

# 1.1 What is data?

## Data

**Data** is information, especially facts or numbers, usually collected or computed for purposes of analysis. Ex: The world population was about **300** million in the year **1000**, **500** million in **1500**, and **7** billion in **2000**. Ex: Analysis of the population data suggests the world's population is growing rapidly, which may influence various decisions like use of natural resources.

### Datum vs. data

*Historically, datum is defined as a single item, and data as the plural of datum. Language evolves, however, and the use of the term datum is diminishing. This material follows the increasingly common usage of the term data for both the singular and plural.*

The amount of collected data has grown tremendously. A first reason is because computers became ubiquitous around the **1980**'s and can easily record data. More recently, the world-wide web became ubiquitous in the early **2000**'s, transforming how people do business, communicate, and recreate, in ways such that data is easily recorded and analyzed. Smartphones and tablets of the **2010**'s provide nearly continuous computer/web access. Plus, numerous items like streets, cars, and buildings have recently been equipped with sensors and cameras and allowing for more data collection. Some estimates are that **90%** of all data ever collected was generated in just the past couple years.

The figure below shows the worldwide data collected per year, in zettabytes. A **zettabyte** is one sextillion or  $10^{21}$  bytes. The table below lists common sources of data.

Figure 1.1.1: Worldwide data collected per year.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

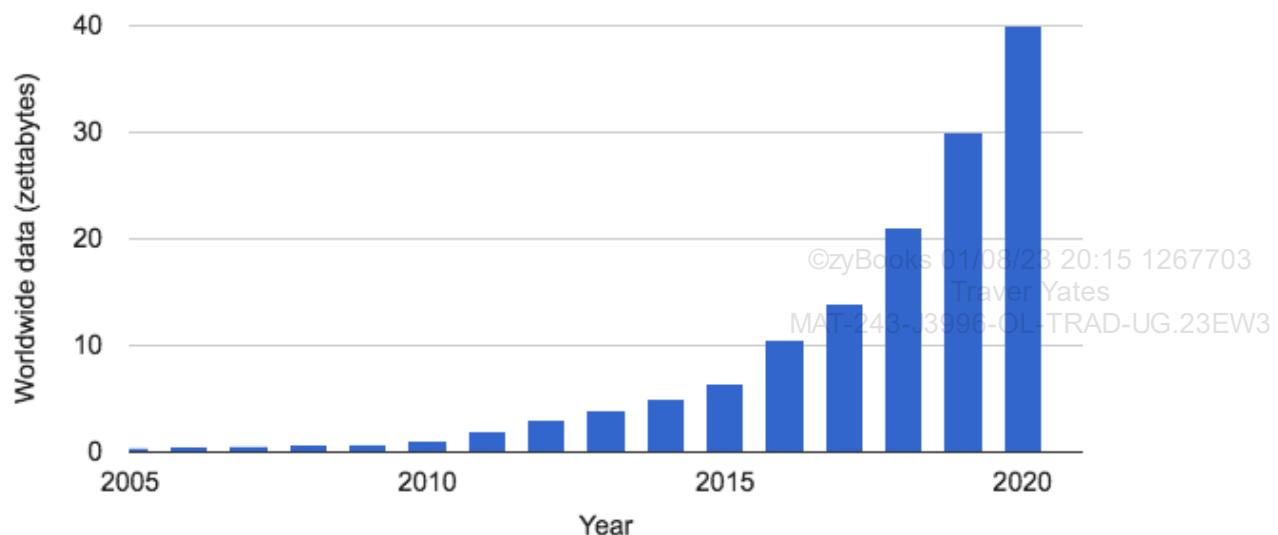


Table 1.1.1: Common sources of data.

Social networks	Traditional business systems	Internet of things
<p>Human-generated data</p> <ul style="list-style-type: none"> <li>• Social Networks: Facebook, Twitter, etc.</li> <li>• Blogs and comments</li> <li>• Personal documents</li> <li>• Pictures: Instagram, Snapchat, etc.</li> <li>• Videos: YouTube etc.</li> <li>• Internet searches</li> <li>• Mobile data: text messages</li> <li>• User-generated maps</li> <li>• E-mail</li> </ul>	<p>Data produced by Public Agencies</p> <ul style="list-style-type: none"> <li>• Medical records</li> </ul> <p>Data produced by businesses</p> <ul style="list-style-type: none"> <li>• Commercial transactions</li> <li>• Banking/stock records</li> <li>• E-commerce</li> <li>• Credit cards</li> </ul>	<p>Data from sensors</p> <ul style="list-style-type: none"> <li>• Fixed sensors <ul style="list-style-type: none"> <li>◦ Home automation</li> <li>◦ Weather/pollution sensors</li> <li>◦ Traffic sensors/webcam</li> <li>◦ Scientific sensors</li> <li>◦ Security/surveillance videos/images</li> </ul> </li> <li>• Mobile sensors (tracking) <ul style="list-style-type: none"> <li>◦ Mobile phone location</li> <li>◦ Cars</li> <li>◦ Satellite images</li> </ul> </li> </ul> <p>Data from computer systems</p> <ul style="list-style-type: none"> <li>• Logs</li> <li>• Web logs</li> </ul>

Source: [United Nations Statistics Division, 2015](#)<sup>1</sup>

**PARTICIPATION ACTIVITY****1.1.1: Data.**

- 1) The amount of data collected worldwide in **2016** is about \_\_\_\_ zettabytes.

©zyBooks 01/08/23 20:15 1267700  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- 1
- 10
- 100

- 2) A zettabyte is \_\_\_\_ bytes.



- one trillion
- 1,000 trillion
- one sextillion

- 3) Which was a substantial source of computerized data before the year **1990**?



- Social network data (like Facebook)
- Internet search data
- Medical records
- None of the above

## Data analytics

The abundance of collected data provides new opportunities for analysis. **Data analytics** is the field of analyzing data to gain insight, draw conclusions, or make decisions.

With so much data being collected today, one can imagine that data analytics is a growing field with increasing job opportunities. Big data is a term commonly used to refer to data analytics on large amounts of data, which is the form in which much data exists today. **Big data** refers to very large data sets that cannot be processed by traditional methods, and is characterized by high volume, rapid velocity of collection, and variety in type and quality. Articles summarizing jobs in big data are abundant, and summaries and predictions both describe large increases in job opportunities, such as this [2014 Forbes article on big data jobs](#).

©zyBooks 01/08/23 20:15 1267700  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Below are real-world applications of data analysis.

## Example 1.1.1: Data analysis catches cheating teachers.

Standardized exams are commonly given to students in public schools. The average scores for a teacher's students are commonly used to evaluate a teacher or a school. A researcher performed data analysis to detect whether some teachers were cheating. For example, if a particular teacher's students answered the last 10 or so questions correctly/more frequently than for another teacher, one might assume that the teacher filled in those last questions (correctly) for students who didn't complete the exams. Or, if a teacher's students did well above average one year, but those same students performed below average the year before and the year after, one might assume that the teacher gave students the answers.

In the book *Freakonomics*, Steven Levitt described analyses he performed on several years of exam data from Chicago public schools. He found that at least 5% of teachers were cheating. As a result of his data analysis, several teachers were fired, and cheating subsequently decreased.

[Levitt's paper describing the analysis](#)

<sup>2</sup> [Two-minute video of Levitt discussing the analysis](#)

## Example 1.1.2: Sports and data analytics.

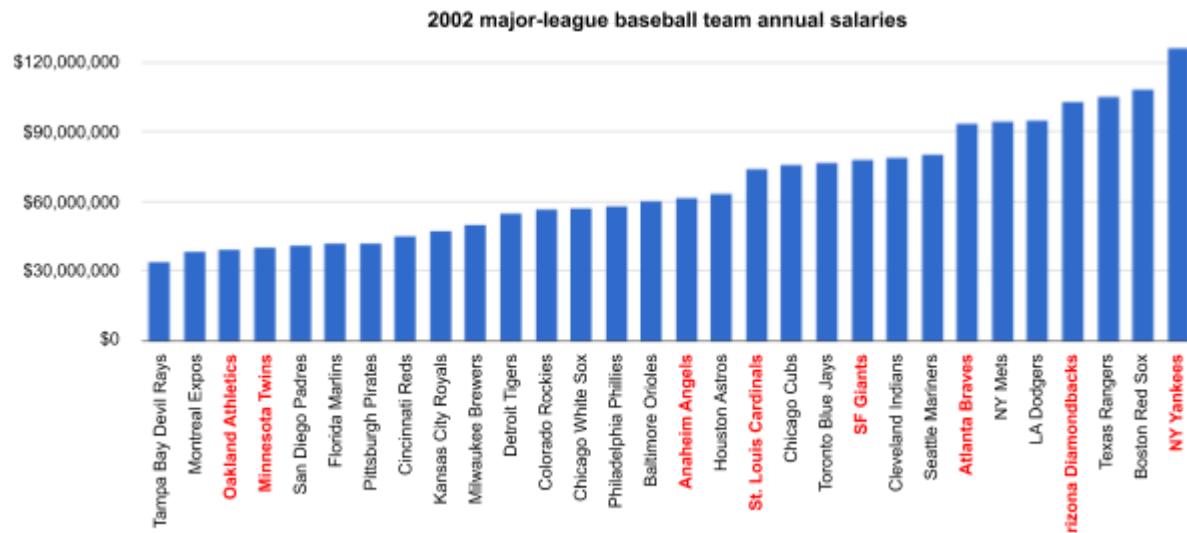
In the early 2000's, the Oakland Athletics had one of the smallest budgets in professional baseball. The team leaders used data analytics to gain an edge. Traditionally, baseball players were sought based on widely-known factors like a player's batting average (how often the player got a hit), runs batted in (how many runs the player caused by making a hit), and similar numbers. Instead, through data analysis, the team leaders found less-popular factors were more important, like on-base percentage. The team thus hired players strong in those less-advertised factors and paid such players with lower salaries due to not being in high demand. The technique worked, and the Oakland Athletics made the playoffs, both in 2002 and 2003, despite having nearly the lowest salaries in the league.



Source: Oakland A's stadium ([Travis Wise / CC-BY-SA-2.0](#) via Flickr)

<sup>3</sup>

This real story is the basis of the popular [2014 movie Moneyball](#) starring Brad Pitt. Many teams, in baseball and other sports, have since adopted such data-analytic techniques. With the advent of computers, and thus more recording of data and more ability to analyze such data, data-analytic techniques are used in more arenas to gain insight and achieve better results. Ex: Online dating sites, stock market investing, language translation, and much more.



©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**2002** playoff teams (in red above): NY Yankees, Anaheim Angels, *Oakland Athletics*, Minnesota Twins, Atlanta Braves, San Francisco Giants, Arizona Diamondbacks, St. Louis Cardinals.

## Types of data analytics

Three types of data analytics exist:

- **Descriptive** data analytics seeks to describe data, providing insight and knowledge. Ex: Based on collected data, the world population in **2015** is about **7** billion.
- **Predictive** data analytics seeks to make predictions from data. Ex: Using models based on birth rates, death rates, medical care improvements, and other data, the United Nations predicts the world population will reach **11.2** billion in **2100**.
- **Prescriptive** data analytics seeks to make decisions (prescriptions) based on data. Ex: Population predictions for specific countries help the United Nations decide where to focus agricultural development efforts.

**PARTICIPATION ACTIVITY**

1.1.2: Worldwide population.

**Animation content:**

undefined

**Animation captions:**

1. Descriptive analytics: Describes the data, perhaps to provide insight. Ex: From census and other data, U.N. estimates 2015 world population as **7** billion.
2. Predictive analytics: Based on models, predict future information. Ex: Based on models, the U.N. predicts world population in 2100 could be as high as **16** billion.
3. Prescriptive: Make decisions based on descriptive/predictive analyses. Ex: Growth is great in India, so U.N. may focus more agricultural efforts there.

**PARTICIPATION ACTIVITY**

1.1.3: Descriptive, predictive, and prescriptive analytics.



Select the definition that matches each term

1) Descriptive analytics

- Strives to make decisions or recommendations based on data.
- Given data and a model, strives to determine future values.
- Given existing data, strives to summarize the data, perhaps to gain insight.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

2) Predictive analytics

- Strives to make decisions or recommendations based on data.
- Given data and a model, strives to determine future values.
- Given existing data, strives to summarize the data, perhaps to gain insight.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### 3) Prescriptive analytics

- Strives to make decisions or recommendations based on data.
- Given data and a model, strives to determine future values.
- Given existing data, strives to summarize the data, perhaps to gain insight.

Reset

## Types of data

### Variables

Data is typically represented using variables. A **variable** is an item that can have different ("varying") values. Ex: A person's age is a variable and can have the value **10, 33, 99**, or other values. Variables are often considered as being of two possible types:

- A **quantitative variable** can take on a numeric value (quantitative data) that can be measured and ordered. Ex: A person's age, the outside temperature, and a meal's price are quantitative variables. Example numeric values are an age of **33** or **99** years, a temperature of **40** or **45** degrees, and a price of **12** or **15** dollars.
- A **categorical variable** can take on the value (usually a label) of one of several categories. Ex: A person's gender, seasons, and U.S. companies are categorical variables. Gender can be male or female, seasons can be fall, winter, spring, or summer, and U.S. companies can be Wal-Mart, McDonalds, UPS, etc. A categorical variable is often called a **qualitative variable** (known by qualities, rather than quantities).

Most numbers represent quantitative data, but exceptions exist. Ex: A person's phone number is a number but is not quantitative data; a phone number isn't measured, nor ordered; people don't say: "Joe's phone number is greater than Mary's." In general, if adding the numbers makes sense, the variable is likely quantitative, else categorical. (People may add ages but don't add phone numbers.)

A reason for distinguishing variable types is that each type is handled differently in data analytics. Ex: A categorical variable typically involves counting the instances of each category, often then depicted with a bar chart or pie chart. But a quantitative variable is commonly plotted versus another quantitative variable, often depicted with a scatter plot or line chart. Those chart types are described in other sections.

**PARTICIPATION ACTIVITY****1.1.4: Quantitative vs. categorical variables.**

©zyBooks 01/08/23 20:15 1267703

Traver Yates



MAT-243-J3996-OL-TRAD-UG.23EW3

1) A car's age.

- Quantitative
- Categorical



2) A car's maker.

- Quantitative
- Categorical



3) A house's square footage.

- Quantitative
- Categorical



4) A house's color.

- Quantitative
- Categorical



5) A house's address.

- Quantitative
- Categorical



6) "Qualitative variable" is likely another term for which type?

- Quantitative
- Categorical

©zyBooks 01/08/23 20:15 1267703

Traver Yates



MAT-243-J3996-OL-TRAD-UG.23EW3

7) "Numerical variable" is likely another term for which type?

- Quantitative
- Categorical



## Types of categorical variables

Two types of categorical variables are often distinguished:

- A **nominal variable**'s categories have no ordering, existing in name only, like apples, oranges, and grapes. ("Nominal" means "in name only").
- An **ordinal variable**'s categories have an ordering, like disagree, neutral, and agree.

©zyBooks 01/08/23 20:15 1267703

The difference is sometimes relevant. Ex: On a chart, the ordinal variables would almost always be sorted along the x-axis, listed as "small medium large" rather than arbitrarily as "small large medium."

### PARTICIPATION ACTIVITY

#### 1.1.5: Categorical variables: Nominal versus ordinal.



1) A car comes in 5 possible colors: red, grey, brown, black, and white.



- Nominal
- Ordinal

2) A movie has 5 possible ratings: G, PG, PG-13, R, and NC-17. (See [movie ratings](#) if unfamiliar.)



- Nominal
- Ordinal

3) An Amazon product has 5 possible ratings: 1 star, 2 stars, 3 stars, 4 stars, or 5 stars.



- Nominal
- Ordinal

4) A survey asks users to enter a number indicating political affiliation: 1 for Libertarian, 2 for Democratic, 3 for Republican, and 4 for Other.



- Nominal
- Ordinal

5) A form asks a person to indicate a country of birth.



- Nominal
- Ordinal

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Types of quantitative variables

Two types of quantitative variables are often distinguished:

- A **continuous variable**'s values are infinite along a continuum of values within a range, typically real numbers. Continuous variables usually represent measurements, like height (0.00104 meters) or temperature (98.6 degrees).
- A **discrete variable**'s values are finite within a range, typically integers. Discrete variables usually represent countable items, like people in a family (5) or cars in a city (502, 434). Generally, if "number of" can be added to the beginning, the variable is discrete, like "number of people in a family", but not "number of height". Note: "Discrete" means separate or distinct, not to be confused with "discreet" which means careful or unobtrusive.

### PARTICIPATION ACTIVITY

#### 1.1.6: Continuous vs. discrete quantitative variables.



Indicate whether the variable is continuous or discrete.

1) Width of a house.



- Continuous
- Discrete

2) Height of a human.



- Continuous
- Discrete

3) Gallons in a car's gas tank.



- Continuous
- Discrete

4) Fingers on a human's hands.



- Continuous
- Discrete

5) Hairs on a human's head.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



- Continuous
- Discrete

6) Air molecules in a house.



- Continuous

## References

- (\*1) United Nations Global Working Group on Big Data for Official Statistics Task Team on Cross-Cutting Issues. "Deliverable 2: Revision and Further Development of the Classification of Big Data." *United Nations Statistics Division*. 12 October 2015, [unstats.un.org/unsd/trade/events/2015/abudhabi/gwg/GWG%202015%20-%20item%202%20\(iv\)%20-%20Big%20Data%20Classification.pdf](http://unstats.un.org/unsd/trade/events/2015/abudhabi/gwg/GWG%202015%20-%20item%202%20(iv)%20-%20Big%20Data%20Classification.pdf). ©zyBooks 01/08/23 20:15 1267703 Traver Yates MAT-243-J3996-OL-TRAD-UG.23EW3
- (\*2) Jacob, Brian A. and Steven D. Levitt. "ROTTEN APPLES: AN INVESTIGATION OF THE PREVALENCE AND PREDICTORS OF TEACHER CHEATING." *Quarterly Journal of Economics*. Volume 118, Issue 3, 1 August 2003, Pages 843-877, doi.org/10.1162/00335530360698441.
- (\*3) Wise, Travis. "Oakland A's." *Flickr*. 12 May 2007, [www.flickr.com/photos/photographinatravis/16666072878](http://www.flickr.com/photos/photographinatravis/16666072878).

## 1.2 What is data visualization?

### Introduction to data visualization

**Data visualization** is the display of data in a format, such as a table or chart, that seeks to achieve a goal of conveying particular information to a viewer. Data presented in a text-only format often does not convey information well. Ex: Given this text-only data on 2013 median house prices in southern California counties, finding the price for a particular county is inconvenient: Los Angeles \$405,000; Orange \$661,000; Riverside \$306,000; San Bernardino \$192,000; San Diego \$473,000; Ventura \$464,000.

Instead, displaying the data visually as a table better conveys the information. A **table** displays data using rows and columns.

Table 1.2.1: Southern California median house prices by county (2013).

County	Median house price
Los Angeles	\$405,000
Orange	\$661,000
Riverside	\$306,000

San Bernardino	\$192,000
San Diego	\$473,000
Ventura	\$464,000

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

As another example, the following data represents California median house prices from 2000-2010: 2000 \$241,000; 2001 \$262,000; 2002 \$316,000; 2003 \$372,000; 2004 \$451,000; 2005 \$523,000; 2006 \$556,000; 2007 \$560,000; 2008 \$348,000; 2009 \$275,000; 2010 \$305,000. A table conveys the information better than text, but if the goal is to illustrate the housing price "bubble" that grew and then burst in 2008, a chart is even better.

Table 1.2.2: California median house prices, 2000-2010.

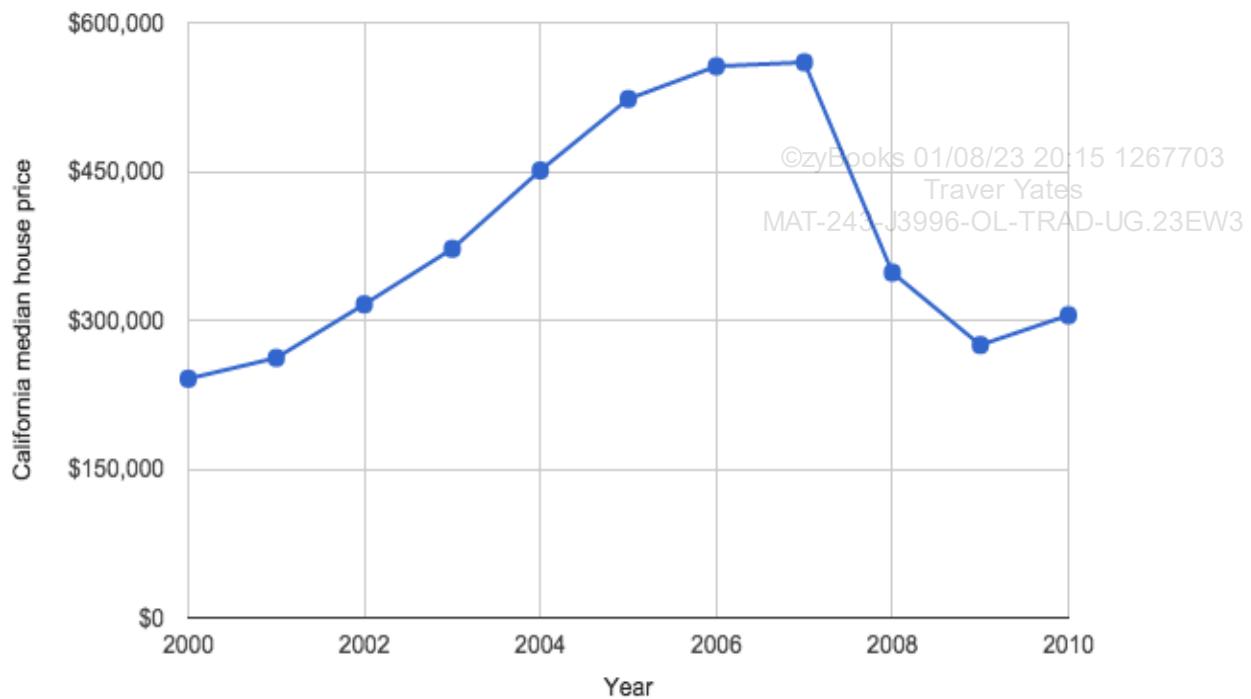
Year	California median house price
2000	\$241,000
2001	\$262,000
2002	\$316,000
2003	\$372,000
2004	\$451,000
2005	\$523,000
2006	\$556,000
2007	\$560,000
2008	\$348,000
2009	\$275,000
2010	\$305,000

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Figure 1.2.1: California median house prices, 2000-2010.

**PARTICIPATION ACTIVITY**

1.2.1: Data visualization.



Refer to the tables and charts above.

- 1) Refer to the *table* above showing California house prices by county (2013). A company is considering moving offices to San Bernardino county. What is the median house price in that county?

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- 2) Refer to the *table* above showing California house prices from 2000-2010. In what year did the price bubble burst? That is, in what year was the price drastically lower than the previous year?



Check Show answer

- 3) Refer to the *figure* above showing California house prices from 2000-2010. In what year was the peak of house prices?

 ©zyBooks 01/08/23 20:15 1267703 Traver Yates MAT-243-J3996-OL-TRAD-UG.23EW3 Check Show answer

- 4) Referring to the *figure* above showing California house prices from 2000-2010, what was the relative difference between the highest and lowest California house prices? Answer with: double, triple, or quadruple.

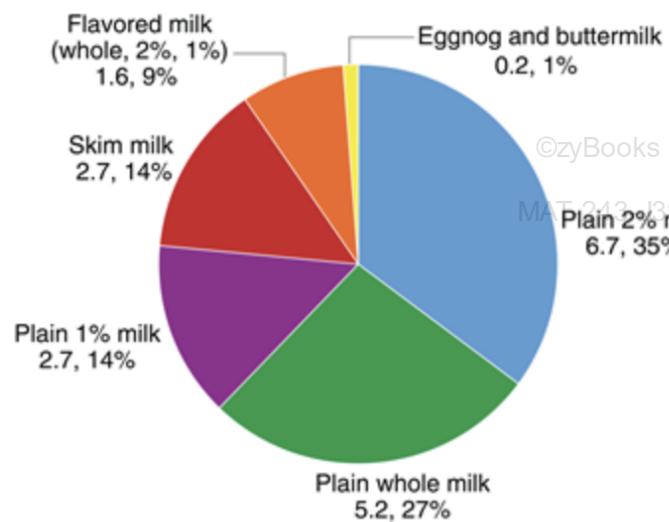
 Check Show answer

## Uses of data visualization

Expressing data as a table or chart allows the viewer to comprehend data more quickly than data presented as a list of numbers. A chart is particularly helpful in analyzing large datasets where a list, or even a table of the data would be incomprehensible. Visual representation is also more intuitively grasped than numbers. The pie chart below shows the per capita availability of milk in the United States in 2013. The viewer is able to quickly grasp that plain 2% milk has the greatest availability, and gains an intuitive sense of how much more 2% milk is available than any other category.

Figure 1.2.2: Charts allow for quick analysis of data.

 ©zyBooks 01/08/23 20:15 1267703 Traver Yates MAT-243-J3996-OL-TRAD-UG.23EW3

**U.S. per capita fluid milk availability, gallons, shares, 2013**

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Source: USDA, Economic Research Service, Food Availability Data.

Credit: Brilliant.org <sup>1 2</sup>

A chart can help the viewer see trends in the data. The chart below shows the price of gold from 1971 to 2019. The overall trend is upward with the exception of market crashes in 1981 and 2013 and short downward trends in that period. Thus, someone looking to invest in gold might conclude that gold is generally a good long-term investment if timed properly.

Figure 1.2.3: Charts help identify trends.



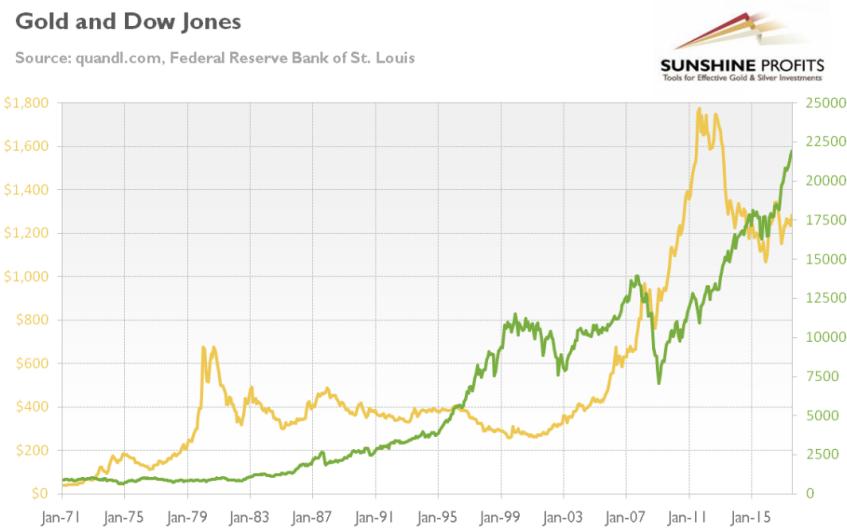
©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Credit: Macrotrends<sup>3</sup>

Charts also allow the viewer to identify relationships and patterns in the data. The chart below shows the gold prices, in yellow, and stock prices, in green, from 1971 to 2017. An investor might note that gold and stock prices do not always move together, and that, as of the end of 2011, gold was relatively expensive compared to stock prices. This type of chart might be used to ascertain whether the time is right to buy gold, or if the price of gold is at a peak and likely to come back down.

©zyBooks 01/08/23 20:15 1267703  
MAT-243-J3996-OL-TRAD-UG.23EW3

Figure 1.2.4: Charts help identify relationships between data.



Credit: Sunshine Profits<sup>4</sup>

### PARTICIPATION ACTIVITY

#### 1.2.2: Using data visualizations.



- 1) Charts are useful for small datasets, but become too crowded when used with large datasets.



- True
- False

- 2) Charts can be used to identify relationships between different variables.



- True
- False

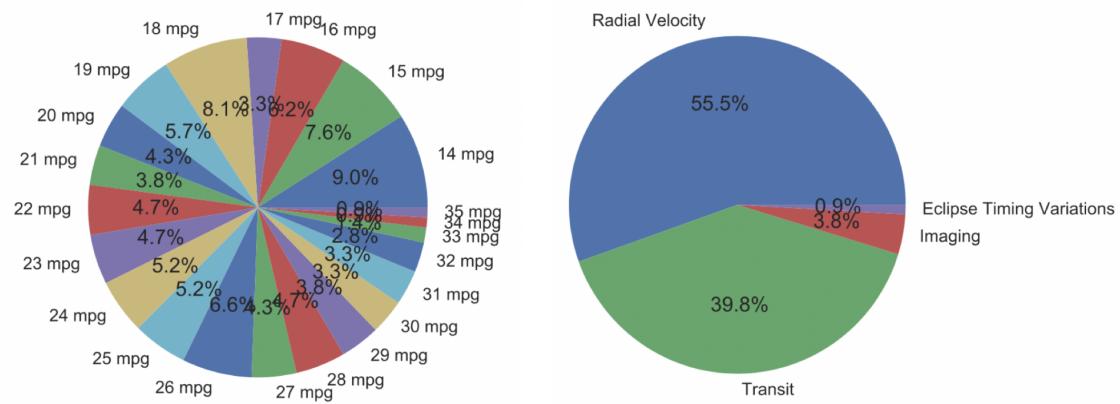
©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## Considerations for data visualization

While data visualization is useful, and even necessary, in exploring and understanding large datasets, a variety of tables and charts are available, and a number of factors must be considered when choosing how to present the data.

