are called **associated** variables. Sometimes referred to as **dependent**.

Definition

If there is a downward trend in the scatter plot, the variables are said to be **negatively associated**.

Definition

If there is an upward trend in the scatter plot, the variables are said to be **positive associated**.

Definition

If variables are not associated, then they are said to be **independent**.

Definition

When two variables show some connection with one another, they are called **associated** variables. Sometimes referred to as **dependent**.

Definition

If there is a downward trend in the scatter plot, the variables are said to be **negatively associated**.

Definition

If there is an upward trend in the scatter plot, the variables are said to be **positive associated**.

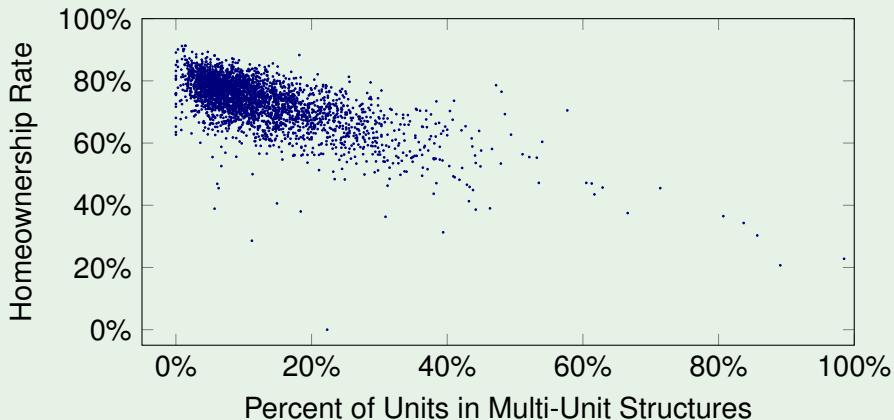
Definition

If variables are not associated, then they are said to be **independent**.

Note

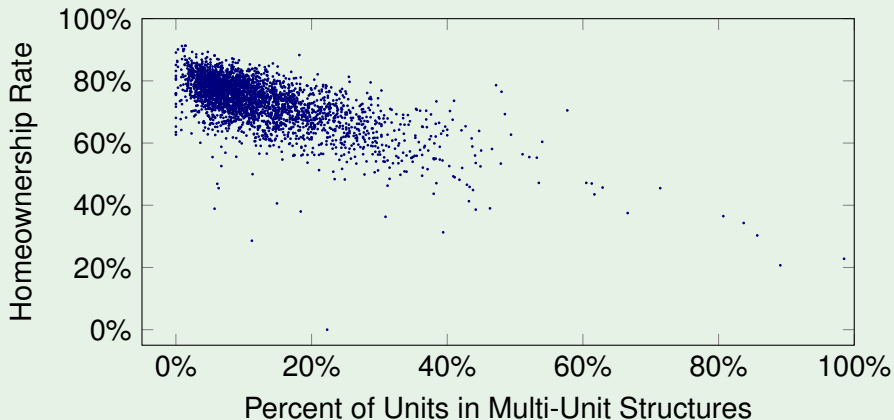
It is impossible to be both associated and independent.

Example 20



Are the Multi-Unit Structure Rate and the Homeownership Rate associated?

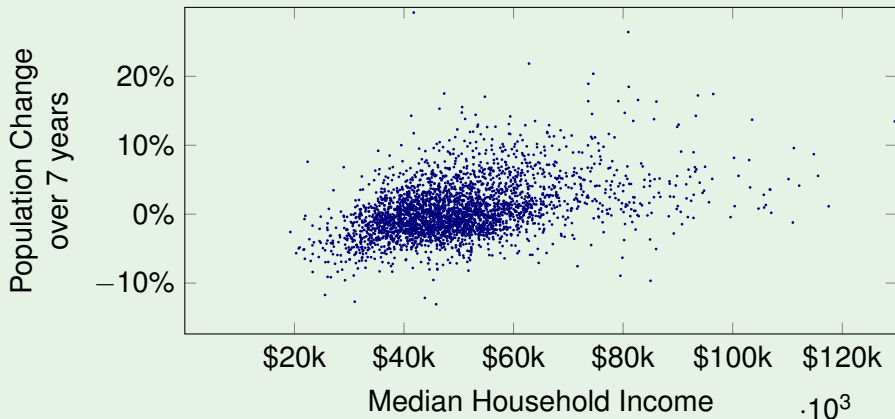
Example 20



Are the Multi-Unit Structure Rate and the Homeownership Rate associated?

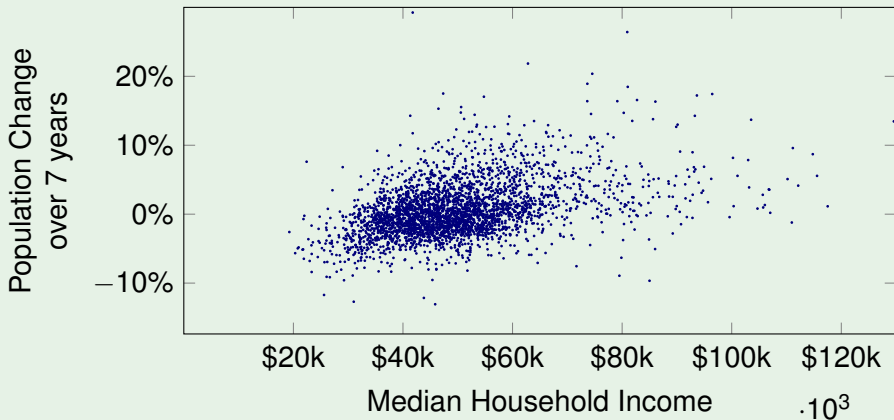
Yes, they are negatively associated.