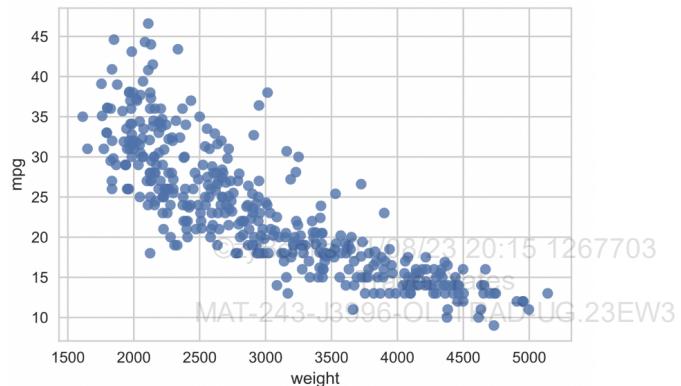
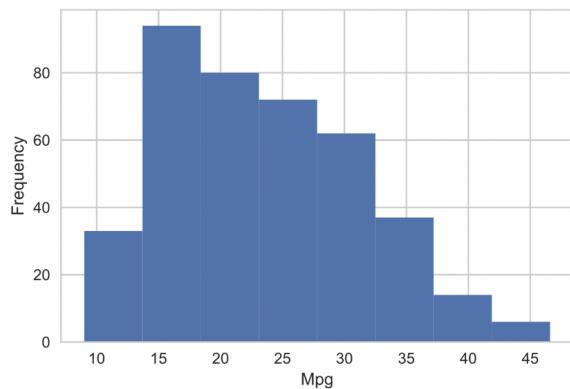
First, the size and cardinality of the dataset must be considered. **Cardinality** is the number of unique elements in a dataset. Ex: the set of student IDs of students in a class has high cardinality, since each ID is unique, whereas the set of student ages will have lower cardinality, since many students will have the same ages. Certain chart types, such as pie charts or bar charts, are well-suited to data with low cardinality, but not well-suited for high-cardinality data, as illustrated by the pie charts below. The chart on the left displays high-cardinality data, and is difficult to read, while the chart on the right displays low-cardinality data.

Figure 1.2.5: Pie charts are better suited for low-cardinality data than high-cardinality data.



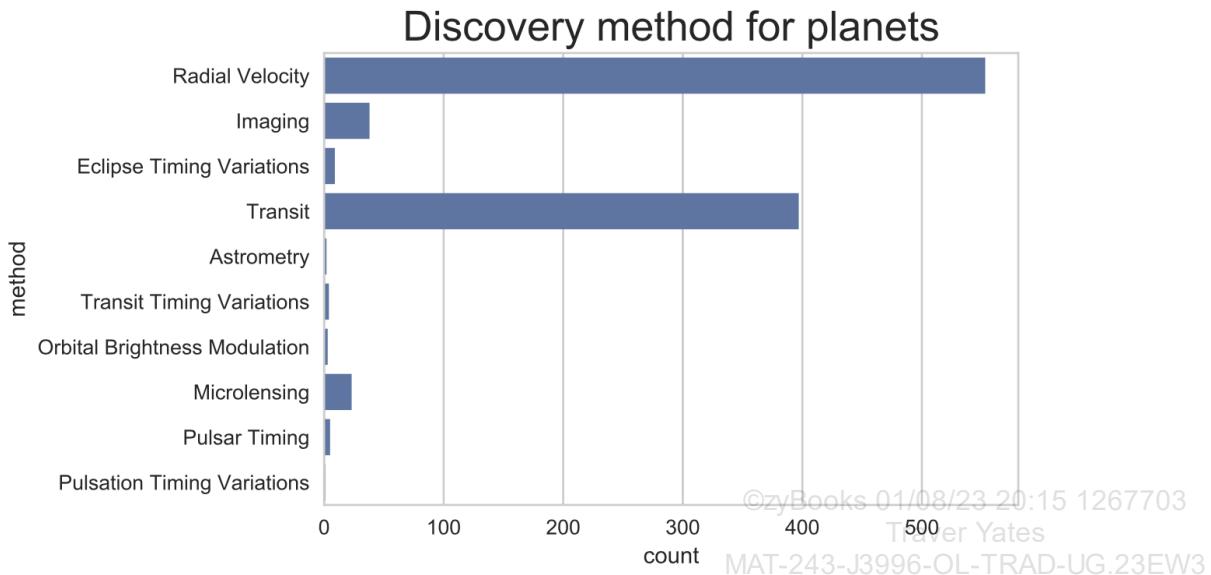
Other charts, such as scatter graphs, line charts, and histograms, work very well for high-cardinality data. In the figure below, the histogram collects the data into eight equal sized bins and shows the distribution of a large number of unique data points. The scatter plot shows the relationship between two variables with high cardinality.

Figure 1.2.6: Histograms and scatter plots are well-suited for high-cardinality data.



The type of chart used also depends on the kind of data being presented, and the information to be conveyed. In the case of a dataset that has only one variable, or where only one variable needs to be presented, can be visualized using a pie chart, histogram, or box plot. A dataset with two or more variables that are related may be best suited to visualization with a type of scatter plot or line chart. A dataset in which one of the variables is categorical works with a bar graph, pie chart, or violin plot. Ex: The bar chart below shows the number of exoplanets discovered using each of ten methods. The method type is a categorical variable.

Figure 1.2.7: Bar charts are appropriate for plotting categorical data.

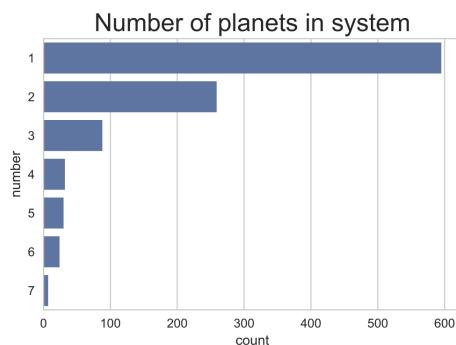


### PARTICIPATION ACTIVITY

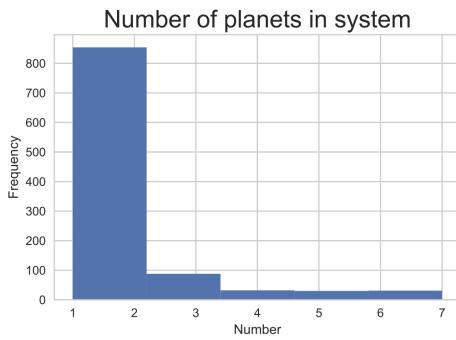
1.2.3: Choosing an appropriate chart.



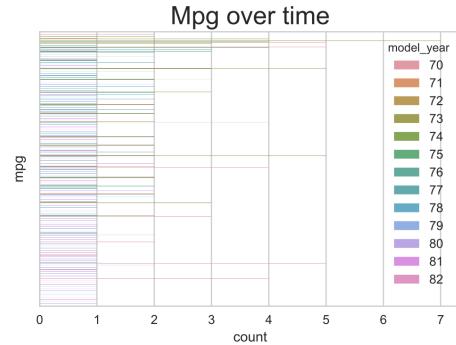
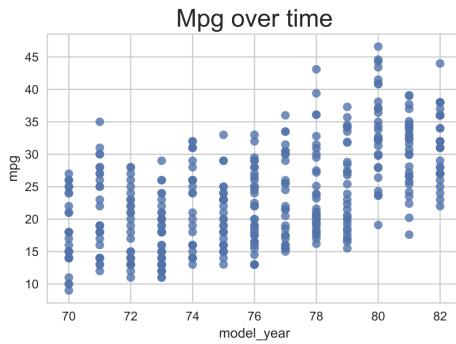
- 1) Which chart better conveys the most common number of discovered planets in a star system?



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



- 2) Which chart better conveys how gas mileage for cars has changed over time?



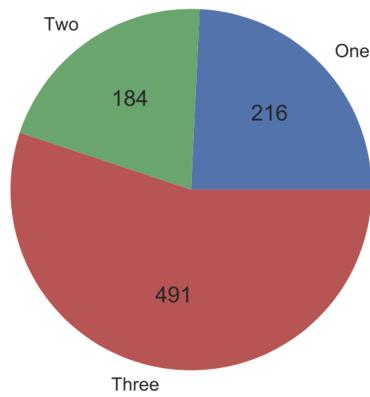
©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- 3) Which chart is more appropriate for showing data from three unrelated categories?





Pie chart



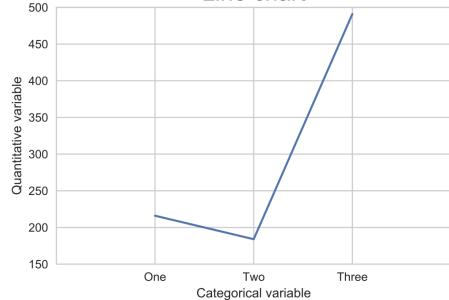
©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



Line chart



## References

(\*) Moore, Karleigh, et al. "Data Presentation - Pie Charts." *Brilliant.org*. Retrieved 16 July 2018, [brilliant.org/wiki/data-presentation-pie-charts/](https://brilliant.org/wiki/data-presentation-pie-charts/).

(\*) USDA Economic Research Service. "Table. dymfg." *United States Department of Agriculture*. [www.ers.usda.gov/data-products/food-availability-per-capita-data-system/](https://www.ers.usda.gov/data-products/food-availability-per-capita-data-system/).

(\*) Macrotrends. "Gold Prices - 100 Year Historical Chart". Last accessed: 1 August 2019. <https://www.macrotrends.net/1333/historical-gold-prices-100-year-chart>.

(\*) Sunshine Profits. "Precious metals investment terms A to Z." Last accessed: 1 August 2019. <https://www.sunshineprofits.com/gold-silver/dictionary/dow-iones-aold/>.

## 1.3 Python for data visualization

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Introduction to Python

Python is a programming language frequently used for data analysis and data visualization.

This material uses Python 3.7. Two major versions of Python are currently supported: Python 2.7 is the most popular and widely supported, while Python 3.7 has the newest features but is not compatible with Python 2.x. More details are available in the official Python documentation.

- [Python 2.7](#)
- [Python 3.7](#)

## Running Python

The **Python interpreter** is a computer program that executes Python code. An **interactive interpreter** is a program that allows the user to execute one line of code at a time.

In many systems, the Python interpreter is started by entering `python` in the command line. In some cases, additional software or administrator privileges may be required. Other Python development environments exist, such as Jupyter, Anaconda, and Cloud9. Some are web-based and do not require the installation of additional software.

## Data types and data structures

Python uses several data types to represent data. The most important data types in Python are shown in the following table.

Table 1.3.1: Basic data types in Python.

Type	Description	Examples
int	An integer number.	0, -1, 4
float	A decimal number or number in scientific notation.	2.7168, -2.0, 1.618, 6.02e23, 1.60e-19
string	A sequence of characters (stored and output surrounded by single or double quotes).	'Hello', 'pie', '3.14'
boolean	A value that is either true or false.	True, False

Data can be organized into more complex data structures known as containers. Common containers are listed in the table below.

Table 1.3.2: Basic data structures in Python.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Type	Description	Examples
set	An unordered collection of	{1, 2, 3}, { 'First name', 'Last name' }

	items.	
list	An ordered collection of changeable items. Two-dimensional arrays can be formed from lists of lists.	[1, 2, 3], ['Height', 6.1], [[1, 2], [3, 4]] ©zyBooks 01/08/23 20:15 1267703 Traver Yates MAT-243-J3996-OL-TRAD-UG.23EW3
tuple	An ordered collection of unchangeable items.	(1, 2, 3), ('Hello', 'World')
dictionary (or dict)	A collection of mappings between keys and values.	{'Home': '310-555-5555', 'Office': '951-555-5555'}, {'C': 'do', 'D': 're', 'E': 'mi', 'F': 'fa', 'G': 'so'}

◀ ▶

PARTICIPATION ACTIVITY
1.3.1: Data types and data structures.


Select the definition that matches each term

1) int

- 42
- 3.14159
- {'CS1': 'Intro to Computer Science', 'Ma1': 'Calc 1'}
- [1, 2, 3, 4, 5, 6, 7]
- 'I love statistics!'

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

2) float

- 42
- 3.14159
- {'CS1': 'Intro to Computer Science', 'Ma1': 'Calc 1'}

[1, 2, 3, 4, 5, 6, 7] 'I love statistics!'

3) string

 42 3.14159 {'CS1': 'Intro to Computer Science', 'Ma1': 'Calc 1'} [1, 2, 3, 4, 5, 6, 7] 'I love statistics!'

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

4) list

 42 3.14159 {'CS1': 'Intro to Computer Science', 'Ma1': 'Calc 1'} [1, 2, 3, 4, 5, 6, 7] 'I love statistics!'

5) dictionary

 42 3.14159 {'CS1': 'Intro to Computer Science', 'Ma1': 'Calc 1'} [1, 2, 3, 4, 5, 6, 7] 'I love statistics!'**Reset**

## Importing modules

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Many commands for data visualization and analysis in Python require the use of **modules**, or external libraries of commands. Modules are imported using the **import** command. The **as** keyword can be used to create an alias. Ex: The two code blocks below perform the same action.

```
# Imports the pandas library
import pandas

# Directly accesses the read_csv function in pandas
unemployment = pandas.read_csv('unemployment.csv')
```

```
# Imports the pandas library using the alias pd
import pandas as pd

# Uses the alias pd to reference the function in pandas
unemployment = pd.read_csv('unemployment.csv')
```

The following table shows modules commonly used in data visualization and analysis, as well as the aliases used in this material for each module. Depending on the specific system or development environment, some modules may need to be installed before the modules can be imported.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Table 1.3.3: Common Python data visualization and analysis modules.

Module name	Alias	Description	Link to documentation
numpy	np	Mathematical functions, required by other libraries	<a href="#">Link</a>
scipy.stats	st	Statistical functions	<a href="#">Link</a>
pandas	pd	Data frames, statistical functions	<a href="#">Link</a>
matplotlib.pyplot	plt	Data visualization	<a href="#">Link</a>
scikit-learn	sks	Machine learning and data analysis	<a href="#">Link</a>
seaborn	sns	Data visualization	<a href="#">Link</a>
quandl		Data analysis	<a href="#">Link</a>

**PARTICIPATION ACTIVITY**

## 1.3.2: Importing modules.



- 1) Enter the command that imports the **pandas** module using the alias in the table above.

**Check****Show answer**

- 2) Enter the command that imports the **matplotlib.pyplot** module using the alias in the table above.



©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

 Show answer

## The matplotlib module

©zyBooks 01/08/23 20:15 1267703

The **matplotlib** module can be used for plotting in Python. matplotlib, short for MATLAB plotting library, replicates the plotting capability of MATLAB, an engineering-oriented programming language.

matplotlib is not included with Python, but can be downloaded and installed from

<http://matplotlib.org/downloads.html>. matplotlib also requires the [NumPy](#) module, which includes many general mathematical functions.

The `plt.plot()` function plots data onto a graph. Different types of graphs can be specified, as discussed elsewhere. After the data is plotted, the graph can be displayed on the console or saved to a file. The `plt.show()` function displays the graph. The `plt.savefig()` function saves the graph to a file.

The following table shows functions commonly used to label and annotate matplotlib plots.

Table 1.3.4: matplotlib functions for labeling and annotating plots.

Function	Description
<code>title()</code>	Specifies title of plot
<code>xlabel()</code>	Specifies $x$ -axis label
<code>ylabel()</code>	Specifies $y$ -axis label
<code>text()</code>	Creates a text label
<code>annotate()</code>	Creates a text label for a specific data point
<code>legend()</code>	Creates a legend for the plot

©zyBooks 01/08/23 20:15 1267703

Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The following program shows the use of the matplotlib module to plot a simple line graph and add axis labels and titles.

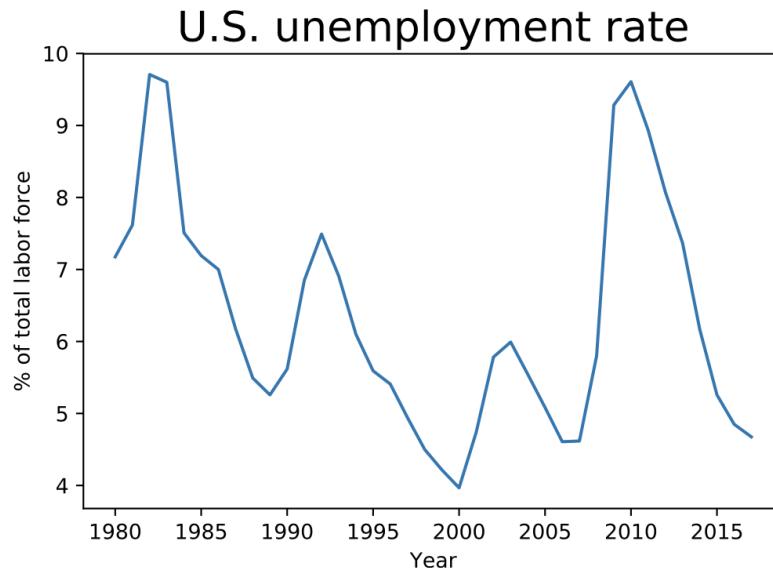
### Python-Function 1.3.1: Unemployment rates.

The code below loads a dataset containing unemployment rates in the United States from 1980 to 2017 and plots the data as a line chart.

```
# imports the necessary libraries
import pandas as pd
import matplotlib.pyplot as plt

# loads the unemployment dataset
unemployment = pd.read_csv('http://data-analytics.zybooks.com/unemployment.csv') 20:15 1267703
# title
plt.title('U.S. unemployment rate', fontsize = 20)
# Traver Yates
# x and y axis labels
plt.xlabel('Year')
plt.ylabel('% of total labor force')
# plot
plt.plot(unemployment["Year"], unemployment["Value"])
# saves the image
plt.savefig("unemployment.png")
# shows the image
plt.show()
```

The resulting graph is shown below.



[Run example](#)

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

#### PARTICIPATION ACTIVITY

1.3.3: Using matplotlib to set plot parameters.



In the following questions, assume that the `matplotlib.pyplot` module has already been imported as `plt`.



- 1) Enter the command that sets the  $x$ -axis label of a plot to 'Height'.

**Check****Show answer**

- 2) Enter the command that sets the title of a plot to 'Height distribution' and sets the font size to 16.

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267700  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## Multiple plots

Multiple sets of data can be plotted on the same axis. The following code plots two categories named `total` and `speeding` from the dataset named `car_crashes` on the same set of axes, forming a stacked bar graph.

### Python-Function 1.3.2: Automobile collisions.

The code below plots the number of automobile collisions related to speeding and the total number of collisions on the same set of axes. The data is presented as a stacked bar chart in which two bar charts are overlaid on each other.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

```
# loads the necessary modules
import matplotlib.pyplot as plt
import seaborn as sns

# initialize figure
f, ax = plt.subplots()

# load dataframe
crashes = sns.load_dataset("car_crashes")
df = crashes.loc[range(5)]

# plot total crashes
sns.set_color_codes("pastel")
sns.barplot(x="total", y="abbrev", data=df,
            label="Total", color="b")

# plot crashes related to speeding
sns.set_color_codes("muted")
sns.barplot(x="speeding", y="abbrev", data=df,
            label="Speeding-related", color="b")

# title
plt.title('Speeding-related automobile collisions', fontsize=20)

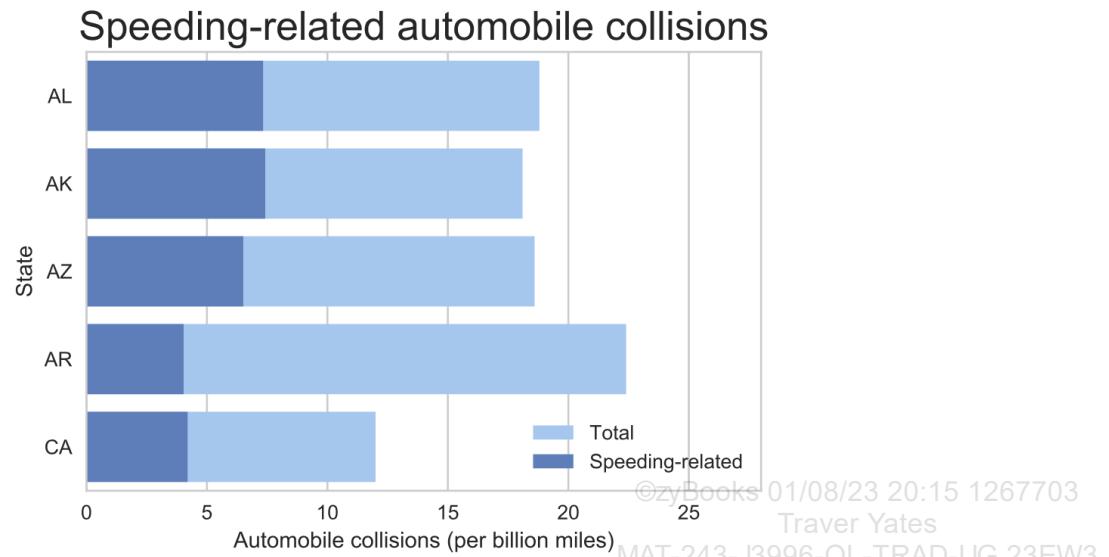
# legend
ax.legend(ncol=1, loc="lower right")
ax.set(xlim=(0, 28), ylabel="State", xlabel="Automobile collisions (per billion miles)");

# saves the image
plt.savefig("stacked.png")

# shows the image
plt.show()
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The resulting stacked bar chart is shown below.



[Run example](#)





1) In the stacked bar graph above, the bars for total collisions are overlaid on top of the bars for speeding-related collisions.

- True
- False

2) The length of a total bar not covered by a speeding-related collisions bar represents collisions not related to speeding.

- True
- False

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



## 1.4 Data frames

### 1 Introduction to data frames

A **data frame**, sometimes typeset as DataFrame, is a two-dimensional tabular data structure with labeled columns and rows. Similar to spreadsheets, a data frame provides a way to store data in a way that's easily sorted and categorized. A data frame has three components: (1) *index*, (2) *columns*, and (3) *values*. A data frame's **index** is the set of row labels. The index can have a name, which would be located in the first row and column of the data frame. A data frame's **columns** are the labels of the column data. The data contained in a data frame are also known as **values**. Although a value may indicate numerical data such as float or int, a value can also be a string or datetime.

PARTICIPATION  
ACTIVITY

1.4.1: Data frames.



### Animation content:

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Animation captions:

1. A data frame is a two-dimensional structure in which each column contains values of one variable and each row contains one set of values from each column.
2. The columns are the labels or names of each column.
3. The index is the set of labels or names of each row.

4. The data contained in each column are called values. The data type of the values in the column size is an integer.

**PARTICIPATION ACTIVITY**
**1.4.2: Data frames.**


Refer to the data frame below lists the mean income in five states in 2005 and 2006:  
Alabama, Alaska, Arizona, Arkansas, and California.

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

<b>GEOID</b>	<b>2005</b>	<b>2006</b>	
<b>State</b>			
<b>Alabama</b>	04000US01	37150	37952
<b>Alaska</b>	04000US02	55891	56418
<b>Arizona</b>	04000US04	45245	46657
<b>Arkansas</b>	04000US05	36658	37057
<b>California</b>	04000US06	51755	55319

<b>Alabama</b>	04000US01	37150	37952
<b>Alaska</b>	04000US02	55891	56418
<b>Arizona</b>	04000US04	45245	46657
<b>Arkansas</b>	04000US05	36658	37057
<b>California</b>	04000US06	51755	55319

- 1) Which of the following is a column label?



- State
- Alabama
- GEOID

- 2) Which of the following is an index label?



- Alaska
- 36658
- 2006

- 3) Which of the following is the data type for GEOID?



- int
- float
- string

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## pandas

**pandas** is a Python library that allows a user to work with data frames by providing tools for reading, writing, subsetting, and reshaping data. **pandas** is especially useful when the dataset has missing and unaligned data, which is common in real-world applications.

To work with data using **pandas**, a file needs to be imported as a **DataFrame** object. Although **DataFrame** objects can be created, this material works with existing data. Sample code to import files of various types is shown below.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

Figure 1.4.1: Importing files as DataFrame objects.MAT-243-J3996-OL-TRAD-UG.23EW3

```
# imports the pandas library
import pandas as pd

# loads a file containing comma-separated values and assigns
# the data frame to variable DataFrame
data_frame1 = pd.read_csv('file.csv')

# loads text file where the values are separated by a space with no column
labels
data_frame2= pd.read_csv('file.txt', sep = ' ', header = None)

# loads an excel file
data_frame3 = pd.read_excel('file.xlsx', sheetname='Sheet1')
```

An object in Python has attributes and methods. An **attribute** is a characteristic of an object. A **method** is a procedure associated with an object. Methods are followed by parentheses `()`. The tables below give a partial list of **DataFrame** methods and attributes. The full list can be obtained in the [pandas documentation](#).

Table 1.4.1: DataFrame attributes.

DataFrame.attribute	Description of output
axes	Index and column labels
columns	Column labels
dtypes	Data types of values in each column
index	Index labels
shape	Ordered pair that gives the number of rows and columns
size	Number of values in the DataFrame

values	Values in the DataFrame
--------	-------------------------

Table 1.4.2: DataFrame methods.

©zyBooks 01/08/23 20:15 1267703

DataFrame.method	Description of output	Traver Yates 3-J3996-OL-TRAD-UG.23EW3
<code>describe()</code>	Summary statistics for numerical columns	
<code>head(), tail()</code>	First/last 5 rows in the DataFrame	
<code>min(), max()</code>	Minimum/maximum of values in a numerical column	
<code>mean(), median()</code>	Mean/median of values in a numerical column	
<code>sample()</code>	Random row	
<code>std()</code>	Standard deviation of values in a numerical column	

## Example 1.4.1: Titanic dataset.

A popular dataset for learning data exploration and analysis is the dataset containing various information about the passengers of the Titanic. `titanic` is one of the datasets included in the `seaborn` package and can be loaded using the code below.

```
# loads the seaborn library
import seaborn as sns

# loads the titanic dataset
titanic = sns.load_dataset("titanic")
```

What commands are needed to do find the following?

1. Number of rows
2. Column names
3. Data types of columns

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Solution**

1. To find the number of rows, the command `titanic.shape` can be used to find the number of rows and the number of columns. The output is the ordered pair **(891, 15)** where **891** is the number of rows and **15** is the number of columns.
2. To find the column names, the command `titanic.columns` is used, which gives the output

```
Index(['survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare',
       'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town',
       'alive', 'alone'],
      dtype='object')
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

3. To find the data types of the columns , the command `titanic.dtypes` is used, which gives the output

<code>survived</code>	<code>int64</code>
<code>pclass</code>	<code>int64</code>
<code>sex</code>	<code>object</code>
<code>age</code>	<code>float64</code>
<code>sibsp</code>	<code>int64</code>
<code>parch</code>	<code>int64</code>
<code>fare</code>	<code>float64</code>
<code>embarked</code>	<code>object</code>
<code>class</code>	<code>category</code>
<code>who</code>	<code>object</code>
<code>adult_male</code>	<code>bool</code>
<code>deck</code>	<code>category</code>
<code>embark_town</code>	<code>object</code>
<code>alive</code>	<code>object</code>
<code>alone</code>	<code>bool</code>
<code>dtype:</code>	<code>object</code>

In pandas, the data types int and float are `int64` and `float64` respectively. **64** is a reference to memory allocation, which is beyond the scope of this material. Similarly, the string data type is called `object` in pandas. `bool` stands for boolean, which can have a value of either true or false (equivalently, `1` or `0`).

#### PARTICIPATION ACTIVITY

#### 1.4.3: Titanic data set.



Refer to the example above that includes a `DataFrame` variable named `titanic`.

- 1) Write a command that finds the mean of all numerical variables in `titanic`.



`titanic.`  `//`  
( )

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

**Check**

**Show answer**



- 2) Write a command that finds the minimum values of all columns in **titanic**.

```
titanic.  
()
```

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- 3) Using the **code** from the question above, find the minimum age of all passengers in **titanic**. Type as:

```
#.##
```

**Check****Show answer**

## Subsetting data

**Subsetting** is the process of retrieving parts of a data frame. A data frame can be subset in a variety of ways, the most common of which involve selecting a range of rows and columns, or selecting columns by column label. Ex: In the Titanic data set, a data exploration project may involve finding the sex of people on deck B, so a data analyst would filter the rows for which the value of deck is B.

### Python-Function 1.4.1: Selecting columns and rows.

To select a column or columns of a data frame, the command `data_frame["column"]` where `data_frame` is the `DataFrame` object and `column` is the name of the selected column. The column name is enclosed within single or double quotes. Alternatively, the command `data_frame.column` can be used as long as the column name is not the same as an existing method such as `rank()`. The data type of the output is a series, which can be thought of as a list with labels.

To select more than one column or make the output a data frame rather than a list, double brackets should be used with the command above.

`data_frame[["column1", "column2", ...]]` returns a data frame with `column1`, `column2`, ... included.

The code below returns a data frame with only the columns `sex` and `survived` from the Titanic data set.

```
titanic[["sex", "survived"]]
```

The first 5 rows of the resulting data frame are shown in the output below.

```
titanic[["sex","survived"]].head()
```

	sex	survived
0	male	0
1	female	1
2	female	1
3	female	1
4	male	0

©zyBooks 01/08/23 20:15 1267703

To select a row by position, the command `data_frame[a:b]` where `a` and `b` are the initial and final rows included in the output.

MAT-243-J3996-OL-TRAD-UG.23EW3

The code below selects the first 5 rows of the titanic dataset and returns a data frame.

```
titanic[0:5]
```

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	\
0	0	3	male	22.0	1	0	7.2500	S	Third	
1	1	1	female	38.0	1	0	71.2833	C	First	
2	1	3	female	26.0	0	0	7.9250	S	Third	
3	1	1	female	35.0	1	0	53.1000	S	First	
4	0	3	male	35.0	0	0	8.0500	S	Third	
	who	adult_male	deck	embark_town	alive	alone				
0	man	True	NaN	Southampton	no	False				
1	woman	False	C	Cherbourg	yes	False				
2	woman	False	NaN	Southampton	yes	True				
3	woman	False	C	Southampton	yes	False				
4	man	True	NaN	Southampton	no	True				

To select only the rows where one column has a particular value, the command `data_frame[data_frame.column == x]` is used, where `x` is the desired value in `column`.

The code below selects all the rows where the column `pclass` has a value of 3.

```
titanic[titanic.pclass == 3]
```

	survived	pclass	sex	age	sibsp	parch	fare	...	class	who	adult_male
deck	embark_town	alive	alone								
0	0	3	male	22.0	1	0	7.2500	...	Third	man	True
NaN	Southampton	no	False								
2	1	3	female	26.0	0	0	7.9250	...	Third	woman	False
NaN	Southampton	yes	True								
4	0	3	male	35.0	0	0	8.0500	...	Third	man	True
NaN	Southampton	no	True								
5	0	3	male	NaN	0	0	8.4583	...	Third	man	True
NaN	Queenstown	no	True								
7	0	3	male	2.0	3	1	21.0750	...	Third	child	False
NaN	Southampton	no	False								
...	...	...	...	...	...	...	...	...	...	...	...
882	0	3	female	22.0	0	0	10.5167	...	...	...	...
NaN	Southampton	no	True								
884	0	3	male	25.0	0	0	7.0500	...	...	...	...
NaN	Southampton	no	True								
885	0	3	female	39.0	0	5	29.1250	...	Third	woman	False
NaN	Queenstown	no	False								
888	0	3	female	NaN	1	2	23.4500	...	Third	woman	False
NaN	Southampton	no	False								
890	0	3	male	32.0	0	0	7.7500	...	Third	man	True
NaN	Queenstown	no	True								

Similarly the comparison operators `>`, `<`, `>=`, and `<=` can be used to select rows where a column has values more or less than a given value.

The code below selects all the rows where the column `fare` has a value less than `10.0`.

```
titanic[titanic.fare < 10.0]
```

	survived	pclass	sex	age	sibsp	parch	fare	...	class	who	adult_male
0	0	3	male	22.0	1	0	7.2500	...	Third	man	True
NaN	Southampton	no	False								
2	1	3	female	26.0	0	0	7.9250	...	Third	woman	False
NaN	Southampton	yes	True								
4	0	3	male	35.0	0	0	8.0500	...	Third	man	True
NaN	Southampton	no	True								
5	0	3	male	NaN	0	0	8.4583	...	Third	man	True
NaN	Queenstown	no	True								
12	0	3	male	20.0	0	0	8.0500	...	Third	man	True
NaN	Southampton	no	True								
...	...	...	...	...	...	...	...	...	...	...	...
...	...	...	...	...	...	...	...	...	...	...	...
877	0	3	male	19.0	0	0	7.8958	...	Third	man	True
NaN	Southampton	no	True								
878	0	3	male	NaN	0	0	7.8958	...	Third	man	True
NaN	Southampton	no	True								
881	0	3	male	33.0	0	0	7.8958	...	Third	man	True
NaN	Southampton	no	True								
884	0	3	male	25.0	0	0	7.0500	...	Third	man	True
NaN	Southampton	no	True								
890	0	3	male	32.0	0	0	7.7500	...	Third	man	True
NaN	Queenstown	no	True								

[336 rows x 15 columns]

## Python-Function 1.4.2: loc() and iloc()

The `loc()` is used to select a range of rows and/or a subset of columns. Ex:

`titanic.loc[0:5,[ "pclass", "age"]]` returns a DataFrame containing the first 6 rows and the columns `pclass` and `age` of the `titanic` data set as shown below.

	pclass	age
0	3	22.0
1	1	38.0
2	3	26.0
3	1	35.0
4	3	35.0
5	3	NaN

The `iloc()` is used to select a range of rows and/or columns. Ex:

`titanic.iloc[0:5,0:5]` returns a DataFrame containing the first 5 rows and the first 5 columns of the titanic data set as shown below.

	survived	pclass	sex	age	sibsp
0	0	3	male	22.0	1
1	1	1	female	38.0	1
2	1	3	female	26.0	0
3	1	1	female	35.0	1
4	0	3	male	35.0	0

**PARTICIPATION ACTIVITY**

## 1.4.4: Subsetting data.



1) Which command selects the column `age` from the Titanic data frame and returns a data frame?

- `titanic.age`
- `titanic["age"]`
- `titanic[["age"]]`
- `titanic[[age]]`

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



2) What pandas object is returned by the command

`titanic[["age", "fare", "alive"]]`

- `Series`
- `DataFrame`



3) What command returns the first 10 rows of the `titanic` data frame?

- `titanic[1:10]`
- `titanic[0:10]`
- `titanic[[0:10]]`



4) What command returns the first 5 rows of the data frame and the columns `age`, `fare`, and `class`?

- `titanic.loc[0:4, ["age", "fare", "class"]]`
- `titanic.iloc[0:4, ["age", "fare", "class"]]`
- `titanic[0:4, ["age", "fare", "class"]]`

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



## Reshaping data

Two common forms of a data frame exist: *long form* and *wide form*. A data frame is in ***long form*** when each column is a variable and each row gives non-repeated data. A data frame in long form is also referred to as unstacked or in record form. A data frame is in ***wide form*** if each data variable is in a different column. A data frame in wide form is also referred to as stacked. **Reshaping data** involves

converting a data frame from one form into the other. **Pivoting** converts a data frame from long form to wide form. **Melting** converts a data frame from wide form to long form.

**PARTICIPATION  
ACTIVITY**

## 1.4.5: Pivoting.

**Animation content:**

undefined

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Animation captions:**

1. In long form, each row represents an observation. Ex: Each row in the data frame is a test taken by one student on a particular date.
2. Pivoting converts a data frame in long form into wide form. Ex: The values of the variable "Student" is used as columns and dates as index labels.
3. The values are from a specified column. Ex: The values in the data frame are from the column "Score".

**Python-Function 1.4.3: pivot().**

The command used to pivot the data frame, df, in the animation above is as follows.

```
pd.pivot(df, index = "Date", columns = "Student", values= "Score")
```

The index labels of the pivoted data frame are the unique values of the column **Date**. The columns of the pivoted data frame are the unique values of the column **Student**. The values of the pivoted data frame are the values from the column **Score**.

The code to generate the data frame in long form and the corresponding output are given below.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

```
# Creates a data frame
df = pd.DataFrame({
...
...
    'Student': {0: "Michael", 1: "Arushi", 2: "Roberta",
    3: "Michael", 4: "Arushi", 5: "Roberta"}, 
    'Score': {0: 90, 1: 98, 2: 92, 3: 90, 4: 98, 5: 92}, 
    'Date': {0: "2018-01-03", 1: "2018-01-03", 2: "2018-01-03", 3:
    "2018-01-13",
    ...
    4: "2018-01-13", 5: "2018-01-13"}}, 
    columns=['Date','Student','Score'])
```

	Date	Student	Score
0	2018-01-03	Michael	90
1	2018-01-03	Arushi	98
2	2018-01-03	Roberta	92
3	2018-01-13	Michael	90
4	2018-01-13	Arushi	98
5	2018-01-13	Roberta	92

©zyBooks 01/08/23 20:15 1267703

The output of the `pivot()` command is given below.

Student	Arushi	Michael	Roberta
Date			
2018-01-03	98	90	92
2018-01-13	98	90	92

## PARTICIPATION ACTIVITY

#### 1.4.6: Melting.

## **Animation content:**

undefined

## Animation captions:

1. In wide form, each data variable is in a different column. Ex: Each row gives data for each student where the columns are "Student" and "SAT"/"ACT" scores.
  2. Melting converts a data frame in wide form into long form. Ex: Each row represents a different test taken by a student and the columns are "Student", "Test", and "Scores".
  3. New column labels are specified. Ex: "Test" is the column label for values SAT and ACT and the label for the values from the pivoted data frame is "Scores".

## Python-Function 1.4.4: melt()

©zyBooks 01/08/23 20:15 1267703

The command used to melt the data frame `df` in the animation above is as follows.

```
pd.melt(df, id_vars = "Student", var_name = "Test", value_vars = ["SAT", "ACT"], value_name = "Scores")
```

The `id_vars` parameter indicates a string, tuple, or list to use as an identifier variable . The `var_name` parameter indicates the column label to use for the variable column. The parameter `value_vars` indicates the labels of the columns to unpivot. The `value_name` indicates the column label to use for the value column.

The code to generate the data frame in wide form and the corresponding output are given below.

```
df = pd.DataFrame({'Student': {0: "Michael", 1: "Arushi", 2: "Roberta"},  
...                 'SAT': {0: 1480, 1: 1520, 2: 1460},  
...                 'ACT': {0: 34, 1: 32, 2: 32}},  
...                 columns=['Student','SAT','ACT'])
```

	Student	SAT	ACT
0	Michael	1480	34
1	Arushi	1520	32
2	Roberta	1460	32

The output of the `melt()` command is given below.

	Student	Test	Scores
0	Michael	SAT	1480
1	Arushi	SAT	1520
2	Roberta	SAT	1460
3	Michael	ACT	34
4	Arushi	ACT	32
5	Roberta	ACT	32

#### PARTICIPATION ACTIVITY

#### 1.4.7: pivot() and melt().

- 1) What command is needed to reshape the first data frame `df` into the second data frame?

	bar	baz	foo
0	A	1	one
1	B	2	one
2	C	3	one
3	A	4	two
4	B	5	two
5	C	6	two

	foo	one	two
bar			
A		1	4
B		2	5
C		3	6

```
(df,  
index="bar", columns="foo",  
values="baz")
```

**Check**

**Show answer**

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- 2) What command is needed to reshape the first data frame, `df`, into the second data frame?

	A	B	C
0	a	1	2
1	b	3	4
2	c	5	6

	A variable	value
0	a	2
1	b	4
2	c	6

```
(df,
id_vars="A", value_vars= "C")
```

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- 3) To use the label **Counts** as the label for the value column, what should be added to the command above?

**Check****Show answer**

## 1.5 Bar charts



This section has been set as optional by your instructor.

### Bar charts

A **bar chart** depicts data values for a categorical variable, using rectangular bars having lengths proportional to category values. The chart is drawn using two axes: a category axis that displays the category names and a value axis that displays the counts. Ex: The animation below shows the number of employees of each of the 4 largest private employers in the United States in 2017.<sup>1</sup>

**PARTICIPATION  
ACTIVITY**

1.5.1: Bar chart.



### Animation content:

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

undefined

### Animation captions:

1. A bar chart shows the counts in each category. The categories and the corresponding counts are obtained from a table.

2. The category axis of a bar chart provides a label for each category.
3. The value axis indicates the data values.
4. Each category's value is shown using a bar with appropriate height.
5. All charts should have a title.

Categories are commonly ordered along the category axis. In the animation, the nominal variable Company's categories were ordered by each category's data value, highest (Wal-Mart's **2,300,000**) to lowest (Kroger's **449,000**). If instead the categories represented years (1970, 1980, etc.) or some other measure of time, such an ordinal variable's categories would be ordered with time increasing to the right.

Each listed category has a **category label**, such as "Amazon" above. If labels don't fit when written horizontally, the labels can be rotated, such as rotated 30 degrees as above. Rotations of 60 or 90 degrees are also common.

The appropriate increment for the value axis is important for readability. Using small increments clutters the chart with too many values: Above, an increment of **100,000** would yield **25** values. Too few, like increments of **1,000,000** above, can make visually estimating a category's value difficult. Above, using increments of **500,000** leads to values that can easily be estimated (Wal-Mart can be seen to have a value of about **2,300,000**) without clutter.

Grid lines help the viewer estimate the value for a category. If easier estimation is desired, additional grid lines can be drawn between number increments (but kept minimal to reduce clutter). If precise values need to be conveyed, data values known as **data labels** can be shown next to the bars, or even inside the bars, as in the chart further below. However, precise values are not typically needed if the information goal is to show relative differences among categories.

PARTICIPATION ACTIVITY

1.5.2: Bar chart basics.



1) A bar chart excels at showing precise values.



True

False

2) A bar chart excels at showing relative values.



True

False

3) The more gridlines drawn in a bar chart, the better.



True

False

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- 4) A data label is a category name, such as "IBM" in the above example.

True  
 False

- 5) In the above bar chart on U.S.

employers, the category axis is named "Employees".

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

True  
 False

- 6) For the above U.S. employers chart, categories were ordered by value. An alphabetical ordering would have been just as good.



True  
 False

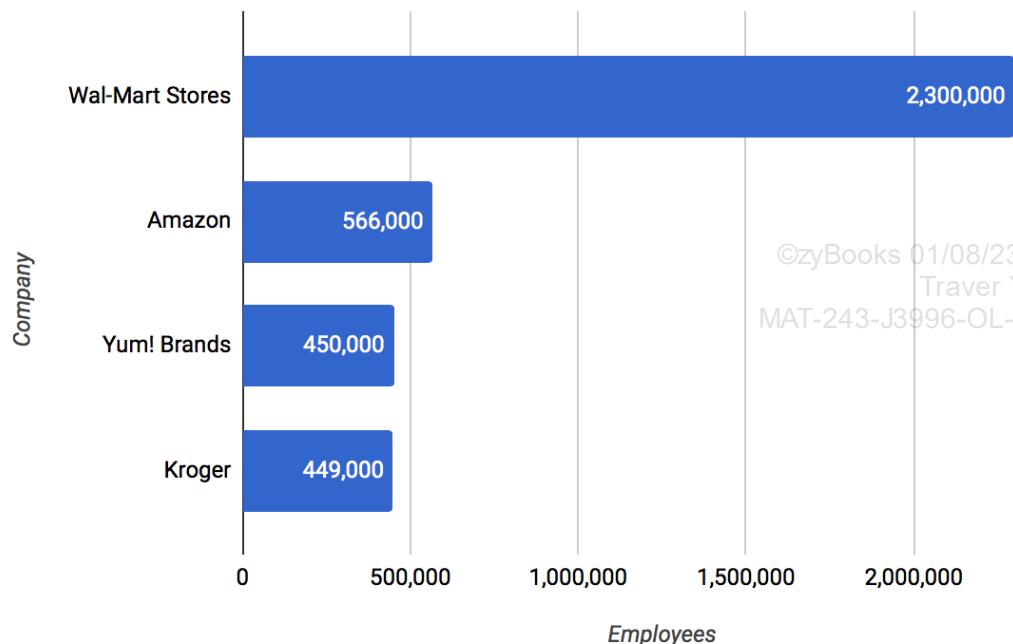
## Horizontal bar charts

A bar chart can be drawn vertically or horizontally. A horizontal bar chart is useful for long labels, like "Wal-Mart Stores", which need not be written at an angle as was done above. A horizontal bar chart is also useful when numerous categories exist because the categories increase the height rather than width, and due to the nature of paper and computers, width is usually more limited while height is less limited. In contrast, a vertical chart is often preferred due to "height" intuitively representing amount. This preference is especially the case when negative values are shown (which would appear going downwards).

Some authors and tools use the term "bar chart" to refer exclusively to a horizontal bar chart. In that case, a **column chart** is a term used for a vertical bar chart. However, the term bar chart is widely used for vertical charts by many respected authors and tools. Thus, this material uses the term bar chart for either orientation, adding the word "horizontal" or "vertical" as appropriate.

Figure 1.5.1: A horizontal bar chart with data labels for the largest private U.S. employers in 2016.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

**PARTICIPATION ACTIVITY****1.5.3: Horizontal bar charts.**

- 1) A horizontal bar chart may be preferable if the category labels are long.



True  
 False

- 2) A horizontal bar chart may be preferable if many categories exist.



True  
 False

- 3) A horizontal bar chart may be preferable to depict different companies' annual profits, some of which are negative.



True  
 False

- 4) A horizontal bar chart may be preferable for conveying the number of floors for the world's 5 tallest buildings.



- True
- False

## Python-Function 1.5.1: Bar charts.

©zyBooks 01/08/23 20:15 1267703

The `sns.countplot(x = "category", data = DataFrame, ...)` function from the `seaborn` library plots categorical data with respect to different categories. For bar charts with a single category that do not include multiple groups, setting `x` equal to a category creates a vertical bar chart and setting `y` equal to a category creates a horizontal bar chart.

```
# loads the necessary modules
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# loads the titanic dataset
titanic = sns.load_dataset("titanic")

# sets the style of the bar charts
sns.set(style="whitegrid", color_codes=True)

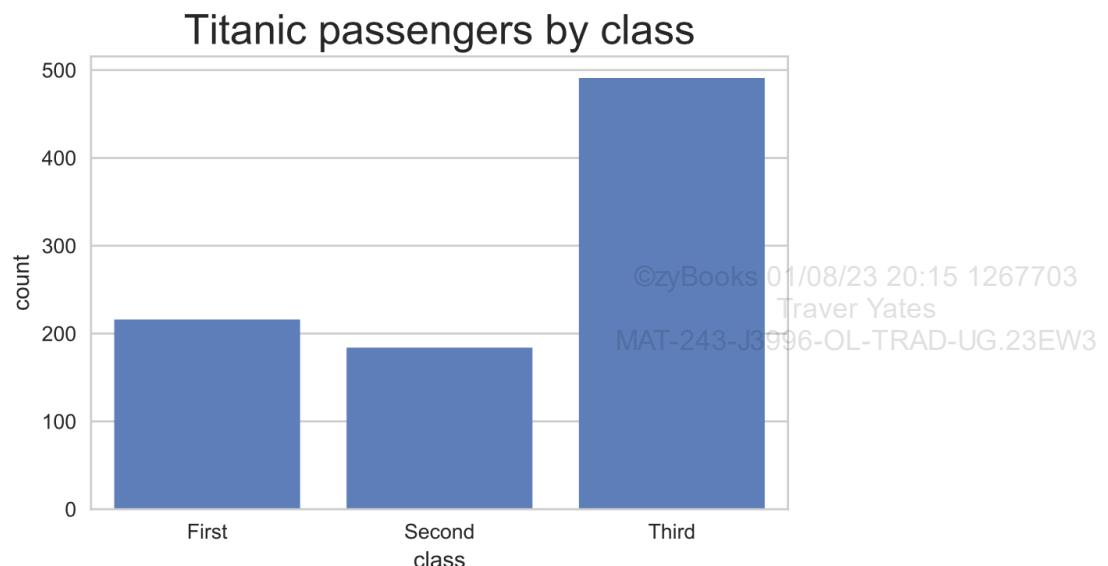
# title
plt.title('HMS Titanic passengers by class', fontsize=20)

# plots a vertical bar chart
sns.countplot(x="class", color="b", data=titanic);

# saves the image
plt.savefig("verticalbarchart.png")

# shows the image
plt.show()
```

The resulting vertical bar chart is shown below.

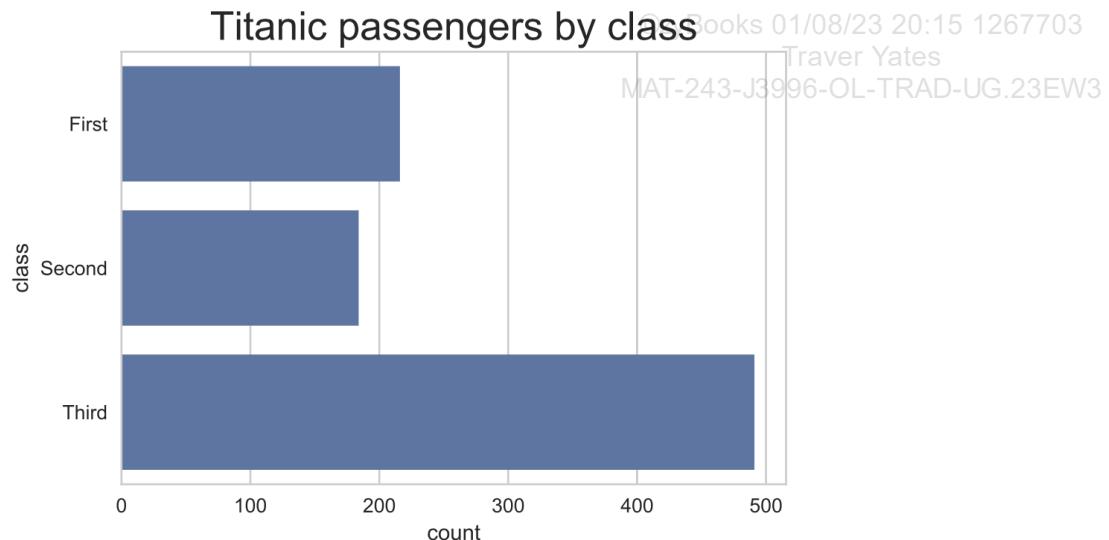


[Run example](#)

The code below generates a horizontal bar chart.

```
# plots a horizontal bar chart
sns.countplot(y="class", color="b", data=titanic);
```

The resulting horizontal bar chart is shown below.



[Run example](#)

## Relative-frequency bar chart

A basic bar chart's value axis provides the raw data value for each category. Instead, a **relative-frequency bar chart** shows each category's portion of the total data, typically as a percentage. The data total is first computed, then the percentage for each category is computed, and finally those percentages are drawn as a bar chart.

PARTICIPATION  
ACTIVITY

1.5.4: Relative-frequency bar chart.

### Animation content:

undefined

### Animation captions:

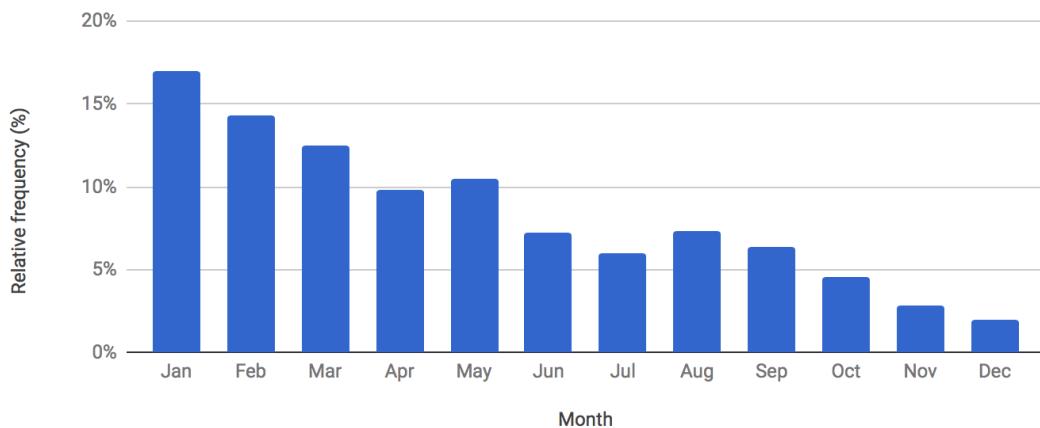
©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

1. A relative frequency bar chart shows the percent that each category is of the total. First, the total is calculated.
2. Next, the value of each category is divided by the total and multiplied by 100 to find the relative frequency as a percentage of the total.
3. The value axis can stop short of 100%, but should usually begin at 0%.
4. Finally, each category's bar is drawn.

## Example 1.5.1: Professional European soccer player birth months.

The following relative-frequency bar chart depicts the birth months for professional European soccer players. Most players were born in the first few months of the year. A possible explanation is the January 1 cutoff date for youth soccer leagues, meaning kids born in January are the oldest on their teams, while kids born in December are the youngest. Older kids are likely to be better players initially, causing coaches to give them more attention and playing time, and also causing those kids to enjoy playing and thus practicing more.<sup>2</sup>

**Professional European soccer player birth months**



The above example illustrates the power of data analytics and data visualization. Having understood such data, many parents now choose to postpone their child's entry onto a sports team so that the child is not the youngest on the team. In fact, similar data exists for school kids: Kids born just after cutoff dates tend to be more successful in school (getting more attention from teachers, and causing those kids to feel smarter and enjoy school more), and thus many parents now delay their child's school enrollment by a year.<sup>3</sup>

**PARTICIPATION ACTIVITY**

1.5.5: Relative-frequency bar charts.



Consider a dataset composed of 4 small shirts, 10 medium shirts, and 6 large shirts.

- 1) If the dataset is represented by a bar chart, what is the height of the bar representing the small shirts?



**Check**

**Show answer**

Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EV



- 2) If the dataset is represented by a relative-frequency bar chart, what is the height of the bar representing the small shirts?

%

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- 3) If the dataset is represented by a relative-frequency bar chart, what is the height of the bar representing the medium shirts?

%

**Check****Show answer**

## Grouped bar chart

A **grouped bar chart** depicts two or more groups on a single bar chart, with each group using a different colored (or shaded) bar. A **legend** indicates what group each color represents in a chart. Ex: The below bar chart shows the number of men vs. women in the U.S. workforce over time. The categories are decades (1970, 1980, ...), the category values are number of people, and the two groups are men and women.

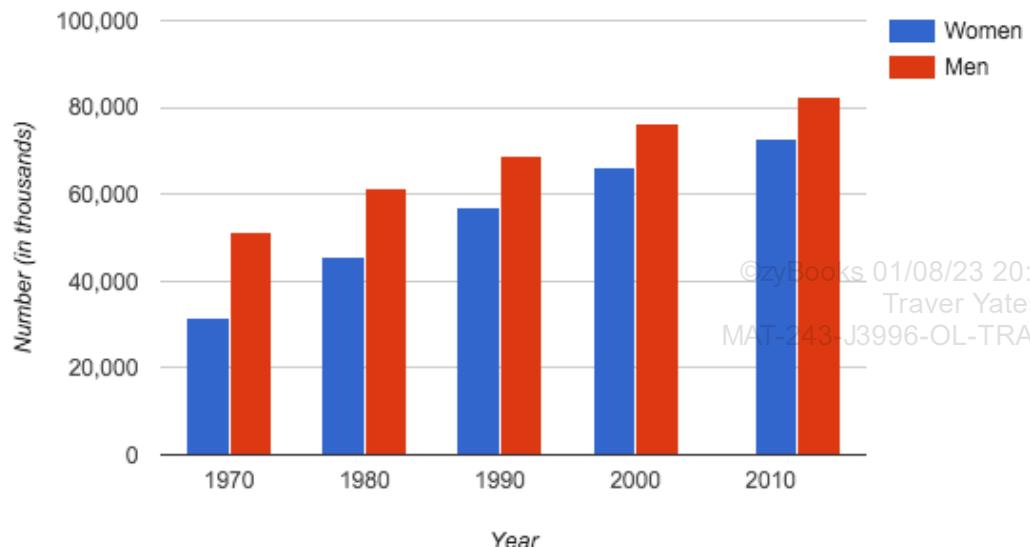
Because the categories represent time (decades), a vertical chart is preferred so that time proceeds to the right.

Figure 1.5.2: Sex of the U.S. workforce.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



Source: U.S. Dept. of Labor <sup>4</sup>

#### PARTICIPATION ACTIVITY

##### 1.5.6: Grouped bar chart.



Consider the above grouped bar chart showing men and women in the U.S. workforce.

- 1) How does the number of women in the workforce compare with the number of men for each decade?



- Fewer
- More
- Same

- 2) How has the difference between men and women changed as time has progressed?



- Increased
- Decreased
- Same

- 3) How has the total number of workers changed as time has progressed?

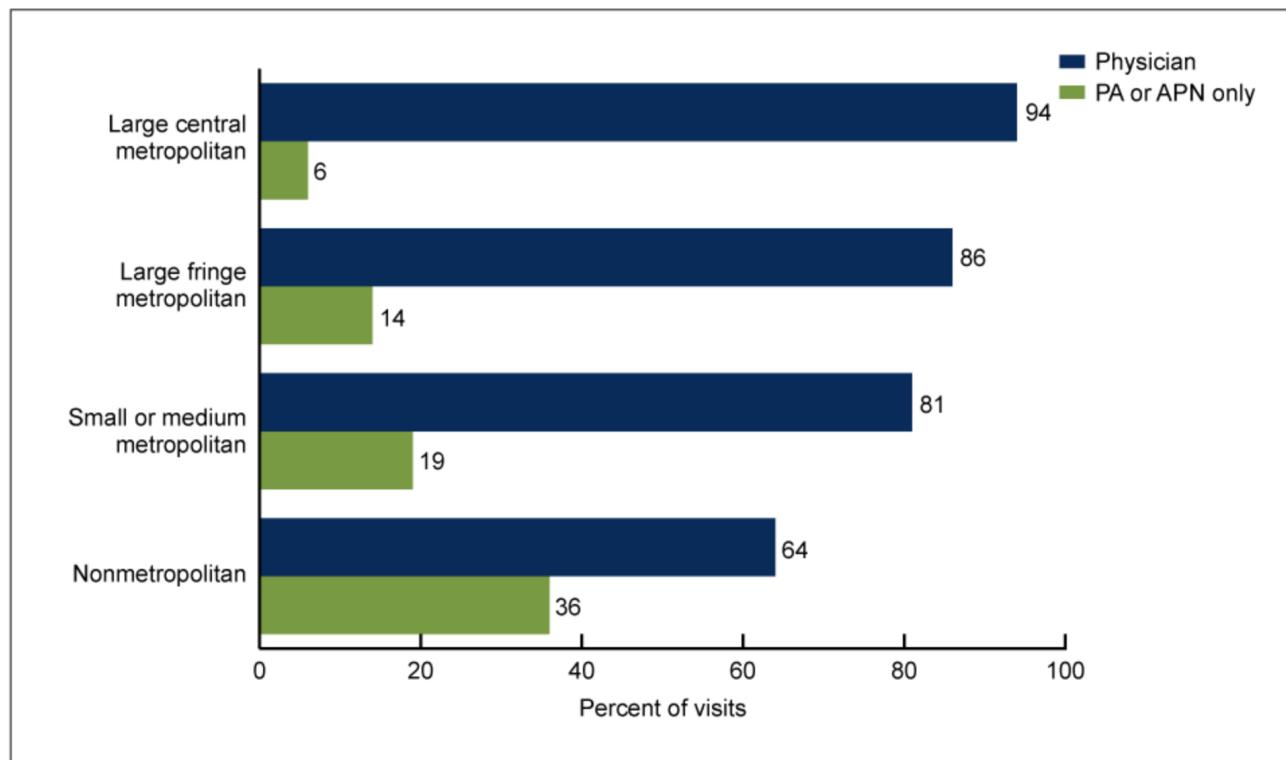


- Increased
- Decreased
- Unable to determine

## Example 1.5.2: Case study: CDC data on the utilization of physician assistants and advance practice nurses.

A physician assistant (PA) is a state-licensed health professional who practices medicine under a physician's supervision. An advanced practice nurse (APN) is a registered nurse with advanced training. A study published by the U.S. Centers for Disease Control on physician assistant and advance practice nurse care in hospital outpatient departments finds that the supply of PAs and APNs are expanding and are playing increasingly diverse roles in the healthcare system.