Example 21



Are the Median Household Income and the Population Change associated?

Example 21



Are the Median Household Income and the Population Change associated?

Yes, they are positively associated.

Example 22

Consider the question:

If there is an increase in the median household income in a county, does this drive an increase in its population?

Example 22

Consider the question:

If there is an increase in the median household income in a county, does this drive an increase in its population?

A higher median income is likely one of the causes of an increased population.

Example 22

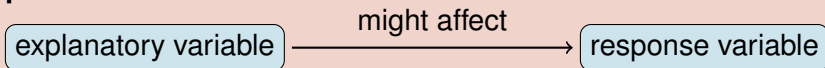
Consider the question:

If there is an increase in the median household income in a county, does this drive an increase in its population?

A higher median income is likely one of the causes of an increased population.

Definition

When we suspect one variable might causally affect another, we label the first variable the **explanatory variable** and the second the **response variable**.



Example 22

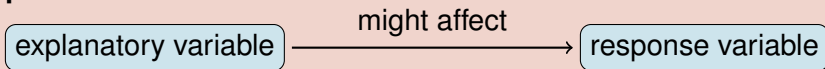
Consider the question:

If there is an increase in the median household income in a county, does this drive an increase in its population?

A higher median income is likely one of the causes of an increased population.

Definition

When we suspect one variable might causally affect another, we label the first variable the **explanatory variable** and the second the **response variable**.



Note

Labeling a variable like this does **nothing** to guarantee that a casual relationship exists!

Definition

A **Placebo** is a treatment that has no medicinal effect. (Such as a sugar pill or saline injection.)

Definition

A **Placebo** is a treatment that has no medicinal effect. (Such as a sugar pill or saline injection.)

Example 23

In 1954, an experiment was designed to test the effectiveness of the Salk vaccine in preventing polio, which had killed or paralyzed thousands of children.

Definition

A **Placebo** is a treatment that has no medicinal effect. (Such as a sugar pill or saline injection.)

Example 23

In 1954, an experiment was designed to test the effectiveness of the Salk vaccine in preventing polio, which had killed or paralyzed thousands of children.

By random selection, 401,974 children were assigned to two groups:

- 200,745 children were given injections of the Salk vaccine.
- 201,229 children were given placebo injections that contained no drug.

Definition

A **Placebo** is a treatment that has no medicinal effect. (Such as a sugar pill or saline injection.)

Example 23

In 1954, an experiment was designed to test the effectiveness of the Salk vaccine in preventing polio, which had killed or paralyzed thousands of children.

By random selection, 401,974 children were assigned to two groups:

- 200,745 children were given injections of the Salk vaccine.
- 201,229 children were given placebo injections that contained no drug.

Among the children given the Salk vaccine, 33 later developed paralytic polio, and among the children given a placebo, 115 later developed paralytic polio.

Definition

In an **experiment**, we apply some treatment and then proceed to observe its effects on the individuals. (The individuals in experiments are called **subjects**.)

Definition

In an **experiment**, we apply some treatment and then proceed to observe its effects on the individuals. (The individuals in experiments are called **subjects**.)

Note

In general, association does not imply causation, and causation can only be inferred from a randomized experiment.

Definition

In an **experiment**, we apply some treatment and then proceed to observe its effects on the individuals. (The individuals in experiments are called **subjects**.)

Note

In general, association does not imply causation, and causation can only be inferred from a randomized experiment.

Definition

In an **observational study**, we observe and measure specific characteristics, but we don't attempt to *modify* the individuals being studied.

Definition

In an **experiment**, we apply some treatment and then proceed to observe its effects on the individuals. (The individuals in experiments are called **subjects**.)

Note

In general, association does not imply causation, and causation can only be inferred from a randomized experiment.

Definition

In an **observational study**, we observe and measure specific characteristics, but we don't attempt to *modify* the individuals being studied.

Note

In general, experiments are preferable to observational studies. But there are times where there are cost, time, or ethical concerns that prohibit the use of an experiment.

Example 24

Suppose we want to determine if listening to music while driving increases the chance of being in an collision.

Example 24

Suppose we want to determine if listening to music while driving increases the chance of being in an collision.

- **Observational study:** If we gathered police reports about collisions and used them to determine if the person was listening to music or not.

Example 24

Suppose we want to determine if listening to music while driving increases the chance of being in an collision.

- **Observational study:** If we gathered police reports about collisions and used them to determine if the person was listening to music or not.
- **Experiment:** We randomly assign subjects to either listen to music while driving or listen to nothing. We then count how many collisions each subject is involved in.

Example 24

Suppose we want to determine if listening to music while driving increases the chance of being in an collision.

- **Observational study:** If we gathered police reports about collisions and used them to determine if the person was listening to music or not.
- **Experiment:** We randomly assign subjects to either listen to music while driving or listen to nothing. We then count how many collisions each subject is involved in.

Is music likely to be the cause of a collision?

Example 24

Suppose we want to determine if listening to music while driving increases the chance of being in an collision.

- **Observational study:** If we gathered police reports about collisions and used them to determine if the person was listening to music or not.
- **Experiment:** We randomly assign subjects to either listen to music while driving or listen to nothing. We then count how many collisions each subject is involved in.

Is music likely to be the cause of a collision?

Automobile collisions happen due to a large number of reasons, many of which have nothing to do with music.