The chart below appears in a section titled "Does PA/APN utilization differ by hospital location?" The categories being grouped are types of hospital locations: large central metropolitan, large fringe metropolitan, small or medium metropolitan, and nonmetropolitan. The groups are physician and PA/APN.



Source: CDC<sup>5</sup>

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

PARTICIPATION ACTIVITY

1.5.7: Utilization of physician assistants and advance practice nurses.



Consider the bar chart above on the utilization of physician assistants and advanced practice nurses.



1) The bar chart is a horizontal bar chart.

- True
- False



2) The groups are PAs and APNs.

- True
- False

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



3) The utilization of PAs and APNs decrease as the hospital's location becomes less metropolitan.

- True
- False



4) Each category shown is distinct.

- True
- False

## Python-Function 1.5.2: Grouped bar charts.

By setting the **hue** parameter equal to a group, a bar chart is created that counts the frequency of a category by group. Ex: The code that counts the number of people who survived and didn't survive the sinking of the Titanic by class is given below.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

```
# Loads the necessary modules
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Loads the titanic dataset
titanic = pd.read_csv("https://static-
resources.zybooks.com/static/titanic.csv")

# sets the style of the bar charts
sns.set(style="whitegrid", color_codes=True)

# title
plt.title('Titanic survivors and deaths by class', fontsize=20)

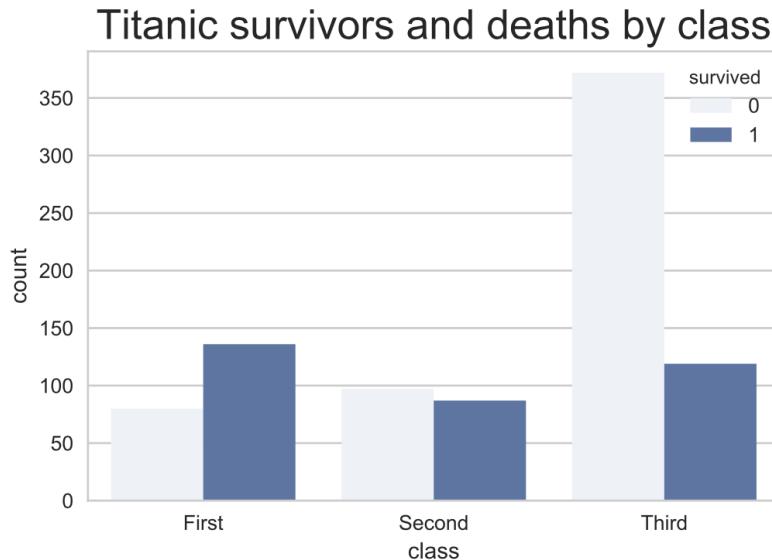
# generates a vertical bar chart
sns.countplot(x="class", hue="survived", color="b", data=titanic);

# saves the image
plt.savefig("groupedbarchart.png")

# shows the image
plt.show()
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The resulting grouped bar chart is shown below.



[Run example](#)

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## Stacked bar chart

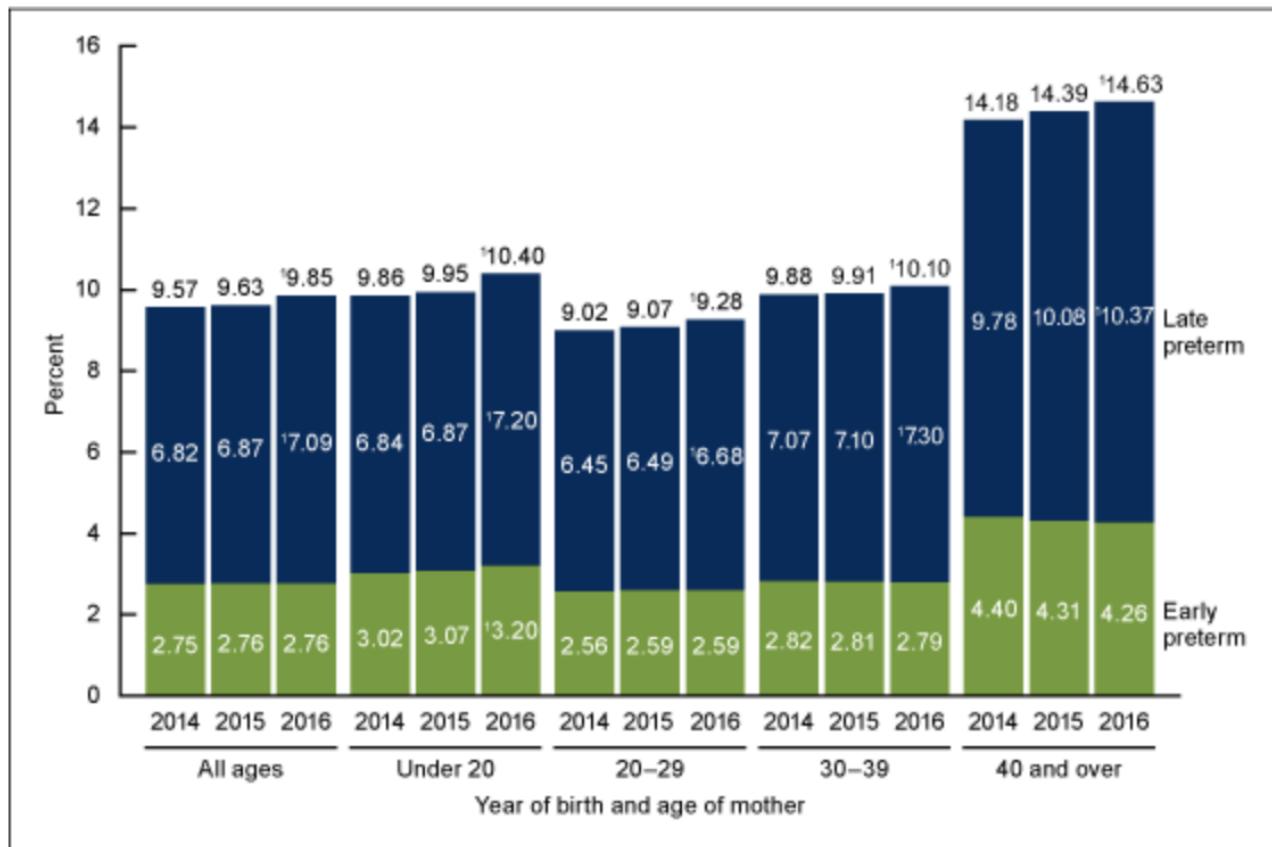
A **stacked bar chart** is a grouped bar chart where the bars are stacked on each other. A stacked bar chart is useful for showing each category's *total*, while still showing the breakdown of groups within each category. However, the relative sizes of each group in a category becomes harder to see due to

not being side-by-side. The following shows a stacked bar chart and grouped bar chart for the same data.

### Example 1.5.3: Case study: CDC data on births in the United States.

A study published by the U.S. Centers for Disease Control on the 2016 birth data presents several health indicators such as fertility rates, birth rates, and multiple birth rates.

The chart below appears in a section titled "The preterm birth rate rose for the second straight year in 2016." Preterm birth refers to a pregnancy that is shorter than normal, especially births that occur after no more than 37 weeks of pregnancy. The chart displays the percentage of preterm births for years 2014, 2015, and 2016 categorized by the mother's age group. The age groups are: under 20, 20-29, 30-39, and 40 and over.



Source: CDC<sup>6</sup>

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

#### PARTICIPATION ACTIVITY

1.5.8: Preterm births in the United States.



Refer to the bar chart above on preterm births in the United States.



1) What is the percentage of preterm births in 2016 for mothers *under 20* years old?

- 3.20%
- 7.20%
- 10.40%

2) What is the percent increase in preterm births for *all* ages between 2015 and 2016?

- 0.22%
- 9.63%
- 9.85%

3) Which age group showed the largest percentage of *late* preterm births in 2016?

- Under 20
- 20 – 29
- 30 – 39
- Over 40

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



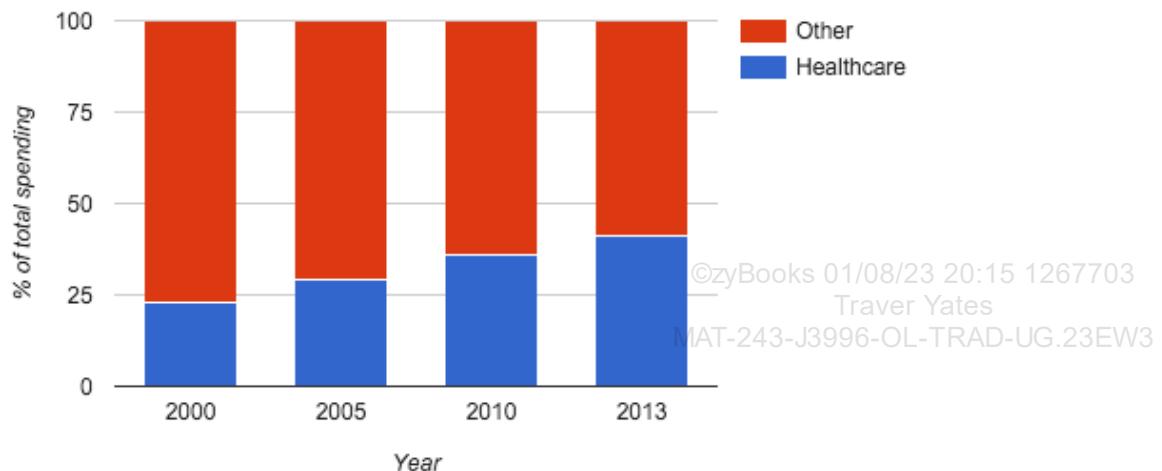
The concept of a relative frequency chart is commonly applied to a stacked bar chart. Such a chart clearly shows how a particular group's proportion of the total changes across categories (such as across years).

Figure 1.5.3: Massachusetts state spending on healthcare versus all other state spending, using a relative frequency stacked bar chart.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



Source: Kaiser Family Foundation

**PARTICIPATION ACTIVITY**

1.5.9: Relative frequency stacked bar chart.



Refer to the above figure showing Massachusetts state spending.

- 1) In 2000, about what percent of state spending was on healthcare?

- 25%
- 40%
- 75%



- 2) Did the relative percentage of healthcare spending to total spending increase or decrease from 2000 to 2013?

- Increase
- Decrease
- Cannot determine



- 3) Based solely on viewing the bar chart, what is the best prediction for relative percentage of healthcare spending to total spending in 2020?

- 25%
- 50%

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



## Python-Function 1.5.3: Stacked bar charts.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

To create a stacked bar chart, the `plt.subplots()` of `matplotlib` is used. Multiple plots share the same set of axes. Ex: The code below plots the number of automobile collisions related to speeding and the total number of collisions in the same set of axes.

```
# loads the necessary modules
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# initialize figure
f, ax = plt.subplots()

# load data frame
crashes = sns.load_dataset("car_crashes")
df = crashes.loc[range(5)]

# plot total crashes
sns.set_color_codes("pastel")
sns.barplot(x="total", y="abbrev", data=df,
            label="Total", color="b")

# plot crashes related to speeding
sns.set_color_codes("muted")
sns.barplot(x="speeding", y="abbrev", data=df,
            label="Speeding-related", color="b")

# title
plt.title('Speeding-related automobile collisions', fontsize=20)

# legend
ax.legend(ncol=1, loc="lower right")
ax.set(xlim=(0, 28), ylabel="State", xlabel="Automobile collisions (per billion miles)");

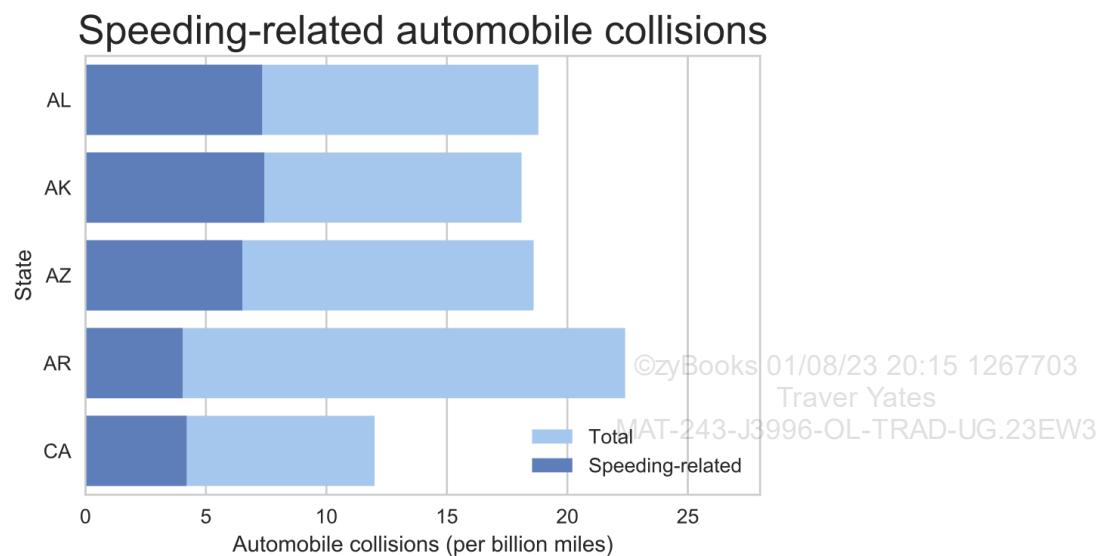
# saves the image
plt.savefig("stacked.png")
```

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

The resulting stacked bar chart is shown below.



[Run example](#)

## References

- (\*)1) "Fortune 500 Lists". [www.fortune.com/fortune500/list](http://www.fortune.com/fortune500/list). Last accessed 24 Sept 2018.
- (\*)2) "The Disadvantages of Summer Babies ". Freakonomics, 2 Nov. 2011, [freakonomics.com/2011/11/02/the-disadvantages-of-summer-babies/](http://freakonomics.com/2011/11/02/the-disadvantages-of-summer-babies/).
- (\*)3) Konnikova, Maria. "Youngest Kid, Smartest Kid". The New Yorker, 19 Nov. 2013, [www.newyorker.com/tech/elements/youngest-kid-smartest-kid](http://www.newyorker.com/tech/elements/youngest-kid-smartest-kid). ©zyBooks 01/08/23 20:15 1267703 Traver Yates
- (\*)4) "Civilian Labor Force by Sex". U.S Dept. of Labor, [www.dol.gov/wb/stats/NEWSTATS/facts/civilian\\_lf\\_sex\\_2016\\_txt.htm](http://www.dol.gov/wb/stats/NEWSTATS/facts/civilian_lf_sex_2016_txt.htm). MAT-243-J3996-OL-TRAD-UG.23EW3
- (\*)5) Esther Hing, M.P.H. and Sayeedha Uddin, M.D., M.P.H. "NCHS Data Brief No. 77: Physician Assistant and Advance Practice Nurse Care in Hospital Outpatient Departments 2008-2009." U.S. Centers for Disease Control.
- (\*)6) Joyce A. Martin, M.P.H, Brady E. Hamilton, Ph.D., and Michelle J.K. Osterman, M.H.S. "NCHS Data Brief No. 287: Births in the United States. 2016". 2011.

## 1.6 Pie charts



This section has been set as optional by your instructor.

### Pie charts

A **pie chart** shows relative frequency for categories using a circle, with each category shown as a slice of appropriate size. The appearance is one of a sliced pie, leading to the chart's name. Because length differences are interpreted more precisely than size differences, bar charts are often preferred. However, pie charts remain common, perhaps in part because curved shapes are more aesthetically pleasing than rectangular shapes.

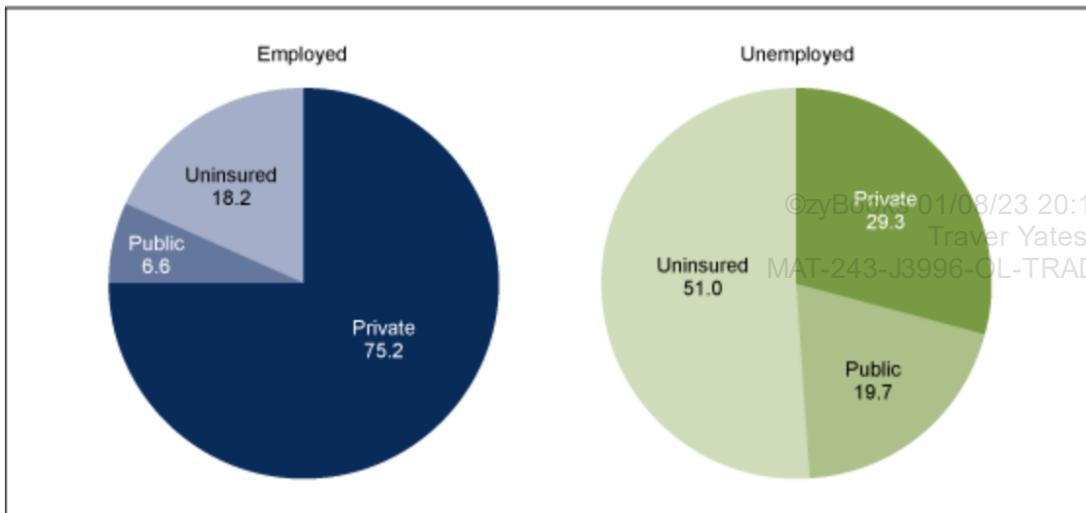
#### Example 1.6.1: Case study: CDC data on health and access to care among employed and unemployed adults.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

A study published by the U.S. Centers for Disease Control identifies the health insurance status and type of insurance among adults in the U.S. between 18-64 years old according to employment in 2009-2010. Key findings showed that unemployed individuals are less likely to have private insurance, less likely to receive needed prescription medications, and less likely to receive needed medical care than those who are employed. The pie chart below

shows the insurance status and type of insurance an individual has according to employment status.

**PARTICIPATION ACTIVITY**

1.6.1: Health and access to care among employed and unemployed adults.



- 1) 6.6% of unemployed individuals had public insurance.



- True
- False

- 2) The percentage of uninsured among unemployed individuals was nearly three times as high as those who are employed.



- True
- False

- 3) The percentage of uninsured individuals in the U.S. is 69.2%.



- True
- False

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Python-Function 1.6.1: Pie charts.

A pie chart is creating by using the `pie(sizes, labels, autopct='%.1f%%')` function in `matplotlib`. `sizes` is a list parameter that consists of frequencies of each

group or category. **labels** is also a list parameter than consists of names of each group and category. **autopct** is an optional parameter that indicates the number of decimal places of the percent labels. Ex: **%1.1f%** means that the labels in the pie chart will have one decimal place with the percent symbol.

```
# loads the necessary libraries
import matplotlib.pyplot as plt
import seaborn as sns

# loads the titanic dataset
df = sns.load_dataset("titanic")

# counts the number of passengers for each class
a = df[df.pclass == 1]["pclass"].count()
b = df[df.pclass == 2]["pclass"].count()
c = df[df.pclass == 3]["pclass"].count()

# data to plot
labels = 'First', 'Second', 'Third'
sizes = [a, b, c]

# plot
plt.pie(sizes, labels=labels, autopct='%1.1f%%')

# title
plt.title('Titanic passengers by class', fontsize=20)

# legend
patches, texts = plt.pie(sizes)
plt.legend(patches, labels, loc="lower right")

# produces a perfectly circular chart
plt.axis('equal')

# saves the image
plt.savefig("piechart.png")

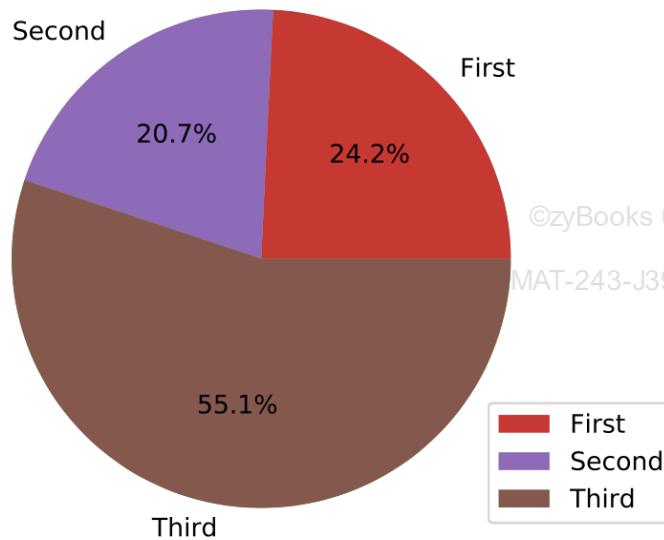
# shows the image
plt.show()
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The resulting pie chart is shown below.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## Titanic passengers by class



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

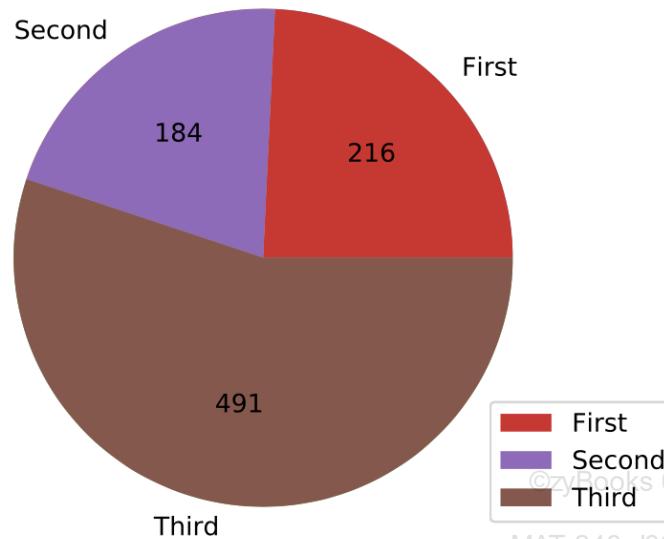
[Run example](#)

The counts, instead of the percentages, can be displayed by setting `autopct` equal to some function as shown in the code below

```
tot = sum(sizes)
autopct = lambda x: "%d" % round(x*tot/100,2)
plt.pie(sizes, labels=labels, autopct= autopct);
```

The result pie chart with counts is shown below.

## Titanic passengers by class



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

[Run example](#)

## Exploded pie charts

An **exploded pie chart** is a pie chart with one or more slices separated from the rest of the pie. By exploding a pie chart, the exploded pieces are emphasized. In practice, the most important category - usually the category with the largest count or percentage - is often exploded. For pie charts with many categories, exploding all slices may enhance readability.

### Example 1.6.2: U.S. Bureau of Labor Statistics time use survey.

©zyBooks 01/08/23 20:15 1267703

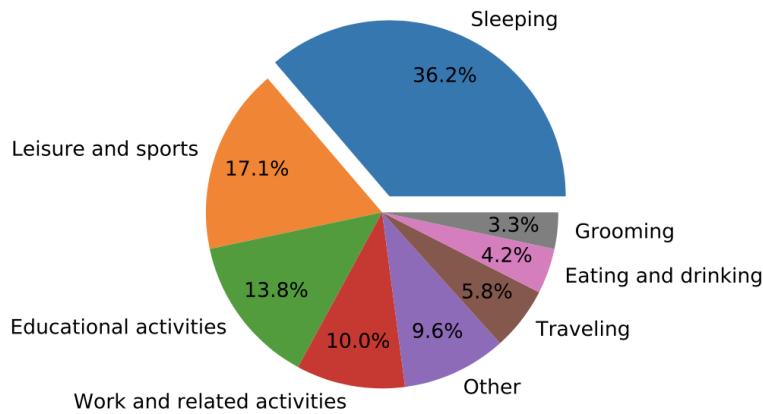
Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Each year, the U.S. Bureau of Labor Statistics (BLS) conducts a survey called the American Time Use Survey (ATUS) that measures the amount of time people spend doing various activities. The BLS website provides interactive charts that displays data by group such as sex, employment status, and age.

The exploded pie chart below illustrates the time use average percentages per weekday for full-time college and university students with an emphasis on time spent sleeping.

**Average weekday time use for full-time college and university students in the U.S.**



Credit: BLS.gov <sup>1</sup>

**PARTICIPATION ACTIVITY**

1.6.2: Exploded pie charts.



Refer to the exploded pie chart above.

- 1) College students spend more of the non-sleeping time on \_\_\_\_ than any other activity.

- work and related activities
- educational activities
- leisure and sports

- 2) The full pie represents \_\_\_\_ hours.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- 24
- 100

## Python-Function 1.6.2: Exploded pie charts.

©zyBooks 01/08/23 20:15 1267703

To explode one or more categories within a pie chart, the parameter `explode` of the `pie` is set to a tuple. Ex: In the tuple `(0.1, 0.0)`, the category "Survived" is exploded by 0.1 units, while the category "Did not survive" is not exploded.

```
# data to plot
a = df[df.survived == 1]["survived"].count()
b = df[df.survived == 0]["survived"].count()

# creates a tuple of the categories
explodeTuple = (0.1, 0.0)

# data to plot
labels = 'Survived', 'Did not survive'
sizes = [a, b]

# plot
plt.pie(sizes, labels=labels, autopct='%1.1f%%', explode = explodeTuple)

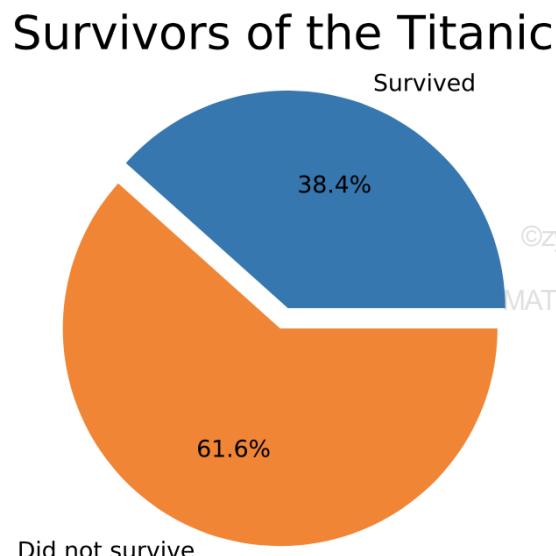
# title
plt.title('Survivors of the Titanic', fontsize=20)

# produces a perfectly circular chart
plt.axis('equal')

# saves the image
plt.savefig("explodedpiechart.png")

# shows the image
plt.show()
```

The resulting exploded pie chart is shown below.



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

[Run example](#)

## Misuse of pie charts

While pie charts can be a useful tool, many situations exist in which pie charts can be difficult to read and interpret. Many optional pie chart effects can also make the pie chart unclear or distort the presented data.

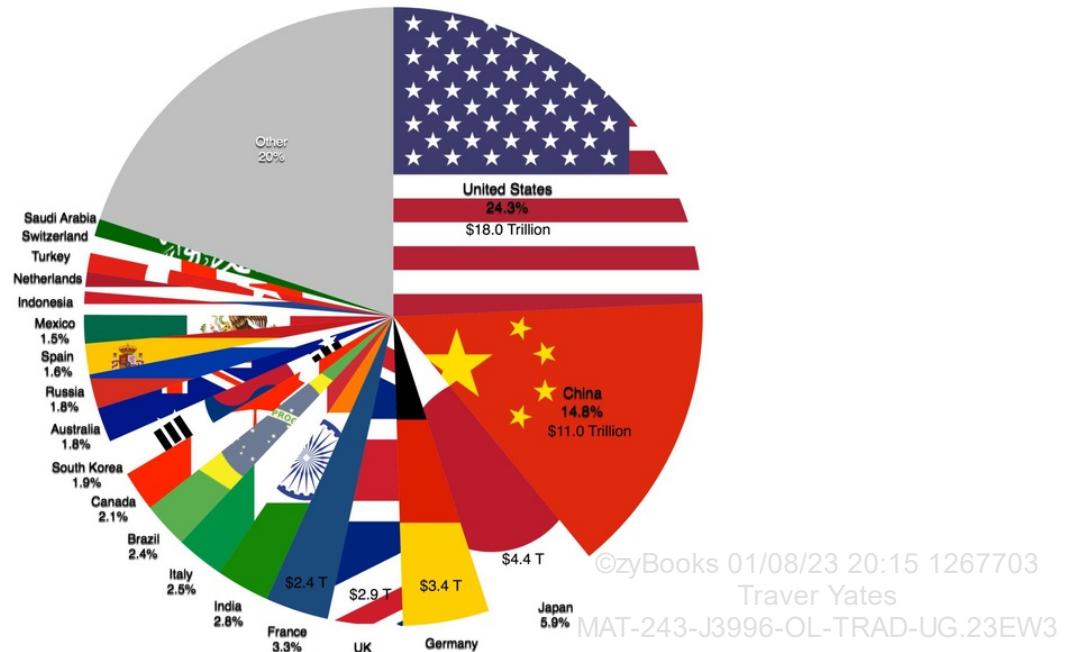
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

### Pie charts with too many slices

A pie chart should not have more than five or six slices. With more slices, the relative percentage of each slice becomes difficult to visually determine. While the percentages can be specified with labels, as in the chart below, too many labels also clutter the chart and make the chart difficult to interpret.

#### Example 1.6.3: Too many slices.

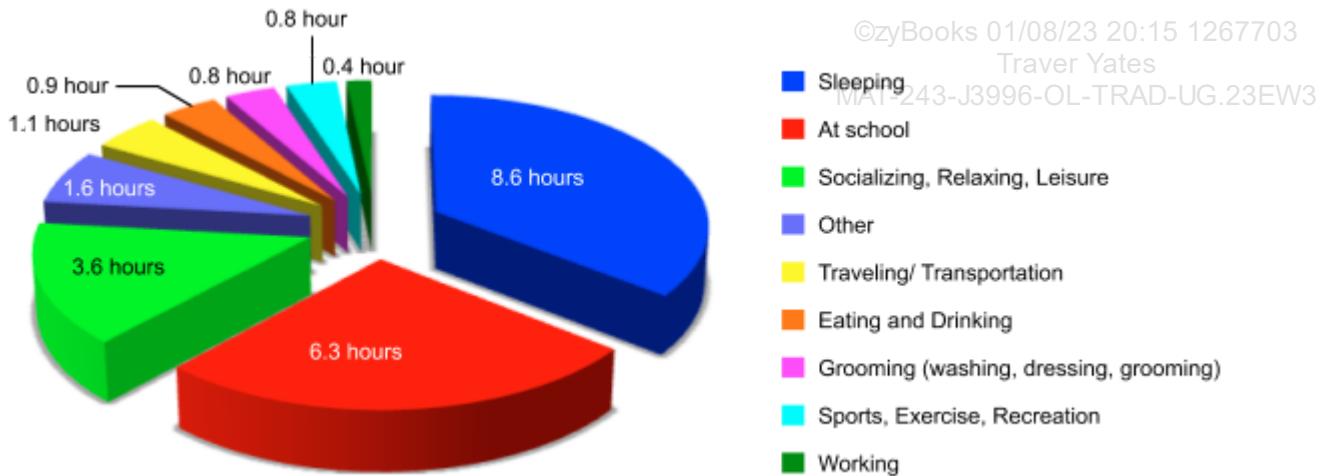
The pie chart below shows the twenty largest economies in the world. While interpreting the largest slices is straightforward, the relative sizes of the slices smaller than Germany and the UK are difficult to compare, and the slices representing Indonesia, Netherlands, Turkey, Switzerland, and Saudi Arabia do not convey much information at all.



Credit: User Wikideas<sup>1</sup> <sup>2</sup> <sup>3</sup>

#### Example 1.6.4: Pie charts with a legend.

To avoid the confusion that can occur with many labels, pie charts are sometimes presented with a legend. Below is a pie chart that shows U.S. high-school student time use. The use of a legend to represent the categories interferes with the intuition provided by graphics, instead requiring a reader to look back and forth to determine which slice corresponds to which activity.



Credit: BLS.gov<sup>4</sup>

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Sleeping

At school

Socializing, Relaxing, Leisure

Other

Traveling/ Transportation

Eating and Drinking

Grooming (washing, dressing, grooming)

Sports, Exercise, Recreation

Working

## Pie charts with poor color choices

Pie charts can also become difficult to read if the colors of the slices are chosen badly. The Largest Economies pie chart above uses the flags of individual nations as the colors of the slices. While the intention may have been to convey more clearly which nation is being represented, the effect is to distract the reader's eye from the sizes of the slices. Colors that are too similar are also bad choices, as in the Time Spent pie chart above, where the two variations of blue used for the Sleeping and the Other categories are difficult to distinguish.

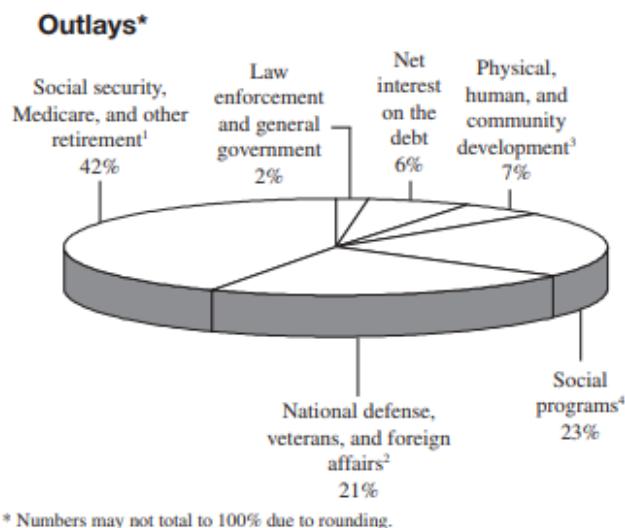
## 3D pie charts

Even a basic, two-dimensional pie chart can be rendered useless if too many pieces, poor labeling, and confusing color choices are used. 3D effects, exploding charts, and unusual shapes often tend to obscure the data being presented even further.

Python allows for the creation three dimensional pie charts. However, 3D pie charts distort the shape of the charts, making estimating the size of the slices difficult, at best, and allowing data to be misinterpreted, at worst.

### Example 1.6.5: 3D pie charts.

The chart below from the Internal Revenue Service shows the relative sizes of the major categories of spending for the federal government in 2016. The distortion of the slices from the 3D effect and the orientation of the chart makes the relative sizes of the slices almost impossible to estimate. Ex: "Law enforcement and general government" and "Physical, human, and community development" appear to be nearly the same size, but the former only makes up 2% of the spending, while the latter makes up 7%. The third slice at the back of the graph, "Net interest on the debt", appears larger than either of the other two slices but only represents 6% of federal spending. Similarly, the "National defense, veterans, and foreign affairs" category at the front of the graph appears considerably larger than the "Social programs" category, even though spending in the former category is 2% less.



Credit: IRS.gov <sup>5</sup>

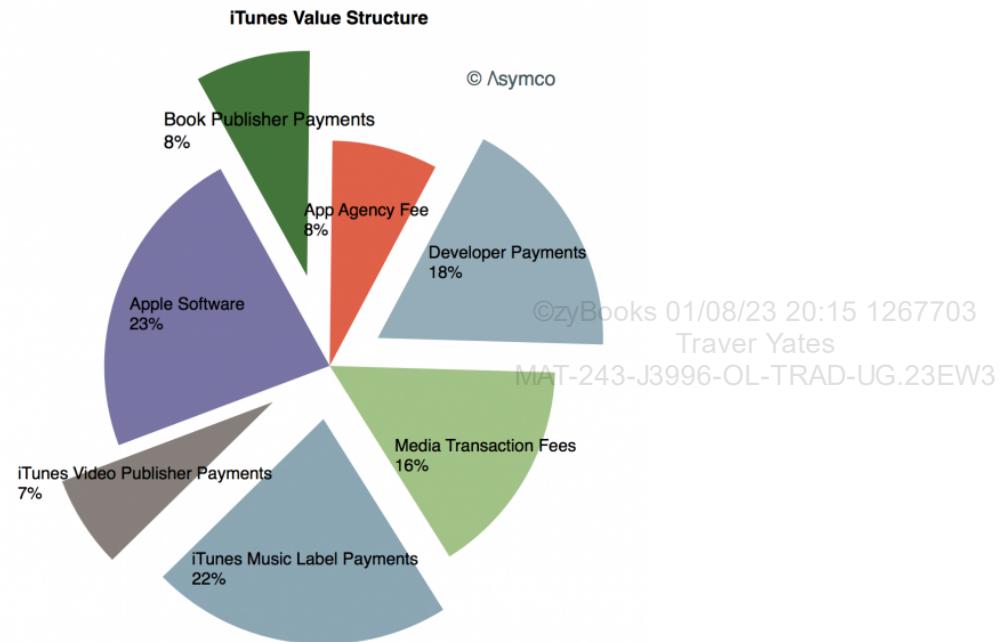
## Exploded pie charts

Exploded pie charts can also be used badly. While using this effect on one or two slices does emphasize the categories being represented, the process also necessarily changes the relative sizes of the slices, and can make the exploded slices appear larger than appropriate. Using the effect on multiple slices increases the difficulty in visually comparing the relative sizes of the slices, and can make the chart disorienting, as in the pie chart below.

### Example 1.6.6: Exploded pie chart.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The exploded pie chart below shows the sources of iTunes revenue. Four of the seven slices have been exploded. Rather than emphasizing those four slices, the unevenness of the white dividing lines and the varying distances that the slices have been moved make interpreting the chart much more difficult.



Credit: Horace Dediu<sup>6</sup>

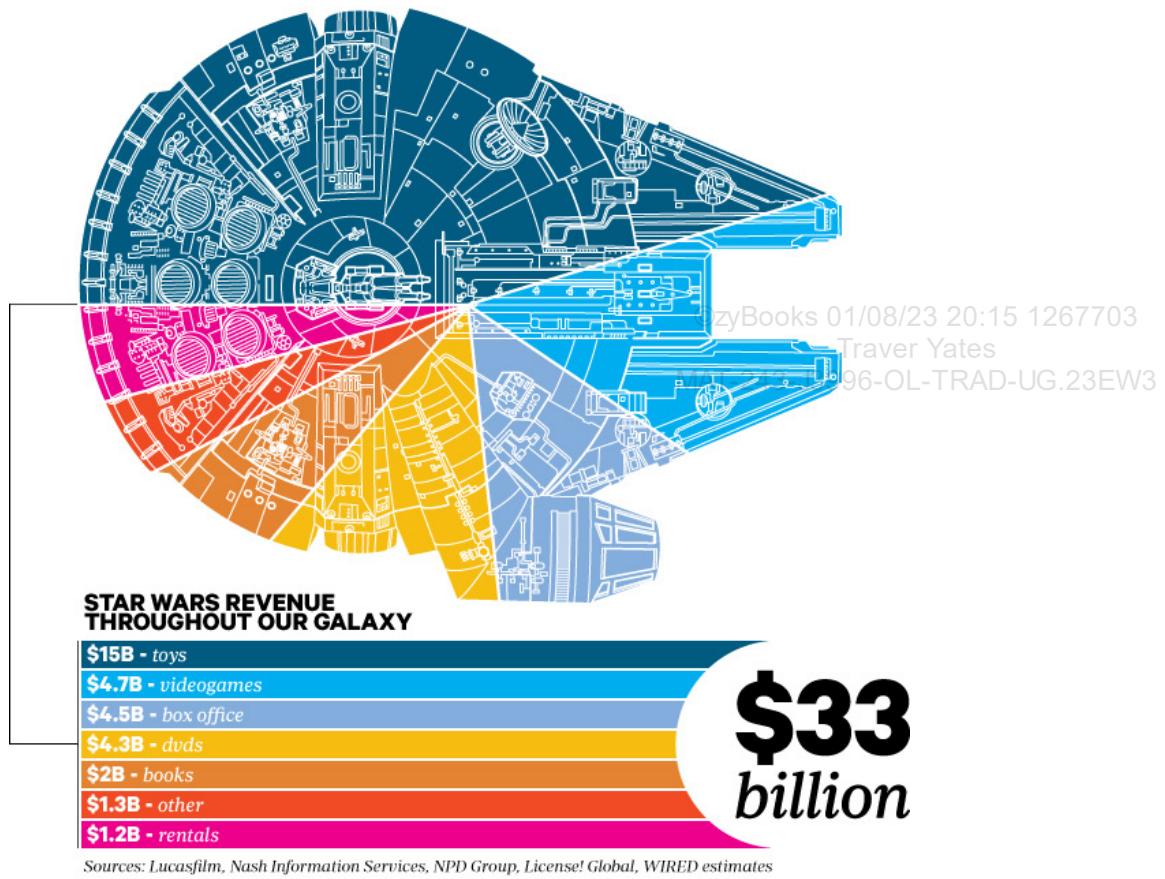
## Non-circular pie charts

Finally, in an attempt to be creative or interesting, non-circular pie charts are sometimes used. Unfortunately, visually comparing the sizes of slices in a non-circular pie chart is often nearly impossible.

### Example 1.6.7

The pie chart of the sources of revenue from the Star Wars movies is shaped like the ship the Millennium Falcon. While the shape is clever, the slices of the "pie" do not have a uniform length, and thus the relative sizes of the slices are even more difficult to interpret. Ex: the slice representing "rentals" is wider than the slice representing "other", even though the "other" slice represents \$100 million more revenue than the Rental slice.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

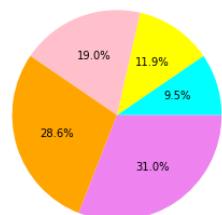
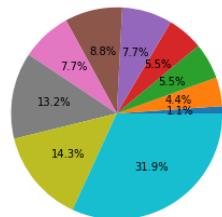


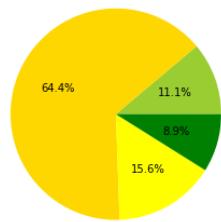
Credit: Michael Cerwonka<sup>7</sup>

### PARTICIPATION ACTIVITY

#### 1.6.3: Identifying misused pie charts.

- 1) Which of the following pie charts is formatted most appropriately?





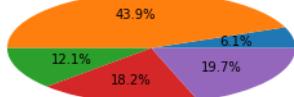
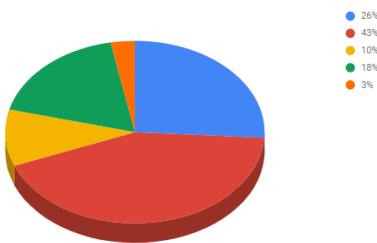
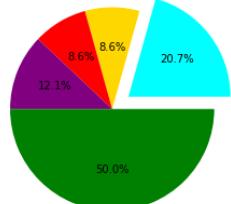
©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- 2) Which of the following pie charts is formatted most appropriately?



## References

(\*)1) "American Time Use Survey." Bureau of Labor Statistics. 20 December 2016, <https://www.bls.gov/tus/charts/students.htm>.

(\*)2) User Wikideas1. "20 Largest Economies Pie Chart." Wikimedia Commons, the free media repository. 25 February 2017, en.wikipedia.org/wiki/File:20\_Largest\_economies\_pie\_chart.pdf. MAT-243-J3996-OL-TRAD-UG.23EW3

(\*)3) "GDP Ranking." The World Bank Data Catalog. 29 June 2018, datacatalog.worldbank.org/dataset/gdp-ranking.

(\*)4) "Time Use Survey activity, How Do You Spend Your Time?." Bureau of Labor Statistics. Retrieved 17 July 2018, www.bls.gov/k12/content/teachers/pdf/atus\_activity1\_intro.pdf.

(\*)5) "Major Categories of Federal Income and Outlays for Fiscal Year 2016." Internal Revenue Service 1040A Instructions 2017. p. 89 Retrieved 17 July 2018, www.irs.gov/pub/irs-pdf/f1040a.pdf.

(\*) Dediu, Horace. "Measuring the iBook Market." Asymco. 28 February 2013, www.asymco.com/2013/02/28/measuring-the-ibook-market/.

(\*) Cerwonka, Michael. "Tell Jabba I've Got His Money: Star Wars Revenue Throughout Our Galaxy." Wired. 25 May 2012, www.wired.com/2012/05/tell-jabba-ive-got-his-money-star-wars-revenue-throughout-our-galaxy/.

## 1.7 Scatter plots

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Scatter plots with quantitative variables

A **scatter plot** depicts the relationship between two variables on a rectangular coordinate system, where each axis corresponds to one variable. Scatter plots are used for both quantitative and categorical data.

In data analytics, scatter plots are especially useful in visualizing the relationship between variables in a multi-dimensional dataset. Ex: A marketing manager for a beverage company may want to check the relationship between average temperature and revenue within a specific time period or between revenue and marketing budget.

Example 1.7.1: Number of engineering faculty versus school rank.

To inform a decision of whether to hire new engineering faculty, a dean collected data showing number of faculty versus engineering school rank (using the U.S. News and World Report ranking, lower rank is better), for eight University of California campuses (UC Berkeley, UCLA, etc.). A table of the data (2014) is shown below. For example, the school with **247** engineering faculty is ranked number **3** in the country, while the school with only **78** engineering faculty is ranked number **81**.

Engineering faculty	USNWR rank
247	3
194	14
155	16
124	19
198	31
115	38

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

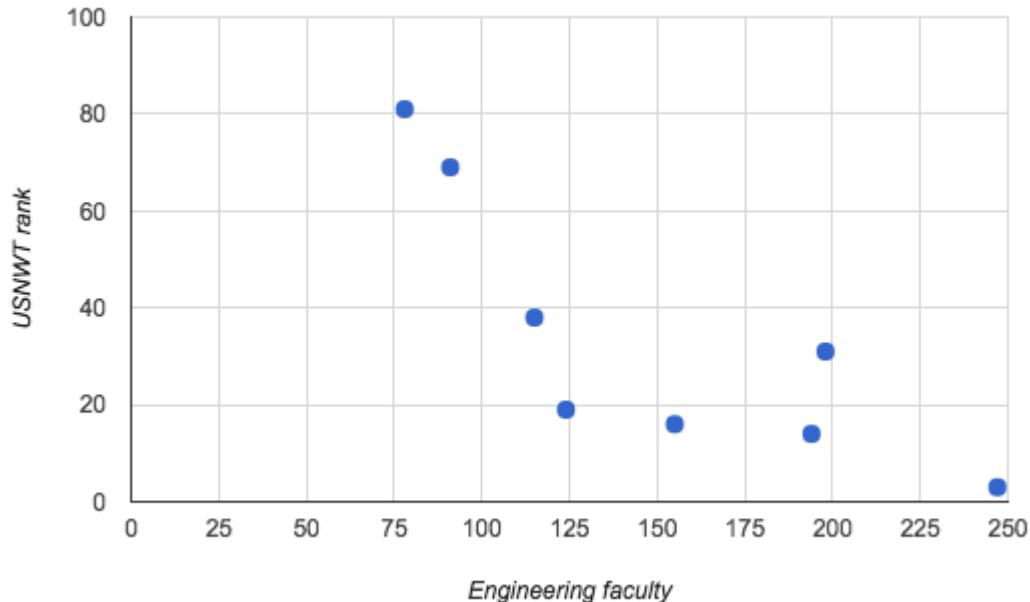
91	69
78	81

Below is a scatter plot showing engineering faculty size vs. engineering school rank in 2014 for the eight campuses in the University of California system. Each row in the above table becomes a coordinate, leading to the following scatter plot.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



The scatter plot clearly shows the relationship between number of faculty and rank. The data suggests, but does not prove, that increasing the number of engineering faculty may be important to improving rank.

A scatter plot often has numerous data points that are "scattered" about the rectangular coordinate system, leading to the name "scatter plot". Below is a scatter plot showing all 128 college football team rankings (lower is better; number 1 is best) and the total salary for each team's head coach<sup>1 2</sup>. (Yes, college head coaches often have multi-million-dollar salaries). A viewer quickly sees that more than half of coaches earn over

(1

million, that many earn

)3

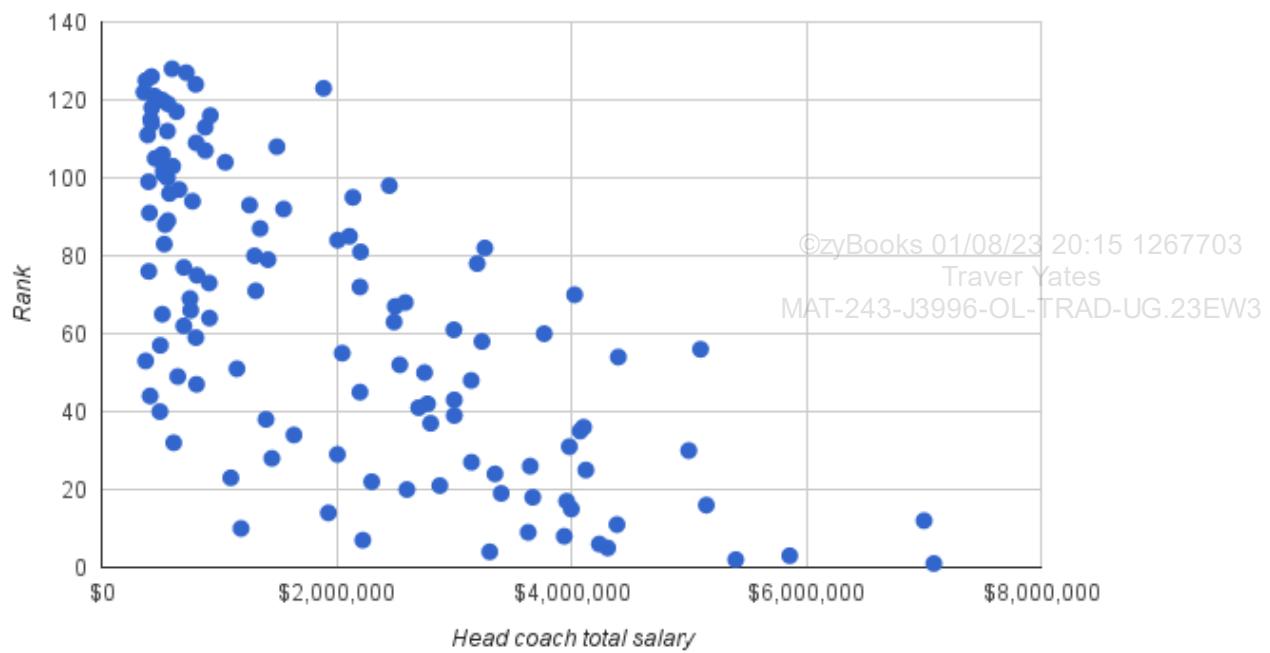
or \$4 million, and that most poorly-ranked team's coaches earn just a few hundred thousand dollars. The viewer can also see that several teams, in the lower left, seem to be getting a great bargain.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Figure 1.7.1: College football rankings vs. head coach salary.

**PARTICIPATION ACTIVITY****1.7.1: Scatter plot.**

Consider the scatter plot above showing number of engineering faculty versus rank for eight UC schools.

- 1) What is the approximate rank for the school having about **150** engineering faculty?



- About 5
- About 20
- About 80

- 2) If a UC school had about **100** engineering faculty, what rank might that school expect to achieve?



- 80
- 55
- 20

- 3) At about **200** on the  $x$ -axis are two schools, one ranks at about **15** and another at about **30**. Which of the following is the best inference?



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- Two schools having about 200 faculty have different ranks.
- One school has two different ranks.
- The data contains an error.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

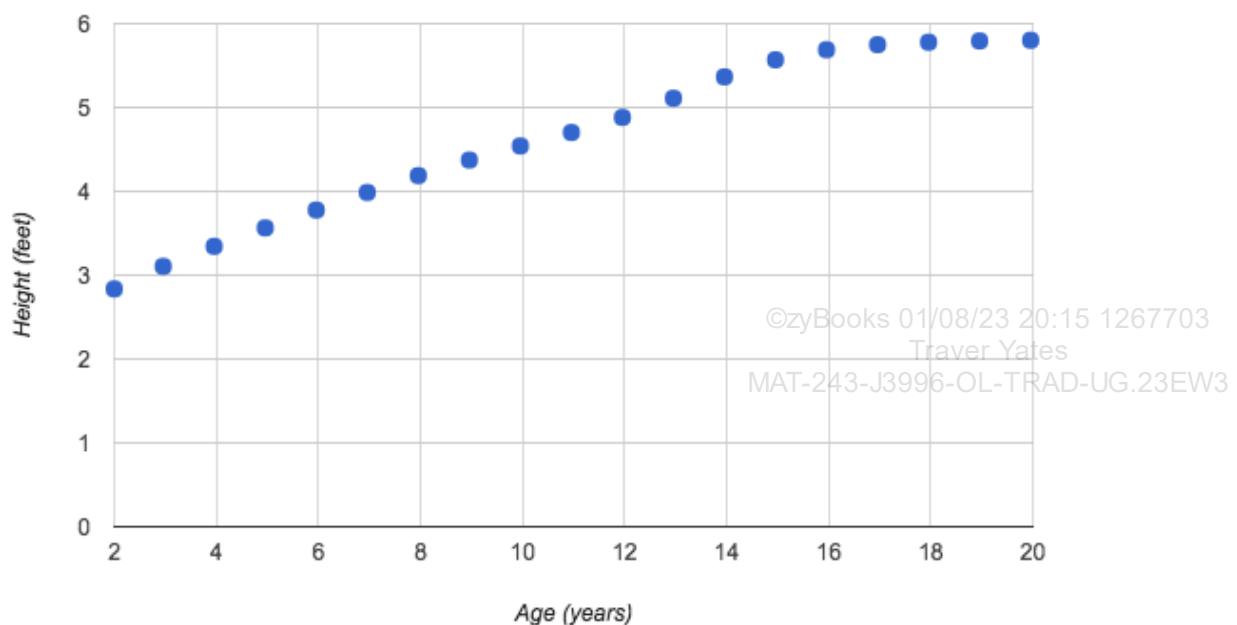
## Independent and dependent variables

Commonly, a scatter plot shows how one variable depends on another. Above, the dean was interested in showing how rank might depend on the number of faculty. To distinguish dependent and independent variables, a question might be: If one variable's value is 5 (or some other value), how is the other variable affected? The variable that is controlled by an observer or is a reason for variation is the **independent variable**, while the variable that is then determined based on that variable is the **dependent variable**. The independent variable is usually plotted on the  $x$ -axis and the dependent variable on the  $y$ -axis. Typically, a desired independent variable is first found along the  $x$ -axis, and the corresponding dependent variable is found along the  $y$ -axis.

Below is an example showing how height depends on age for males aged 2 to 20 years. Ex: a question might be "If a male is aged 12 years, what height might he be?". Age is the independent variable and is thus plotted on the  $x$ -axis.

### Example 1.7.2: Height vs. age for males aged 2 to 20.

The following scatter plot shows the median height for U.S. males aged 2 to 20 years.



**PARTICIPATION ACTIVITY**

## 1.7.2: Independent and dependent variables.



1) A dean wishes to know whether increasing faculty may lead to better rank. The number of faculty is the \_\_\_\_ variable.

- dependent
- independent

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

2) Parents and doctors wish to compare a child's height to other children's heights of the same age. The height is the \_\_\_\_ variable.

- dependent
- independent



3) A news station records temperature and humidity for various days and then creates a scatter plot to determine if any relationship exists. Is either temperature or humidity an obvious independent variable?

- Yes
- No



4) For the above scatter plot of male height, at about what age do males stop growing taller?

- 16
- 20



## Regression curve

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

A **regression curve** is a curve added to a scatter plot that shows the relationship between two variables. A linear regression curve is also called a regression line or a trend line. The first figure below shows a trend line for the earlier data showing the relationship between coach salary and team rank, and the second figure shows a trend line for the earlier data showing the relationship between number of engineering faculty and rank.

Figure 1.7.2: Linear trend line added to scatter plot for data showing the relationship between head coach salary and team ranking.

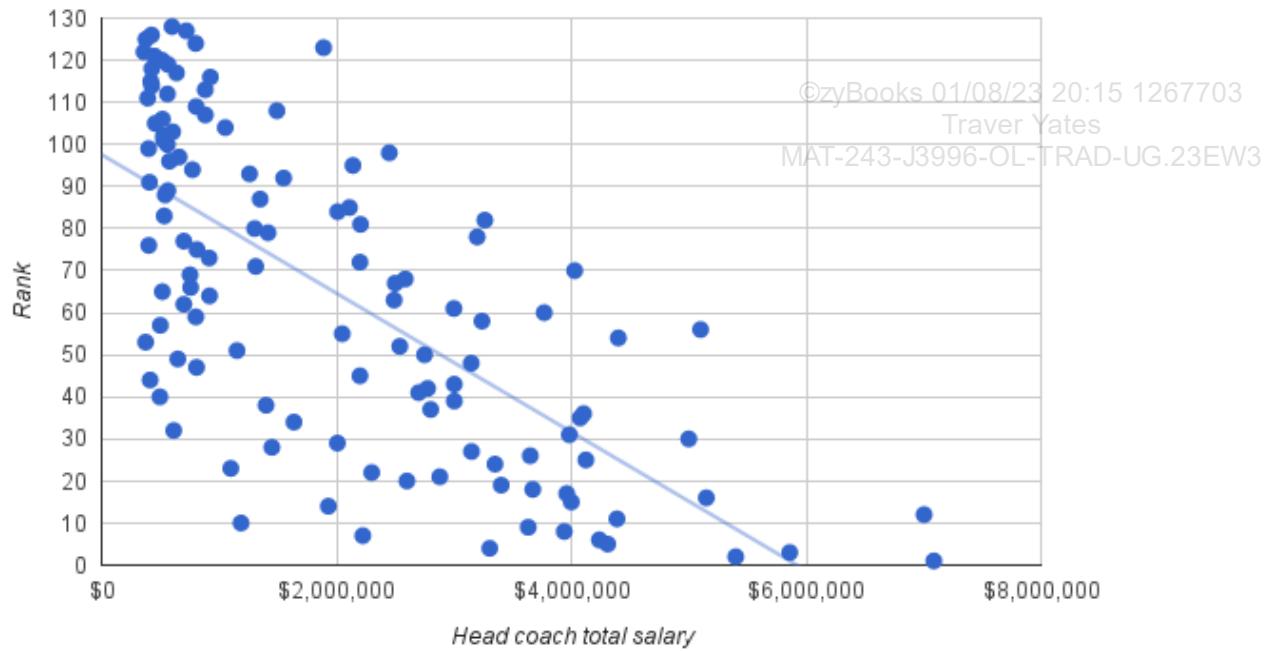
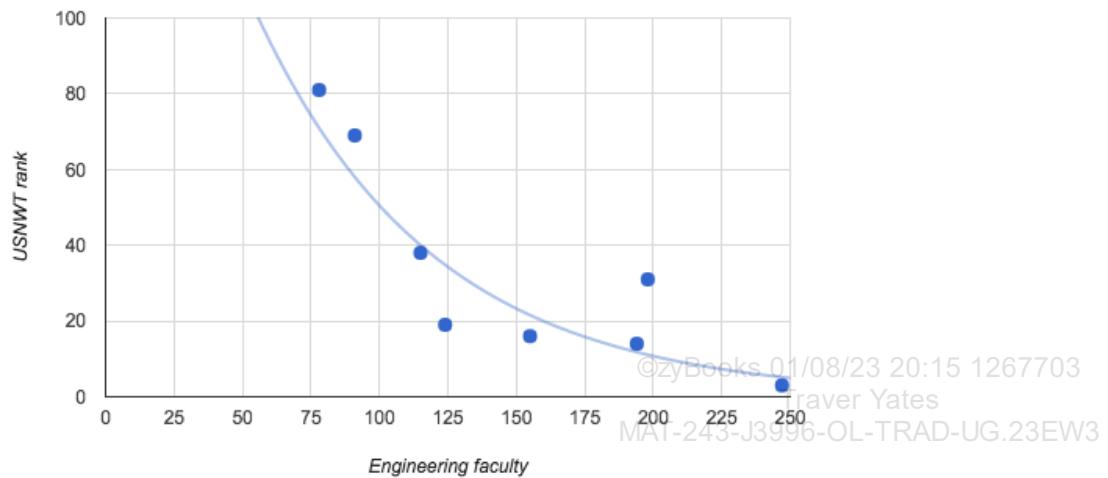


Figure 1.7.3: Exponential regression curve added to scatter plot for data showing the relationship between number of engineering faculty and rank.





- 1) Each point in a scatter plot lies somewhere on the regression curve.

True  
 False

- 2) Given data where the  $x$ -value represents the area to be painted in square feet and the  $y$ -value represents amount of paint in gallons, the regression curve is likely \_\_\_\_.

linear  
 non-linear

- 3) Given data where the  $x$ -value represents a person's height and the  $y$ -value represents a person's age for ages 1-50, the regression curve is likely \_\_\_\_\_.

linear  
 non-linear

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



## Python-Function 1.7.1: Scatter plots.

To graph a scatter plot, the `sns.regplot(x, y, ci = None)` function of the `seaborn` library is used, which takes in two arrays of equal size `x` and `y`. By default, the 95% confidence intervals are displayed, but the parameter `ci` can be set to `None` to disable confidence intervals.

```
# loads the necessary libraries
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

# creates the data points
x = np.array([0, 5, 3, 4, 7, 8, 10])
y = np.array([5, 2, 5, 15, 27, 15, 31])

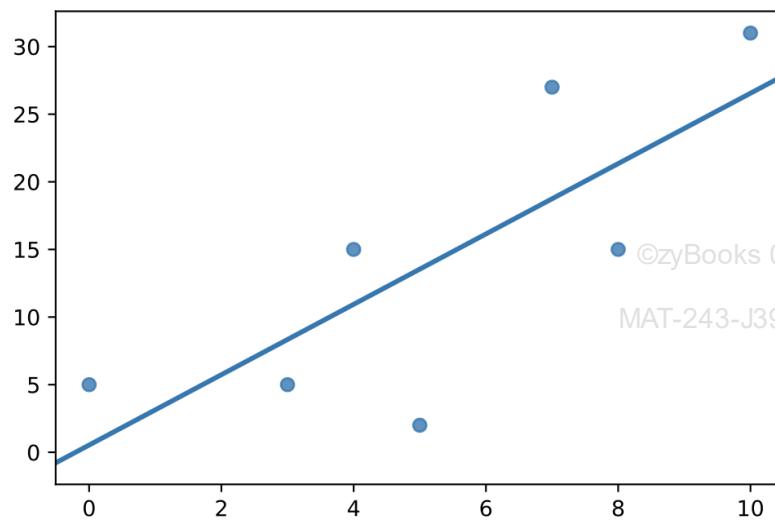
# plot
sns.regplot(x, y, ci=None)

# saves the image
plt.savefig("scatterplot.png")

# shows the image
plt.show()
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The resulting scatter plot is shown below.



[Run example](#)

The `regplot` can also be used with data frames. Ex: Using the `iris` dataset, petal length is plotted with respect to petal width using the code below.

```
# loads the necessary modules
import matplotlib.pyplot as plt
import seaborn as sns

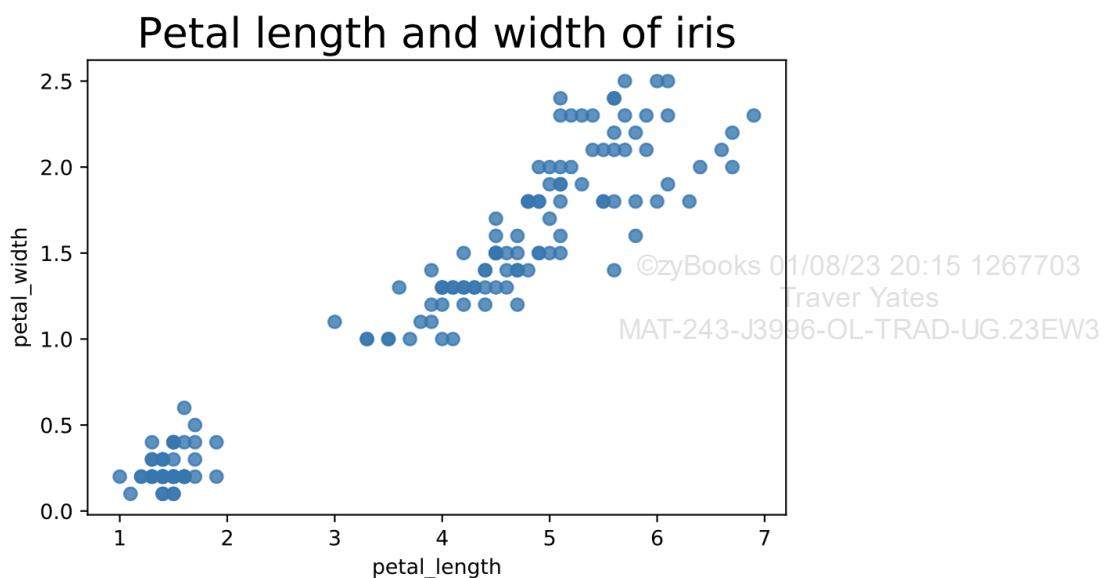
# loads the iris data set
df = sns.load_dataset("iris")

# title
plt.title('Petal length and width of iris', fontsize=20)

# plot
sns.regplot(df["petal_length"], df["petal_width"], ci=None, fit_reg=False);

# saves the image
plt.savefig("irisscatterplot.png")

# shows the image
plt.show()
```



[Run example](#)

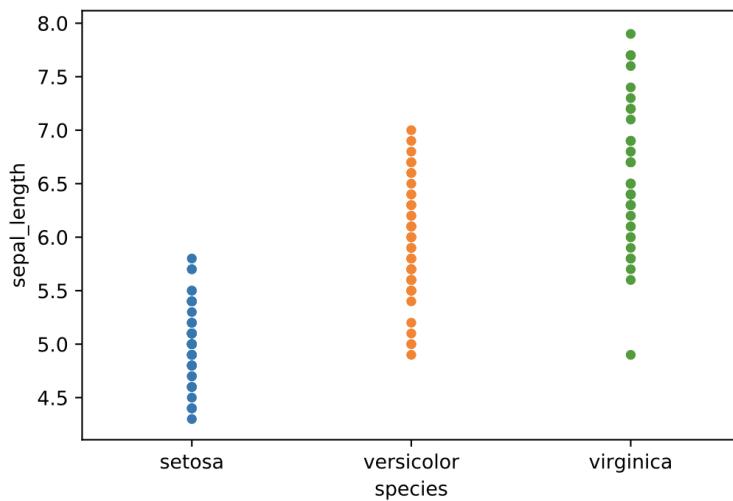
## Scatter plots with categorical variables

Earlier examples showed the use of scatter plots to determine the relationship between two nominal variables with respect to different categories. Ex: The relationship between petal length and petal width of iris is linear. However, when one of the two variables is categorical, other plotting techniques should be used. Three types of plots are commonly used when dealing with categorical variables: (1) strip plots, (2) jittered strip plots, and (3) swarm plots.

### Strip plots

A **strip plot** is a scatter plot where a categorical variable represents an axis and an ordinal variable represents the other. Points are stacked on top of each other and form a single column or strip. A strip plot is useful in summarizing information about the dataset. Ex: The strip plot below shows the relationship between iris species and sepal length. The horizontal axis shows the different species of iris and the vertical axis shows the sepal length in centimeters.

Figure 1.7.4: Strip plot for the sepal lengths of iris species.



### Jittered strip plots

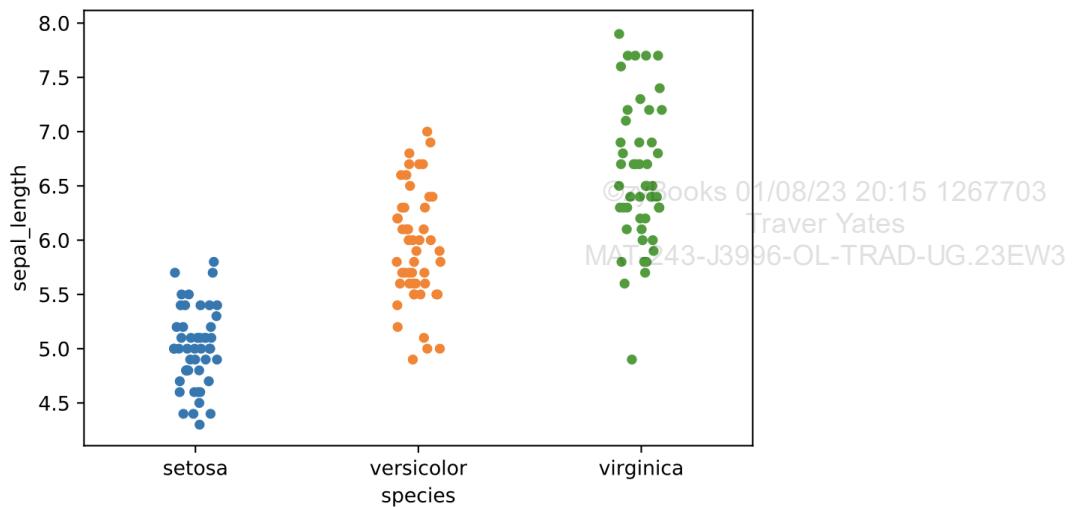
©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Data points may overlap when the dataset is dense. Data points may share the same value or the marks may be too large compared to the plot's resolution. To get a better sense of the dataset, a data analyst might use scatter plot *jittering*. **Jittering** is the addition of random noise to the plot in order to prevent or minimize overlapping data points. Ex: The jittered strip plot below shows 6 samples of *Iris virginica* had sepal lengths between 7.5 cm and 8 cm, which is not clear from the strip plot above.

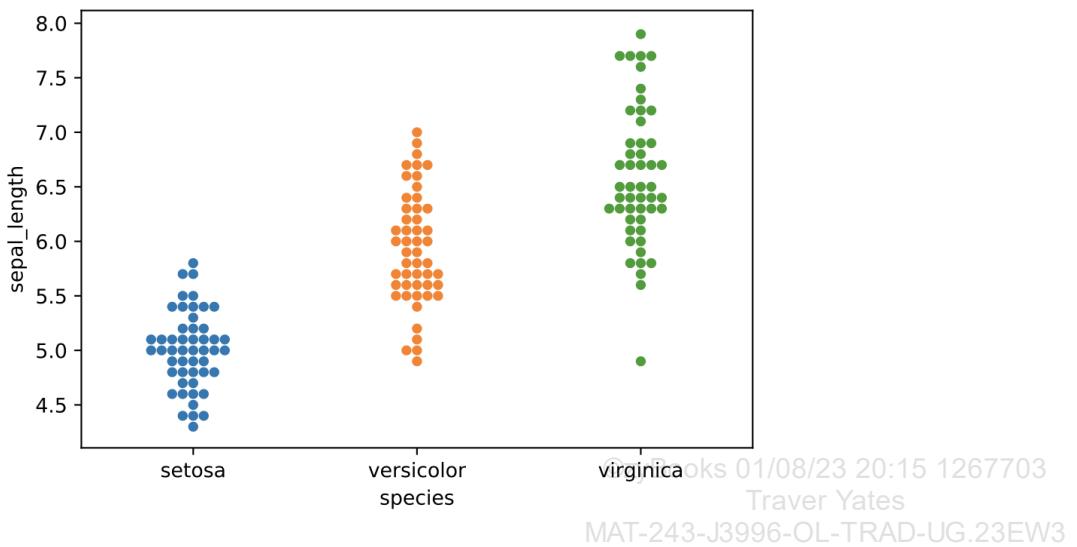
Figure 1.7.5: Jittered strip plot for the sepal lengths of iris species.



## Swarm plots

Jittered plots do not entirely prevent overlapping. An alternative approach to jittering is to use a **swarm plot**. A **swarm plot** uses a random algorithm to set a minimum distance between points.

Figure 1.7.6: Swarm plot for the sepal lengths of iris species.



### PARTICIPATION ACTIVITY

1.7.4: Strip plots and swarm plots.



Refer to the strip plots and swarm plots above.



1) The number of samples of each iris species in the dataset can be determined by looking at a strip plot.

- True
- False

2) Jittering gives a more accurate sense of a dataset's features.

- True
- False

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Python-Function 1.7.2: Strip and swarm plots.

A strip plot is created by using the `sns.stripplot` function of the `seaborn` library. Ex: The code below creates a strip plot where the horizontal axis is the deck to which a passenger of the Titanic is assigned and the vertical axis is the fare paid in dollars.

```
# loads the necessary libraries
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

# loads the titanic dataset
titanic = sns.load_dataset("titanic")

# title
plt.title('Fares paid by passengers of the Titanic by deck', fontsize=20)

# plot
sns.stripplot(x="deck", y="fare", data=titanic);

# saves the image
plt.savefig("titanicstripplot.png")

# shows the image
plt.show()
```

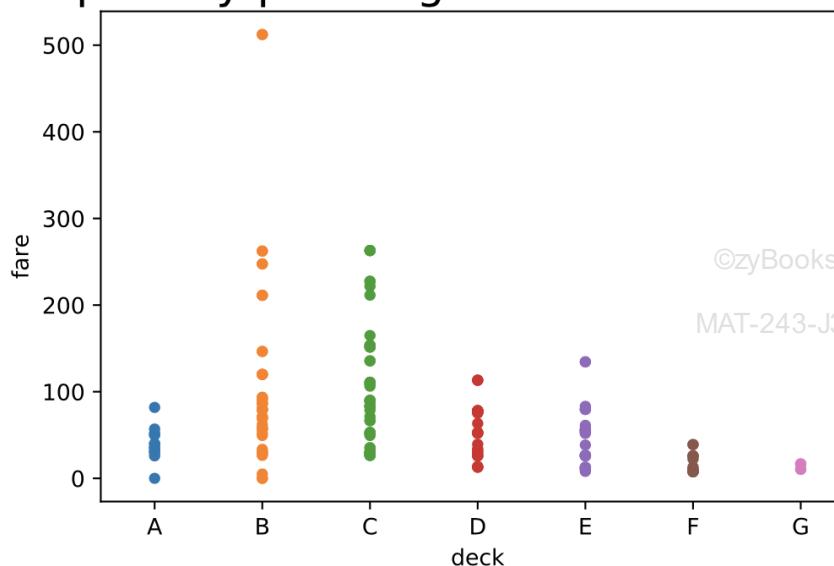
The resulting strip plot is shown below.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Fares paid by passengers of the Titanic by deck

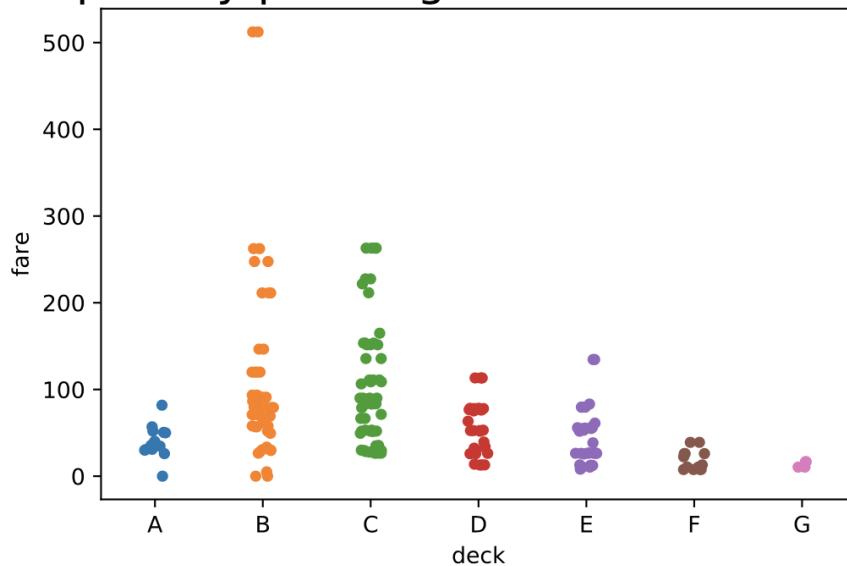


[Run example](#)

To add jittering to a strip plot, the parameter `jitter` should be set to true as shown below.

```
sns.stripplot(x="deck", y="fare", jitter= True, data=titanic);
```

## Fares paid by passengers of the Titanic by deck



[Run example](#)

A swarm plot is created by using the `swarmplot` function of the `seaborn` library. An optional parameter `hue` can also be set to another categorical variable to group each swarm according to the other variable. Ex: The code below displays a swarm plot for the fares paid by passengers by deck where each swarm is grouped by sex.

```
# loads the necessary libraries
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

# loads the titanic dataset
titanic = sns.load_dataset("titanic")

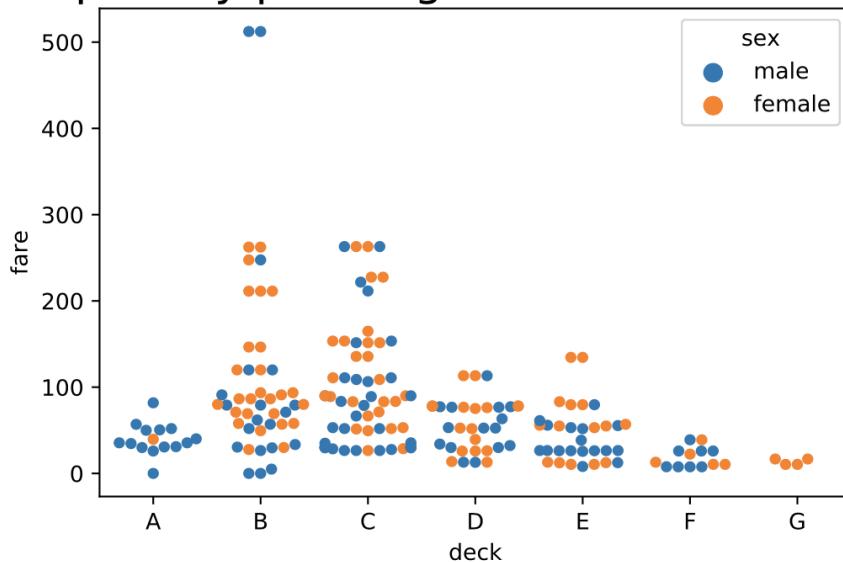
# title
plt.title('Fares paid by passengers of the Titanic by deck', fontsize=20)

# plot
sns.swarmplot(x="deck", y="fare", hue = "sex", data=titanic); ©zyBooks 01/08/23 20:15 1267703
Traver Yates
MAT-243-J3996-OL-TRAD-UG.23EW3

# saves the image
plt.savefig("titanicswarmplot.png")

# shows the image
plt.show()
```

## Fares paid by passengers of the Titanic by deck



[Run example](#)

## References

- (\*) "NCAA College Football Predictive Rankings and Ratings." *TeamRankings*.  
[www.teamrankings.com/college-football/ranking/predictive-by-other](http://www.teamrankings.com/college-football/ranking/predictive-by-other).

- (\*) "NCAA Salaries." *USA Today*. [sports.usatoday.com/ncaa/salaries/](http://sports.usatoday.com/ncaa/salaries/) ©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## 1.8 Line charts

### Introduction to line charts

A **line chart** (or **line graph**) depicts data trends by using straight lines to connect successive data points in a scatter plot. The straight lines show the general direction that data changes over time. Because trends involve time, line charts commonly use a time metric for the horizontal axis. Ex: Given the following data on Apple stock prices from March 2015 to March 2016, the following line chart shows how Apple's stock price changes (vertical axis) as each month passes (horizontal axis).

Table 1.8.1: Apple stock prices.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

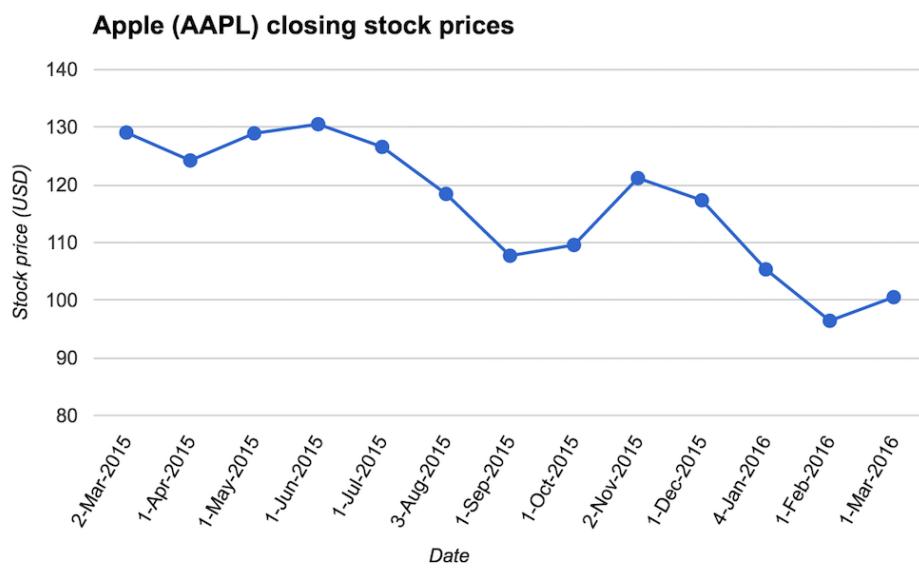
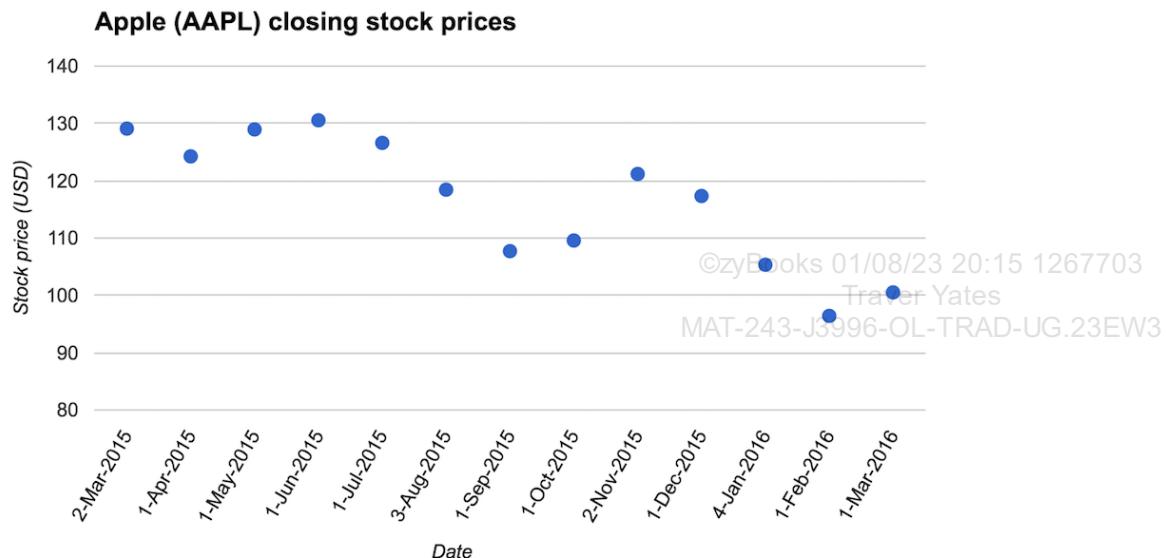
©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Date	Apple stock price (USD)
Mar 2, 2015	129.09
Apr 1, 2015	124.25
May 1, 2015	128.95
Jun 1, 2015	130.54
Jul 1, 2015	126.60
Aug 3, 2015	118.44
Sep 1, 2015	107.72
Oct 1, 2015	109.58
Nov 2, 2015	121.18
Dec 1, 2015	117.34
Jan 4, 2016	105.35
Feb 1, 2016	96.43
Mar 1, 2016	100.53

Source: [Yahoo! Finance, 2016](#)<sup>1</sup>

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Figure 1.8.1: Apple stock charts (March 2015 - March 2016) with and without lines.



The main benefit of a line graph is to quickly convey whether values are increasing, decreasing, or remaining constant between data points. Steeper lines indicate more rapid increases or decreases, while flatter lines indicate little change between data points. Ex: The line graph above clearly shows that the steepest increase in the stock value was from October 2015 to November 2015, which may lead investors to research what happened to Apple in October 2015.

Lines also help convey that values exist between data points. Ex: Although the Apple line chart shows two consecutive data points for July 1 and August 3, the stock price took on many values in between those dates. The line connecting July 1 and August 3 does not represent real data, but rather, a basic trend of the data change between data points.

#### PARTICIPATION ACTIVITY

##### 1.8.1: Interpreting line chart trends.

Consider the following line chart of Alphabet (Google's parent company) stock from March 2015 to March 2016:



- 1) Google's stock price decreased from December 2015 to January 2016.

True  
 False

- 2) Google's stock price never increases more than two months in a row.

True  
 False

- 3) Google's stock price had the largest increase from November to December.

True  
 False

- 4) Google's stock price reached a 52-week low in June.

True  
 False

- 5) Relative to the rest of the graph, Google's stock price remains mostly constant from April 2015 to June 2015.

True



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

False

- 6) \$575 is a likely price for Google stock  
on July 15, 2015.

 True False

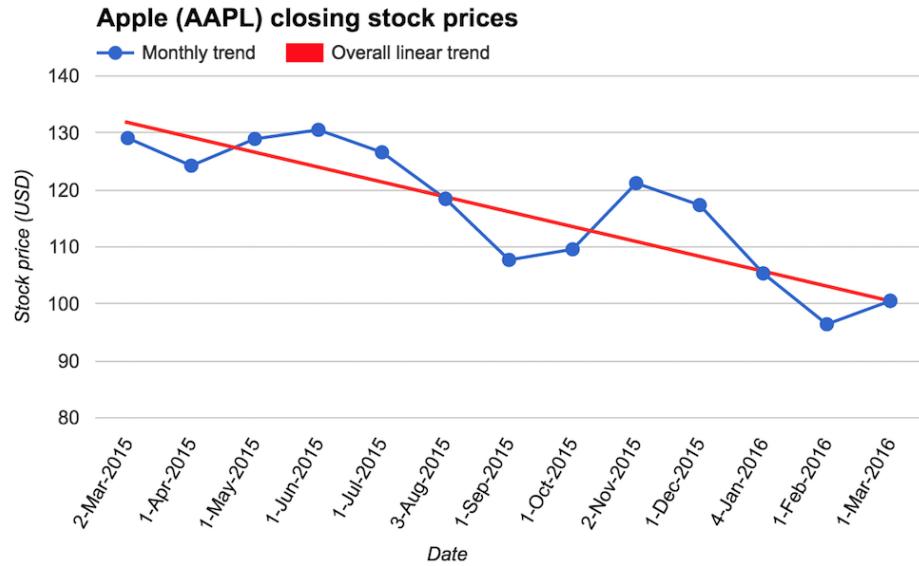
©zyBooks 01/08/23 20:15 1267703

Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

A **linear trend line** is a straight line that depicts the general direction data changes from the first to last data point, often added to summarize the entire chart. A good linear trend line is typically computed using various techniques such as linear regression (discussed elsewhere), and is not a simple connection of the first and last points.

In the Apple stock line chart below, a linear trend line is added in red and starts slightly above the first data point. While the stock price had two large increases from March 2015 to March 2016, the linear trend line clearly shows that the stock price tended to decrease during that time.

Figure 1.8.2: Apple stock prices (March 2015 - March 2016) line chart with overall trend line.



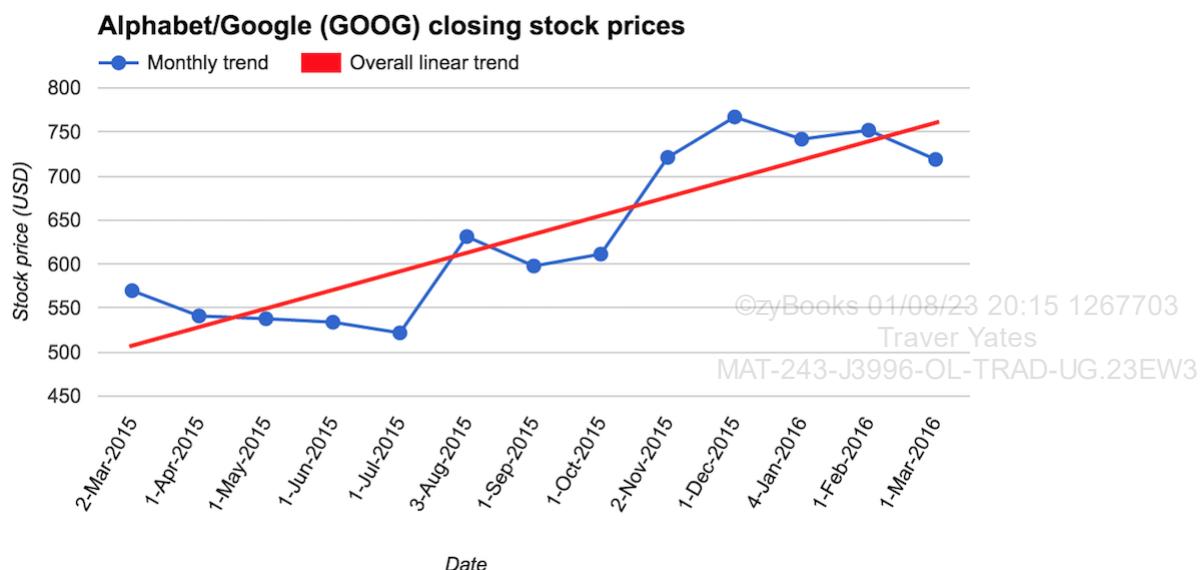
#### PARTICIPATION ACTIVITY

1.8.2: Interpreting linear trend lines on a line chart.

©zyBooks 01/08/23 20:15 1267703

Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Consider the following line chart of Alphabet (owner of Google) stock prices from March 2015 to March 2016:



1) What is Google's stock price in March 2015?

- Around \$500
- Around \$575
- Around \$720

2) What is the difference in Google's stock price between March 2015 and March 2016?

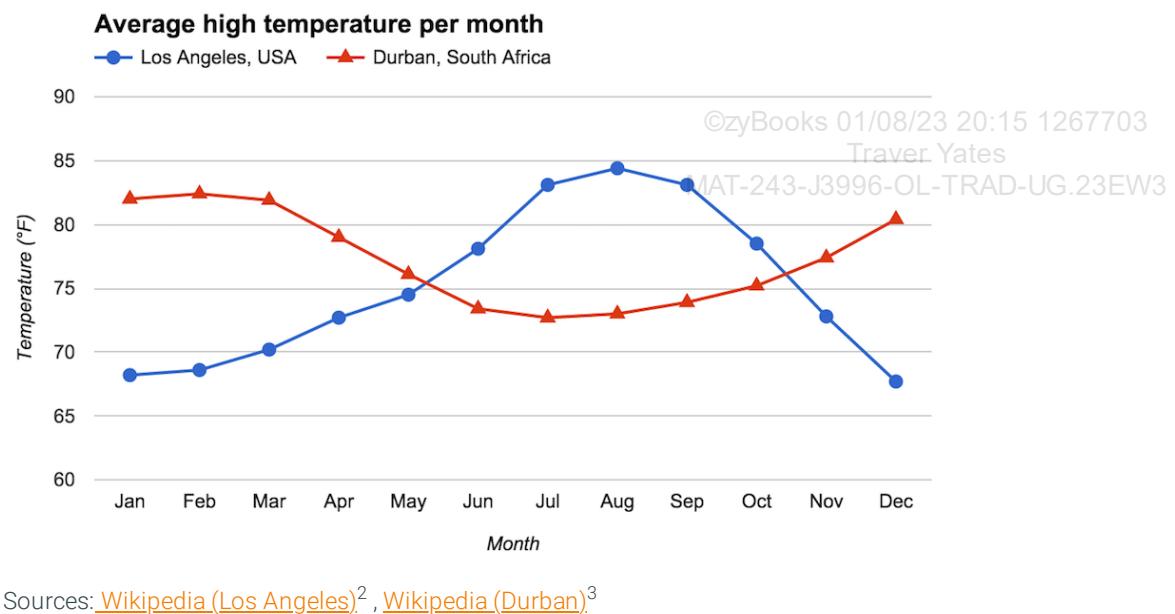
- About \$250
- About \$150

3) Based only on the chart, do Google's stock prices seem more likely to increase or decrease from March 2016 to March 2017 (1 year)?

- Increase
- Decrease

Multiple datasets are commonly shown in one line chart to highlight differences. Each dataset is distinguished by different color, data point shape, and/or line style, as noted by a legend. Ex: The line chart below shows how temperatures in Los Angeles, California, USA and Durban, South Africa differ per month due to being on opposite sides of the equator. The chart also shows how Los Angeles has more extreme temperature swings than Durban.

Figure 1.8.3: Line chart showing multiple datasets: average high temperatures for Los Angeles, USA, and Durban, South Africa.

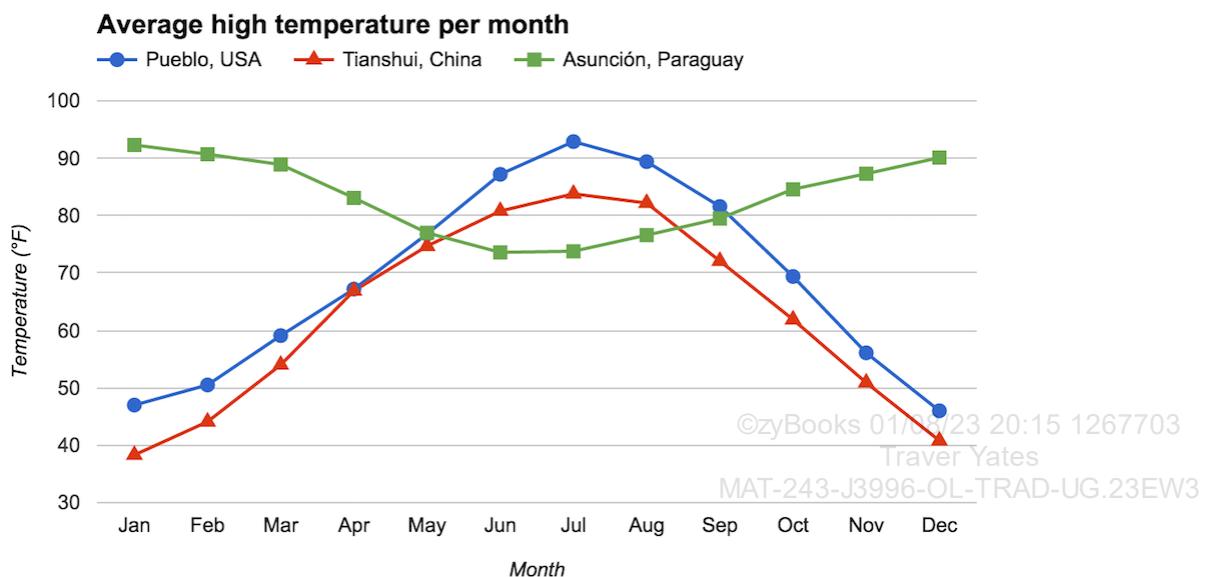


#### PARTICIPATION ACTIVITY

1.8.3: Interpreting multiple datasets on a line chart.



Consider the following line chart of average high temperatures for Pueblo, Colorado, USA, Tianshui, China, and Asunción, Paraguay.



Sources: [Wikipedia \(Pueblo\)](#)<sup>4</sup>, [Wikipedia \(Tianshui\)](#)<sup>5</sup>, [Wikipedia \(Asuncion\)](#)<sup>6</sup>

- What is the average high temperature for Pueblo in March?



- About  $90^{\circ}F$
  - About  $60^{\circ}F$
  - About  $53^{\circ}F$
- 2) Which city has the lowest average high temperature at any point during the year? □
- Asunción
  - Tianshui
  - Pueblo
- 3) In which month are the average high temperatures for Pueblo, Tianshui, and Asuncion the most similar? □
- April
  - May
  - September
- 4) Which city has average high temperatures most unlike the other cities? □
- Asunción
  - Tianshui
  - Pueblo

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Python-Function 1.8.1: Line charts.

The code below loads the dataset containing unemployment rates in the United States from 1980 to 2017 and plots the data as a line chart.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

```
# imports the necessary modules
import pandas as pd
import matplotlib.pyplot as plt

# loads the unemployment dataset
unemployment = pd.read_csv('http://data-analytics.zybooks.com/unemployment.csv')

# title
plt.title('U.S. unemployment rate', fontsize = 20)

# x and y axis labels
plt.xlabel('Year')
plt.ylabel('% of total labor force')

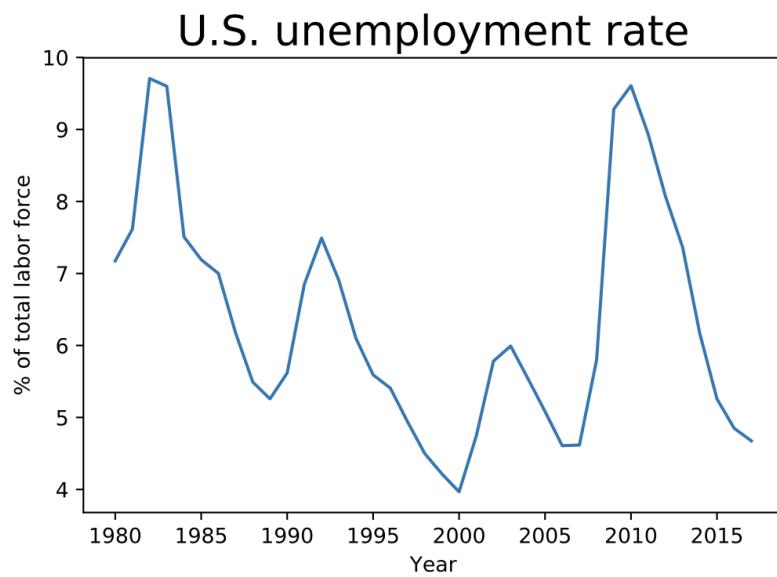
# plot
plt.plot(unemployment["Year"], unemployment["Value"])

# saves the image
plt.savefig("linechart.png")

# shows the image
plt.show()
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The resulting line chart is shown below.



#### [Run example](#)

One of the most useful libraries that includes a multitude of economic data such as stock prices, futures, and unemployment rates is **quandl**. The **quandl** library allows users to preprocess, subset, and change the format of data frames, using a much more convenient syntax than **pandas**, and is especially suited for time series analysis. To use **quandl**, the library needs to be downloaded and installed by running `pip install quandl` from the command line.

For more information about the **quandl** library and specific methods, see the [documentation](#).

The code below creates a line chart for the closing stock prices of Target (TGT) since 1983.

```
# loads the necessary libraries
import quandl
import matplotlib.pyplot as plt

# creates a data frame containing TGT stock data
tgt = quandl.get('WIKI/TGT')

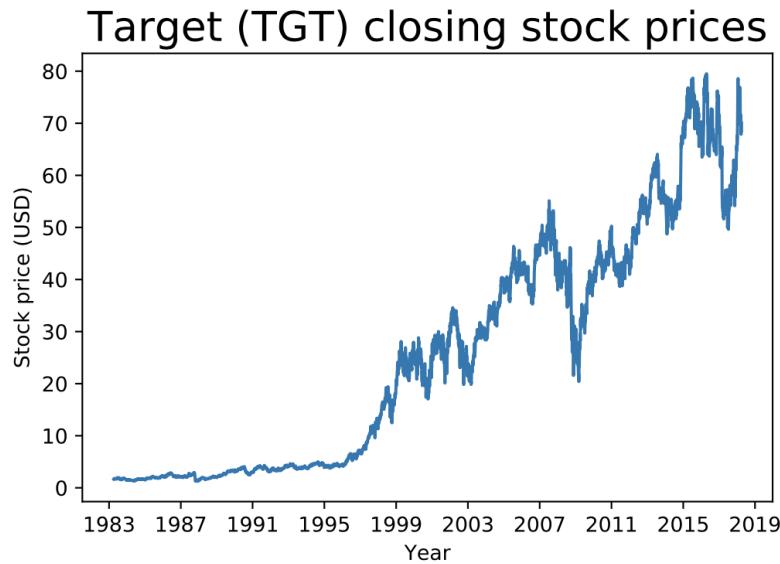
# title
plt.title('Target (TGT) closing stock prices', fontsize=20)

# x and y axis labels
plt.xlabel('Year');
plt.ylabel('Stock price (USD)');

# plot
plt.plot(tgt.index, tgt['Adj. Close'])

# saves the image
plt.savefig("tgtstocklinechart.png")
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

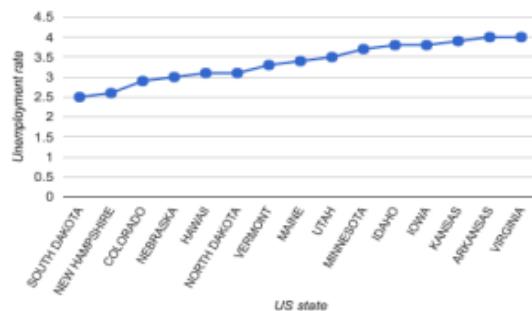


## Misuse of line charts

### Using line charts to represent categorical data

A line chart should not be used for nominal categorical data. Lines suggest some relation from one item to the next, but nominal variables have no ordering so can have no such relation. Ex: The plot below on the left inappropriately shows lines, even though no relationship exists between South Dakota and New Hampshire, for example. However, representing the unemployment rate of each state with a point, without connecting the individual points, would be appropriate. Using a bar chart is also common.

Figure 1.8.4: A line chart is not appropriate for categorical data.



©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**PARTICIPATION ACTIVITY**

## 1.8.4: Line charts.



When is a line chart appropriate?

- 1) The  $x$ -axis is the year ranging from 2000 to 2015, the  $y$ -axis is the amount of rainfall.



- Appropriate
- Not appropriate

- 2) The  $x$ -axis shows a movie rating: G, PG, PG-13, R, and NC-17. The  $y$ -axis is the number of movies released in 2015 of a given rating.



- Appropriate
- Not appropriate

- 3) The  $x$ -axis is a color: black, blue, green, red, silver, or white. The  $y$ -axis is the number of cars of a given color.



- Appropriate
- Not appropriate

©zyBooks 01/08/23 20:15 1267703

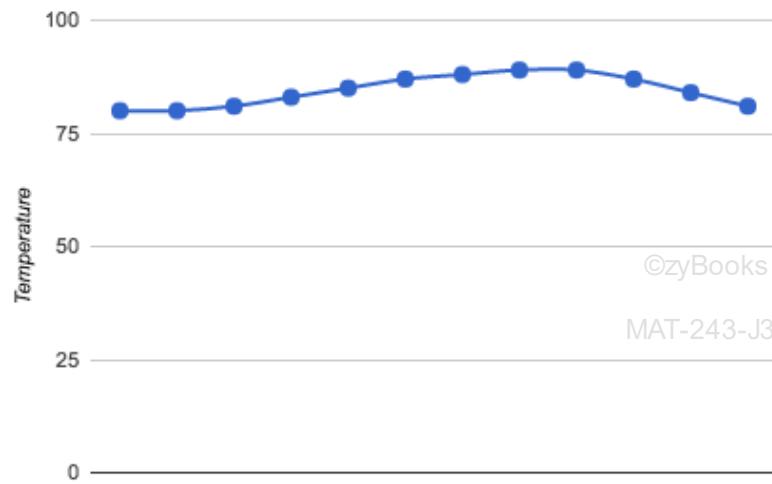
Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Omitting labels**

Common mistakes are to forget to label an axis or to forget to provide units. Without such information, a viewer cannot appropriately interpret the information.

Figure 1.8.5: Average monthly temperatures in Hawaii, with some labels/units missing.



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

**PARTICIPATION ACTIVITY**

## 1.8.5: Missing labels or units.



Refer to the above figure.

1) The units for Temperature are \_\_\_\_.



- present
- missing

2) The label for the  $x$ -axis is \_\_\_\_.



- present
- missing

3) The label for the  $y$ -axis is \_\_\_\_.



- present
- missing

4) If a  $y$ -axis is labeled % or ratio, as in % revenue, or ratio of male/female births, additional units are \_\_\_\_.



- required
- not necessary

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## References

(\*1) "Apple Inc." Yahoo! Finance . 2016, finance.yahoo.com/echarts?s=aapl.

(\*2) Wikipedia Contributors. "Los Angeles." *Wikipedia, The Free Encyclopedia*. Retrieved 17 July 2018, en.wikipedia.org/wiki/Los\_Angeles.

(\*3) Wikipedia Contributors. "Durban." *Wikipedia, The Free Encyclopedia*. Retrieved 17 July 2018, en.wikipedia.org/wiki/Durban.

(\*4) Wikipedia Contributors. "Pueblo, Colorado." *Wikipedia, The Free Encyclopedia*. Retrieved 17 July 2018, en.wikipedia.org/wiki/Pueblo,\_Colorado.

©zyBooks 01/08/23 20:15 1267703

(\*5) Wikipedia Contributors. "Tianshui." *Wikipedia, The Free Encyclopedia*. Retrieved 17 July 2018, en.wikipedia.org/wiki/Tianshui.

MAT-243-J3996-OL-TRAD-UG.23EW3

(\*6) Wikipedia Contributors. "Asuncion." *Wikipedia, The Free Encyclopedia*. Retrieved 17 July 2018, en.wikipedia.org/wiki/Asunción.

## 1.9 Box plots

### Minimum, maximum, and range

In addition to standard deviation and variance, the minimum, maximum, and range of a dataset can describe the spread of the dataset.

The **maximum** of a dataset is the largest value in the dataset. The **minimum** of a dataset is the smallest value in the dataset. The **range** of a dataset is the difference between the maximum and minimum of the dataset.

Example 1.9.1: Minimum, maximum, and range.

Find the minimum, maximum, and range of the dataset  $-5, 3, 0, -1, 4, 7$ .

#### Solution

- The largest value in the dataset is **7**. Thus, the maximum is **7**.
- The smallest value in the dataset is **-5**. Thus, the minimum is **-5**.
- The range is the difference between the maximum and the minimum, or  $7 - (-5) = 12$ .

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Python-Function 1.9.1: min(), max(), and range.

The DataFrame.min() function is used to find the minimum, and the DataFrame.max() function is used to find the maximum.

```
import pandas as pd
scores = pd.read_csv('http://data-analytics.zybooks.com/ExamScores.csv')

# Prints the minimum score for each exam
print(scores.min())

# Prints the minimum score for Exam1 only
print(scores[['Exam1']].min())

# Prints the maximum score for each exam
print(scores.max())

# Prints the maximum score for Exam1 only
print(scores[['Exam1']].max())
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Exam1	59
Exam2	49
Exam3	18
Exam4	55
	dtype: int64
Exam1	59
	dtype: int64
Exam1	100
Exam2	100
Exam3	100
Exam4	100
	dtype: int64
Exam1	100
	dtype: int64

No pre-built function exists for finding the range. However, the range can be found by subtracting the minimum values from the maximum values. Alternatively, a new function can be defined if the range needs to be found repeatedly.

```
import pandas as pd
scores = pd.read_csv('http://data-analytics.zybooks.com/ExamScores.csv')

# Calculating the range by subtracting the minimum from the maximum
score_range = scores.max() - scores.min()
print(score_range)

# Defining a function that can be used repeatedly
def range_of_scores(x):
    return x.max() - x.min()
print(range_of_scores(scores))
print(range_of_scores(scores[['Exam1']]))

print(range_of_scores(scores[['Exam2']]))

print(range_of_scores(scores[['Exam4']])))
```

Exam1	41
Exam2	51
Exam3	82
Exam4	45
	dtype: int64
Exam1	41
Exam2	51
Exam3	82
Exam4	45
	dtype: int64
Exam1	41
	dtype: int64
Exam2	51
	dtype: int64
Exam4	45
	dtype: int64

[Run example](#)

### PARTICIPATION ACTIVITY

1.9.1: Minimum, maximum, and range.



Use the dataset  $-3, 5, 8, 1, -6, 4$  to answer the following.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- What is the maximum?

**Check**

**Show answer**

- What is the minimum?



3) What is the range?



©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Percentiles

The *n*th **percentile** of a dataset is the data value such that *n* percent of the data falls at or below that value. Three percentiles are particularly important.

- The **first quartile** (*Q*<sub>1</sub>) is the 25th percentile. One-quarter of the data fall at or below *Q*<sub>1</sub>. The first quartile is the median of the lower half of the data.
- The **third quartile** (*Q*<sub>3</sub>) is the 75th percentile. Three-quarters of the data fall at or below *Q*<sub>3</sub>. The third quartile is the median of the upper half of the data.
- Because half of the data fall at or below the median, the median is also the 50th percentile of a dataset.

Collectively, the minimum and maximum values, *Q*<sub>1</sub>, median, and *Q*<sub>3</sub> form a set of descriptive statistics called the **five-number summary**.

### Example 1.9.2: Creating a five-number summary.

The number of receptions made by players on a certain American football team are given by the dataset 3, 37, 23, 61, 36, 65, 6, 24, 1, 19, 72, 1, 13, 40, 1. Create a five-number summary of this data.

#### Solution

First, the data should be sorted in ascending or descending order. In ascending order, the dataset is

1, 1, 1, 3, 6, 13, 19, 23, 24, 36, 37, 40, 61, 65, 72

Thus, the minimum is 1 and the maximum is 72. The dataset contains 15 values, so the median is the  $\frac{15+1}{2} = 8$ th value of 23. Thus, the lower half of the data is 1, 1, 1, 3, 6, 13, 19, 23 and the upper half of the data 23, 24, 36, 37, 40, 61, 65, 72 (the

median belongs to both the upper and lower halves). The median of the lower half is

$$\frac{3+6}{2} = \frac{9}{2} = 4.5 \text{ and the median of the upper half is } \frac{37+40}{2} = \frac{77}{2} = 38.5.$$

The five-number summary is

Minimum	1
$Q_1$ (first quartile)	4.5
Median	23
$Q_3$ (third quartile)	38.5
Maximum	72

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## Python-Function 1.9.2: describe().

The `DataFrame.describe()` function generates a five-number summary.  $Q_1$ , the 25th percentile, is shown as 25%, the median, the 50th percentile, is shown as 50%, and  $Q_3$ , the 75th percentile, is shown as 75%.

```
import pandas as pd
scores = pd.read_csv('http://data-
analytics.zybooks.com/ExamScores.csv')

# Prints the summary for each exam
print(scores.describe())

# Prints the summary for Exam1 only
print(scores[['Exam1']].describe())
```

	Exam1	Exam2
Exam3	Exam4	
count	50.000000	50.000000
50.000000	50.000000	
mean	82.700000	79.400000
73.340000	76.500000	
std	9.291756	14.33278
21.754296	8.05656	
min	59.000000	49.000000
18.000000	55.000000	
25%	77.250000	68.500000
59.000000	72.000000	
50%	83.000000	79.500000
74.500000	75.000000	
75%	88.500000	91.750000
94.500000	79.750000	
max	100.000000	100.000000
100.000000	100.000000	
	Exam1	
count	50.000000	
mean	82.700000	
std	9.291756	
min	59.000000	
25%	77.250000	
50%	83.000000	
75%	88.500000	
max	100.000000	

[Run example](#)

**PARTICIPATION ACTIVITY**

## 1.9.2: Five-number summary.



Complete the five-number summary for the dataset  $0, -6, 10, 5, 8, 2, -12, 11, -2$ .

1) Minimum

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Check****Show answer**

2)  $Q1$

**Check****Show answer**

3) median

**Check****Show answer**

4)  $Q3$

**Check****Show answer**

5) Maximum

**Check****Show answer**

Python-Practice 1.9.1: Maximum, minimum, and range. ©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The [rent](#) dataset is first imported.

```
import pandas as pd
rent = pd.read_csv('http://data-analytics.zybooks.com/rent.csv')
print(rent)
```

	Santa Monica CA	Boise ID	Tucson AZ	Detroit MI	Pittsburgh PA	Orlando FL
0	10230	1600	2495	3195	2480	2242
1	10000	1500	2200	2695	2435	2000
2	9000	1029	2150	2595	2405	1912
3	8500	1025	1800	2495	2350	1895
4	8250	980	1650	2495	2320	1800
5	8000	950	1600	1675	2316	1765
6	7000	950	1500	1525	2305	1685
7	6500	925	1500	1480	2290	1670
8	6000	925	1500	1410	2275	1665
9	5815	900	1500	1400	2265	1625

©zyBooks 01/08/23 20.131267703  
Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

The maximum and minimum rent for all cities are found as follows.

<code>print(rent.max())</code>	Santa Monica CA 10230
<code>print(rent.min())</code>	Boise ID 1600
	Tucson AZ 2495
	Detroit MI 3195
	Pittsburgh PA 2480
	Orlando FL 2242
	dtype: int64
	Santa Monica CA 5815
	Boise ID 900
	Tucson AZ 1500
	Detroit MI 1400
	Pittsburgh PA 2265
	Orlando FL 1625
	dtype: int64

The range in rent for all cities is found by subtracting the minimum rents from the maximum rents.

<code>print(rent.max() - rent.min())</code>	Santa Monica CA 4415
	Boise ID 700
	Tucson AZ 995
	Detroit MI 1795
	Pittsburgh PA 215
	Orlando FL 617
	dtype: int64

[Run example](#)

## Introduction to box plots

A **box plot** is a data visualization that uses a box and several lines to depict the distribution of data in a dataset. A box spans the middle 50% of the data, with  $Q1$  as the lower boundary of the box and  $Q3$  as the upper boundary of the box. The median is shown as a line inside the box. Two lines, known as whiskers, extend from the lower boundary of the box to the minimum and from the upper boundary of the box to the maximum. The whiskers represent the lower and upper 25% of the data.

The following animation shows the creation of a box plot using data from a previous example.

### PARTICIPATION ACTIVITY

1.9.3: Creating a box plot.



## Animation content:

undefined

## Animation captions:

1. To create a box plot, the dataset must be sorted in ascending or descending order.
2. An axis with the minimum and maximum data points as the endpoints shows the range of the data.
3. The median is the middle number in the ordered data and is represented as a line inside the box.
4.  $Q_1$  is the median of the lower half of the data and forms the lower boundary of the box.  
$$Q_1 = (3 + 6) \div 2 = 4.5$$
5.  $Q_3$  is the median of the upper half of the data and forms the upper boundary of the box.  
$$Q_3 = (37 + 40) \div 2 = 38.5$$
6. The box is formed and the whiskers are extended from the lower boundary of the box to the minimum and from the upper boundary of the box to the maximum.

## Python-Function 1.9.3: Box plots.

The `boxplot` function of the `seaborn` creates a box plot.

```
# loads the necessary libraries
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd

# loads the ExamScores dataset
scores = pd.read_csv('http://data-analytics.zybooks.com/ExamScores.csv')

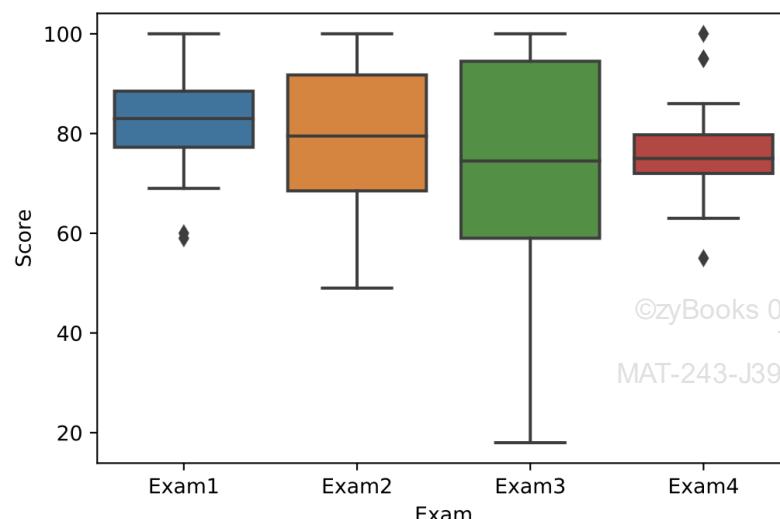
# transforms the data
df = pd.melt(scores, value_name = "Score", var_name = "Exam")

# plot
sns.boxplot(x="Exam", y="Score", data=df);

# saves the image
plt.savefig("Examsboxplots.png")

# shows the image
plt.show()
```

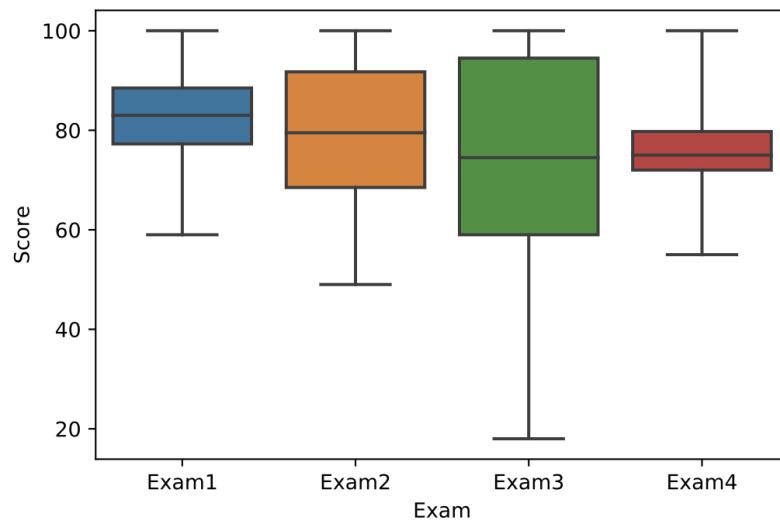
©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

[Run example](#)

To create a box plot that does not identify the outliers, the parameter `whis` can be set to a high number such as `100`. `1.5` is the default value for `whis`.



[Run example](#)

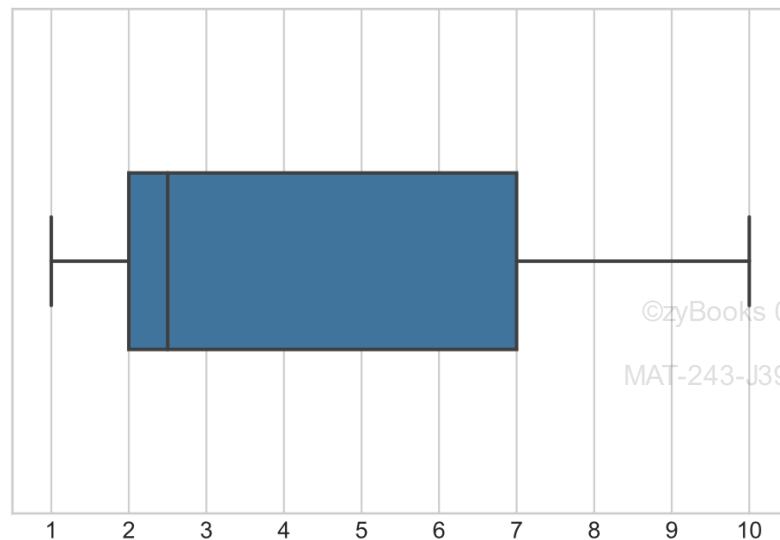
#### PARTICIPATION ACTIVITY

1.9.4: Characteristics of a box plot.



Match the value to the corresponding term, based on the following box plot:

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Select the definition that matches each term

1) 1

- The median value of the data set
- The data points within the box
- The minimum value appearing in the data set
- $Q_3$
- $Q_1$
- The number of data points in the data set
- The maximum value appearing in the data set

2) 10

- The median value of the data set
- The data points within the box
- The minimum value appearing in the data set
- $Q_3$
- $Q_1$
- The number of data points in the data set
- The maximum value appearing in the data set

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

3) 2.5

- The median value of the data set
- The data points within the box
- The minimum value appearing in the data set
- $Q_3$
- $Q_1$
- The number of data points in the data set
- The maximum value appearing in the data set

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

4) Unknown

- The median value of the data set
- The data points within the box
- The minimum value appearing in the data set
- $Q_3$
- $Q_1$
- The number of data points in the data set
- The maximum value appearing in the data set

5) 2

- The median value of the data set
- The data points within the box
- The minimum value appearing in the data set
- $Q_3$
- $Q_1$
- The number of data points in the data set
- The maximum value appearing in the data set

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

6) 7

- The median value of the data set
- The data points within the box
- The minimum value appearing in the data set
- $Q3$
- $Q1$
- The number of data points in the data set
- The maximum value appearing in the data set

7) 50%

- The median value of the data set
- The data points within the box
- The minimum value appearing in the data set
- $Q3$
- $Q1$
- The number of data points in the data set
- The maximum value appearing in the data set

**Reset**

A box plot helps visualize a data set's distribution, giving more information than just the mean or median. The box plot below shows the distribution of the percentages of the total population of the United States for each of the 50 states.<sup>1</sup> The box plot shows the median is 1.37%, whereas the mean is 1.96%. The box plot shows that the upper 25% of states range from 2.2% to 12.15% of the total population, which significantly affects political representation, resource distribution, and other important factors.

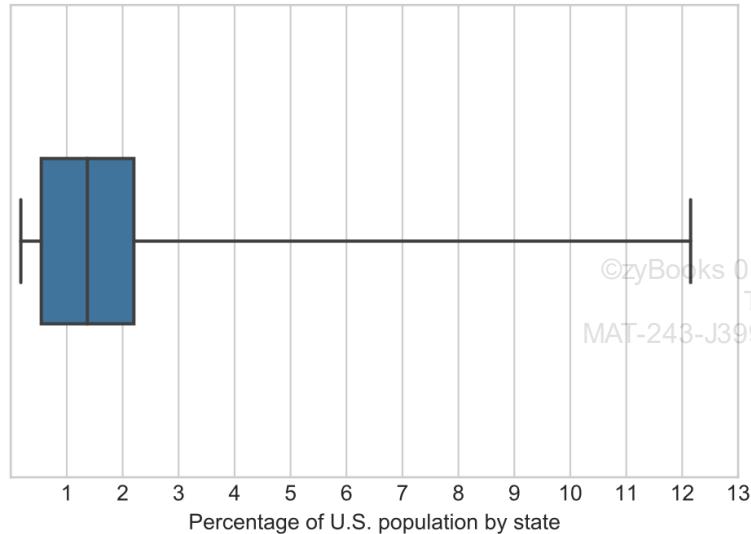
©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

The **skew** is the difference between the mean and the median. A positive skew means that the distribution is skewed to the right, while a negative skew means that the distribution is skewed to the left. In the box plot below, the skew is  $1.96\% - 1.37\% = 0.59\%$ , which means that the distribution is skewed to the right.

Figure 1.9.1: Box plot showing U.S. population distribution by state.



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

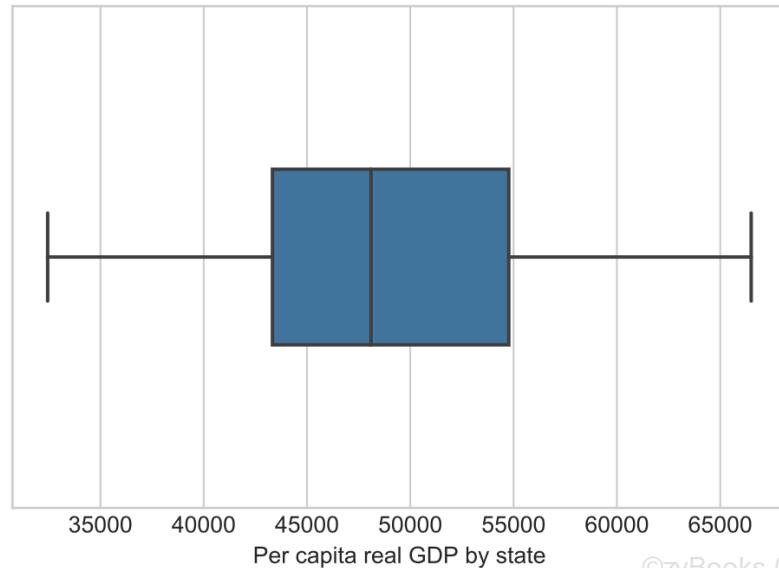
[Source](#) U.S. Census Bureau

PARTICIPATION  
ACTIVITY

1.9.5: Interpreting a box plot.



The following box plot shows the distribution of the per capita real gross domestic product (GDP) in 2017 of each U.S. state<sup>2</sup>. The per capita GDP for the entire U.S. is \$51,749.



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

[Source](#): U.S. Bureau of Economic Analysis

- 1) The real GDP of \$51,749 is a good representation per capita GDP of the United States.

- True
- False





2) The national mean income of **\$51,749** is a good representation of each state's individual mean income.

- True
- False

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Detecting outliers

One way to detect outliers using a box plot is to determine how far each data element is from either  $Q_1$  or  $Q_3$ . The **interquartile range (IQR)** of a dataset is the difference between  $Q_3$  and  $Q_1$  ( $Q_3 - Q_1$ ), or the length of the box in a box plot. A data value greater than  $Q_3 + 1.5(IQR)$  or less than  $Q_1 - 1.5(IQR)$  is considered an outlier. Often, an outlier is not included in either whisker and is instead represented in the plot as a marker such as an open circle or a triangle.

Example 1.9.3: The interquartile range.

For the dataset  $3, 37, 23, 61, 36, 65, 6, 24, 1, 19, 72, 1, 13, 40, 1$ ,  $Q_1 = 4.5$  and  $Q_3 = 38.5$ . What is the IQR, and does the dataset contain any outliers?

### Solution

The IQR is  $38.5 - 4.5 = 34$ . An outlier is either greater than

$Q_3 + 1.5(IQR) = 38.5 + 1.5(34) = 89.5$  or less than

$Q_1 - 1.5(IQR) = 4.5 - 1.5(34) = -46.5$ . Since no data values are greater than 89.5 or less than -46.5, no outliers exist.

Python-Practice 1.9.2: Five-number summary and box plot.

[mtcars](#) is a historical dataset from a 1974 issue of Motor Trend comparing the performance of 32 cars. The dataset should first be imported as a pandas DataFrame using the `read_csv()` function. To see the structure of the dataset, the `head()` function can be used to display the first few lines of the dataset.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

```
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt

# reads the mtcars.csv file into a dataframe called mtcars
mtcars = pd.read_csv("https://data-analytics.zybooks.com/mtcars.csv")

# prints the first few lines of mtcars
print(mtcars.head())
```

@zyBooks 01/08/23 20:15 1267703

Traver Yates

	Unnamed: 0	mpg	cyl	disp	hp	drat	wt	M	qsec	vs	am	gear	R	carb
0	Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4		
1	Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4		
2	Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1		
3	Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1		
4	Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2		

The `describe()` function gives the five-number summary of the weight of the cars.

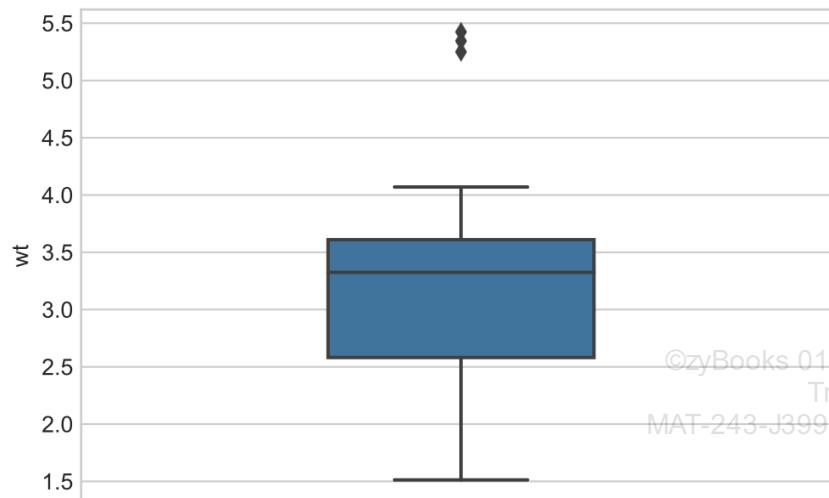
```
mtcars.wt.describe()
```

count	32.000000
mean	3.217250
std	0.978457
min	1.513000
25%	2.581250
50%	3.325000
75%	3.610000
max	5.424000
Name:	wt, dtype: float64

In the command above, `.wt` is needed to display only the summary for the weight. Otherwise, the five-number summaries for each of the attributes or columns are displayed.

The command to display the box plot for the weights of the cars as well as the corresponding output are given below.

```
sns.boxplot(mtcars.wt, width=0.35)
```



@zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

[Run example](#)

**PARTICIPATION  
ACTIVITY**

## 1.9.6: Detecting outliers.

Refer to the Python-Practice above.

- 1) What is the interquartile range for the weights data?

- 3.911
- 0.108
- 1.029

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- 2) What is the upper bound for the whiskers?

- 1.544
- 5.154
- 5.424

- 3) What is the lower bound for the whiskers?

- 1.038
- 2.581
- 1.513

- 4) How many outliers are in the data?

- 0
- 3
- Cannot be determined

- 5) Is a car with a weight of **5.250** (5, 250 lbs) an outlier?

- Yes
- No

- 6) Is a car with a weight of **1.513** (1, 513 lbs) an outlier?

- Yes
- No

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## References

(\*) United States Census Bureau. "State Population Totals and Components of Change: 2010-2017." census.gov. 8 May 2018. Web. 4 Jun. 2018.

(\*) United States Bureau of Economic Analysis. "Per capita real GDP by state (chained 2009 dollars)" bea.gov. 4 May 2018. Web. 4 Jun. 2018.

# 1.10 Histograms

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Histograms with evenly-sized bins

A **frequency distribution** is a table that displays how often an outcome occurs for a sample. To construct a frequency distribution, the data set is divided into mutually exclusive *classes*. A **class** is either a value of a categorical variable or an interval of a continuous variable. The **frequency** of a class is the number of events or values that fall under each class. Ex: An informal poll among a group of friends tallies how many people have  $x$  gaming applications on their phone. The results of the poll can be summarized in the frequency distribution below.

Table 1.10.1: Frequency distribution showing the number of people in a group having  $x$  gaming apps on their phone.

Gaming apps ( $x$ )	Tally	Frequency
0		4
1		7
2		5
3		3
4		2

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

The most common graphical representation of a frequency distribution is a histogram. A **histogram** depicts data values by splitting a continuous variable into a number of **class intervals**, each known as a **bin**. The simplest and most common type of histogram has bins of equal size. When bin sizes are equal, bins have rectangular bars with heights representing the frequency, which is the number of values in a bin.

The  $x$ -axis contains a continuous number line with ticks that represent bin boundaries. A bin includes values equal to or greater than the lower boundary, but less than the upper boundary (lower  $\leq$  value  $<$  upper). Gaps between rectangles are removed to show that the data is continuous.

**PARTICIPATION ACTIVITY**

1.10.1: Histogram of number of tickets per miles per hour over speed limit.


**Animation captions:**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

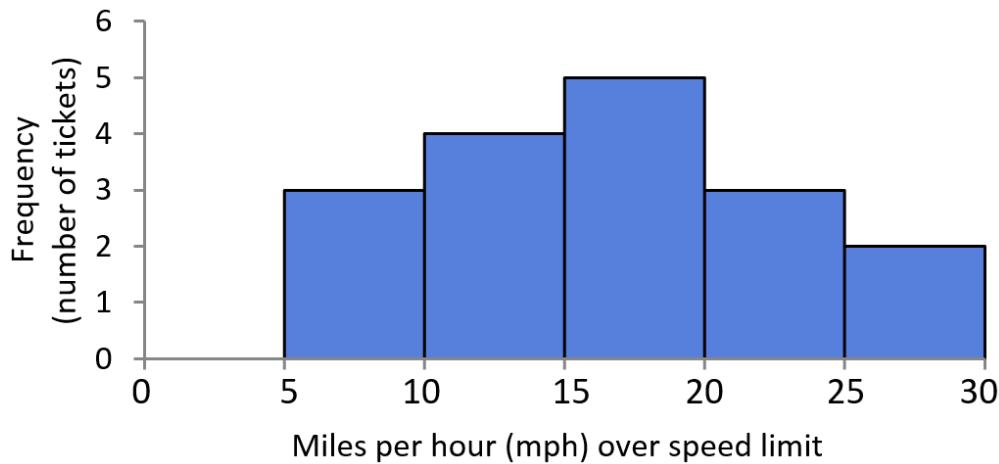
1. Data of MPH over the speed limit for **15** tickets can be represented as a histogram.
2. The  $x$ -axis of the histogram is continuous with evenly-spaced bin intervals. In this case, each interval is **5 MPH**, but other intervals are possible.
3. The  $y$ -axis represents the frequency, or the number of tickets.
4. Each bin frequency is represented by a rectangular column.

**PARTICIPATION ACTIVITY**

1.10.2: Histogram fundamentals.



Consider the following histogram showing speeding tickets issued by Officer Brown.



- 1) **3** speeding tickets were issued between **5 – 10 mph**.

True

False



- 2) The most speeding tickets were given in the **15 – 20 mph** range.

True

False

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



- 3) The number of speeding tickets issued between **0 – 5 mph** is unknown.



True False

- 4) A ticket issued that is 5 mph over the speed limit should be placed in the 0 – 5 mph bin.

 True False

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Histogram bin size

A key goal of a histogram is to estimate the probability density function of the continuous variable on the  $x$ -axis. In short, the goal is to fit a smooth curve over the most rectangles, while minimizing the white space under the curve.

When creating a histogram, multiple bin sizes should be attempted to determine the best distribution of the data. A good rule of thumb is to start with a bin size so that the number of bins is roughly equal to the square root of the number of values. Ex: For Officer Smith's 15 tickets seen in the animation above, a good number of bins to start with are  $\sqrt{15} = 3.9 \approx 4$  bins. Since the tickets are as much as 28 mph over the limit, a good initial bin size would be  $28 \text{ mph} / 4 \text{ bins} = 7 \text{ mph bins}$ .

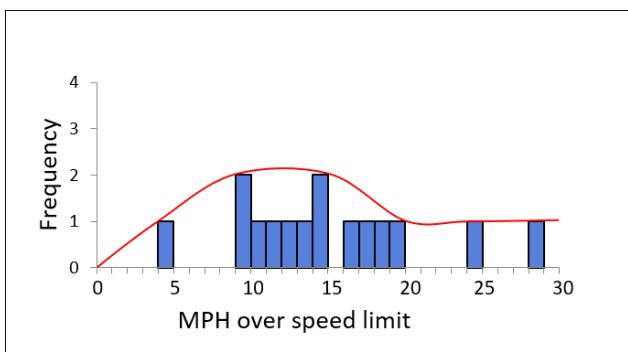
The figure below shows red distribution curves and histograms for Officer Smith's tickets for various bin sizes. Histogram-1 and Histogram-2 contain many gaps between bins, leaving too much white space under the curve. Thus, bin sizes 1 and 2 are not good options. Histogram-15 leaves little white space under the curve, but offers little insight about Officer Smith's ticketing trends. Finally, when compared to Histogram-10, Histogram-5 shows less white space under the curve, and thus, is the best option.

Figure 1.10.1: Histogram bin-size comparison for speeding tickets issued by Officer Smith.

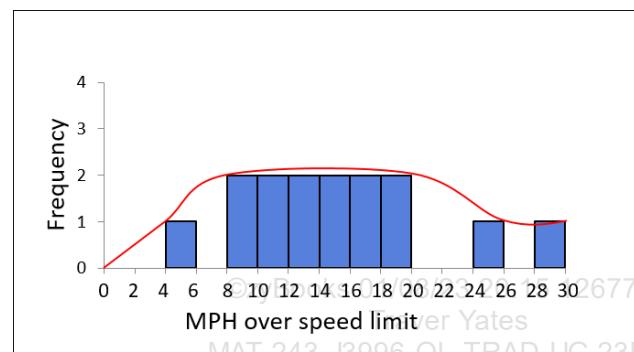
©zyBooks 01/08/23 20:15 1267703

Traver Yates

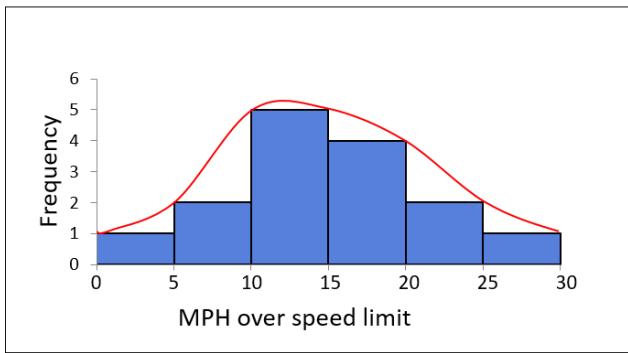
MAT-243-J3996-OL-TRAD-UG.23EW3



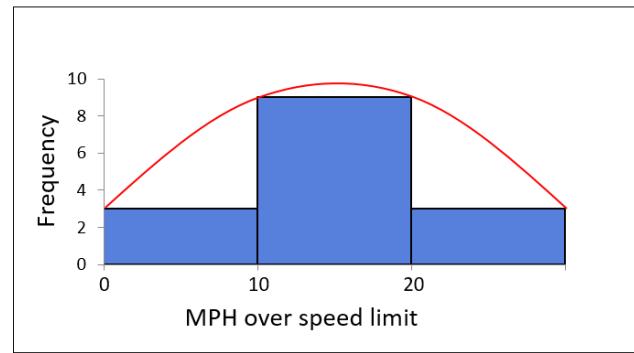
Histogram-1 (Bin size = 1)



Histogram-2 (Bin size = 2)



Histogram-5 (Bin size = 5)



Histogram-10 (Bin size = 10)



Histogram-15 (Bin size = 15)

Several basic distribution patterns should be looked for when selecting a bin size, because some situations are known to follow a certain distribution. The standard bell curve is a unimodal distribution pattern. A **unimodal distribution** occurs when there is one (uni) prevalent peak (mode) in the histogram. Ex: In Histogram-5, the 10 – 15 mph bin has the highest frequency of all bins, and thus, is the single mode.

Other common distribution patterns are listed below.

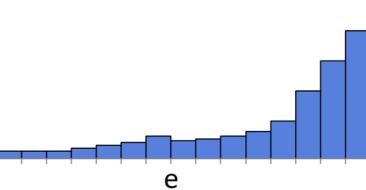
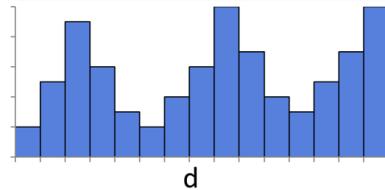
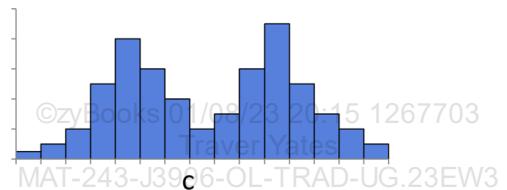
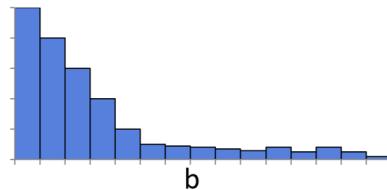
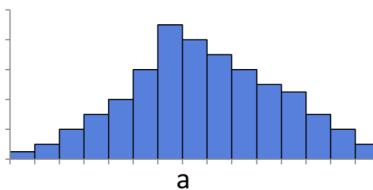
©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- Bimodal: Contains two prevalent modes
- Multimodal: Contains multiple prevalent modes
- Skewed left: Contains a mode on the right with a tail of low-frequency bins on the left
- Skewed right: Contains a mode on the left with a tail of low-frequency bins on the right

**PARTICIPATION ACTIVITY**

## 1.10.3: Common histogram distributions.

The following histograms show common distributions to look for when selecting bin size:



Select the definition that matches each term

1) Histogram a

- Skewed right
- Bimodal
- Skewed left
- Unimodal
- Multimodal

2) Histogram c

- Skewed right
- Bimodal
- Skewed left
- Unimodal
- Multimodal

3) Histogram d

- Skewed right
- Bimodal
- Skewed left
- Unimodal
- Multimodal

4) Histogram e

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- Skewed right
- Bimodal
- Skewed left
- Unimodal
- Multimodal

5) Histogram b

- Skewed right
- Bimodal
- Skewed left
- Unimodal
- Multimodal

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

**Reset**

## Histograms with unevenly-sized bins

Histogram bins are not always equally sized. Ex: Consider the table below containing data for 2014 motor vehicle crash deaths. Most age data is represented in 5 year intervals (green). However, some age intervals are larger or smaller than 5 years (red). Thus, a histogram representing the crash data cannot have equally-sized bins.

Table 1.10.2: Insurance Institute for Highway Safety (IIHS) motor vehicle crash deaths per **100,000** people data, showing unequal interval sizes.

Age	Deaths
0 – 12	872
13 – 15	380
16 – 19	2,243
20 – 24	4,047
25 – 29	3,250
30 – 34	2,567

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

35 – 39	2,155
40 – 44	2,067
45 – 49	2,196
50 – 54	2,712
55 – 59	2,414
60 – 64	1,976
65 – 69	1,517
70 – 74	1,228
75 – 79	1,107
80 – 84	872
85+	985

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Green: 5-year bins  
Red: Non-5-year bins

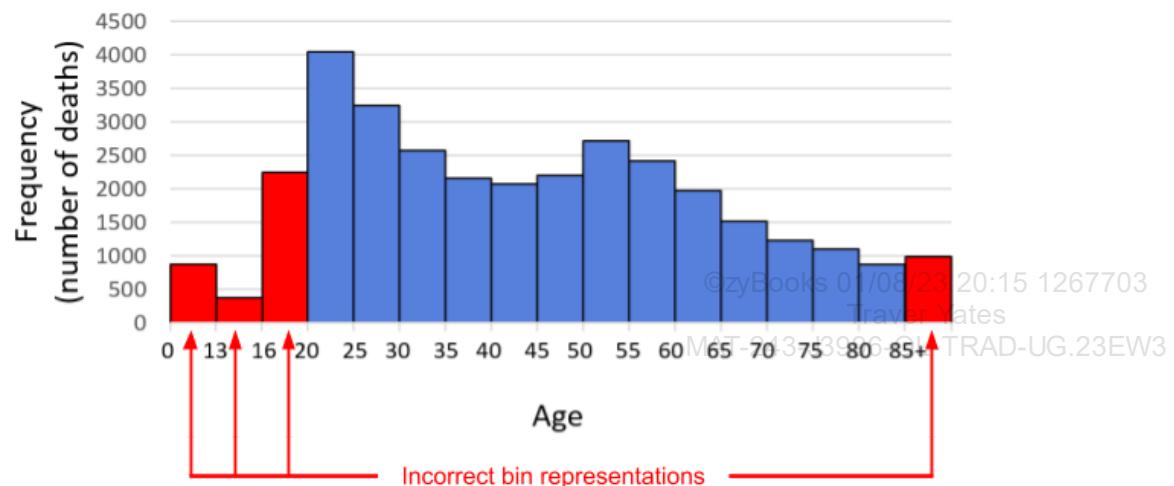
Source: [IIHS fatality facts, 2014](#)<sup>1</sup>

The figure below shows an incorrect histogram for the IIHS crash data with rectangles of equal width, despite different bin sizes. The first red bin represents **13** ages and the second red bin represents **3** ages, while the blue bins represent **5** ages.

By comparing bin heights, the figure's histogram correctly indicates that **20 – 24** year-olds are more likely to die in a crash than **25 – 29** year-olds. Comparing the two bins is reasonable because both bins have the same bin size: **5** years. However, using the same rectangle width to represent bins of different sizes can lead to incorrect conclusions about the likelihood of each bin. Ex: The histogram visually indicates, incorrectly, that a **10** year-old is more likely to die in a crash than a **15** year-old (see correct histogram further below).

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Figure 1.10.2: Incorrect histogram for IIHS crash data, showing 2014 motor vehicle crash deaths per **100,000** people.



To compare likelihoods of two unequally-sized bins, a unit area histogram must be created. A **unit area histogram** has rectangle heights equal to the bin frequency divided by the bin size. The following table shows rectangle heights for a unit area histogram being computed for the IIHS crash data.

Table 1.10.3: IIHS motor vehicle crash deaths per 100,000 people data, showing rectangle heights (deaths per age) being computed for a unit area histogram.

Age	Deaths	Bin size	Deaths per age
0 – 12	872	13	$872/13 = 67$
13 – 15	380	3	127
16 – 19	2,243	4	561
20 – 24	4,047	5	809
25 – 29	3,250	5	650
30 – 34	2,567	5	513
35 – 39	2,155	5	431
40 – 44	2,067	5	413
45 – 49	2,196	5	439

50 – 54	2,712	5	542
55 – 59	2,414	5	483
60 – 64	1,976	5	395
65 – 69	1,517	5	303
70 – 74	1,228	5	246
75 – 79	1,107	5	221
80 – 85	872	5	174
85+	985	15	66

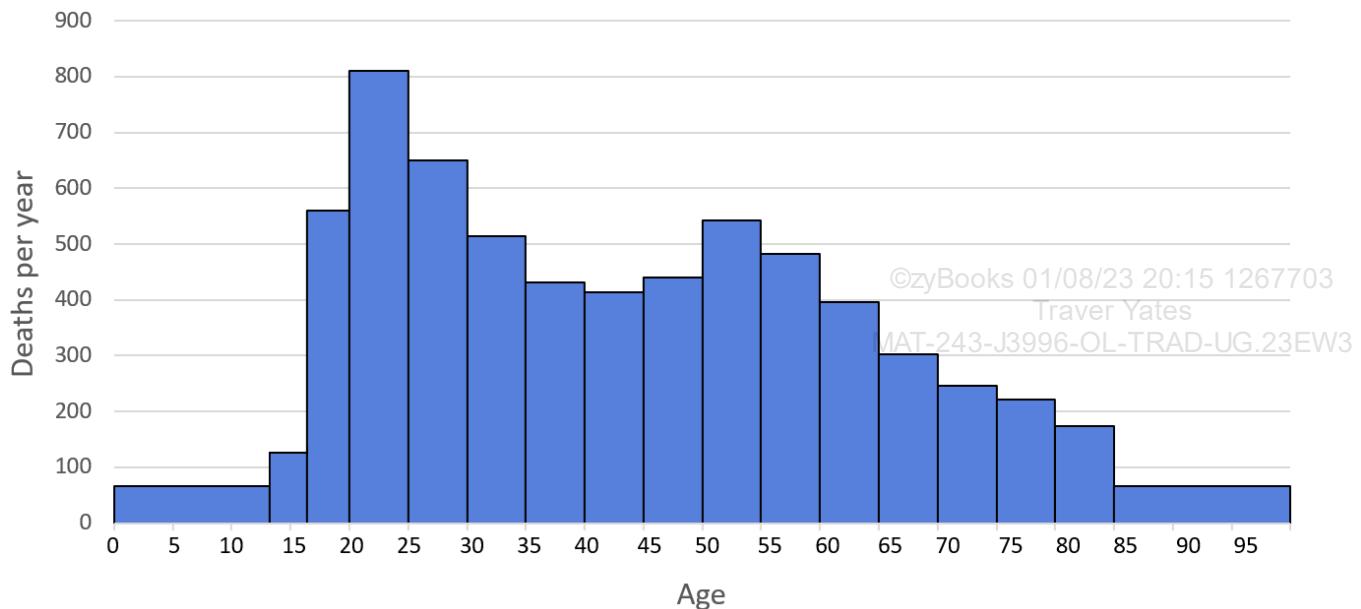
©zyBooks 01/08/23 20:15 1267703  
 Traver Yates  
 MAT-243-J3996-OL-TRAD-UG.23EW3

The figure below shows the unit area histogram for the IIHS crash data. Notice that the y-axis has changed to reflect the Deaths per age metric and that some bins are different sizes. With different bin sizes, bin frequency is determined by rectangle area, instead of rectangle height. Ex: The unit area histogram shows 650 deaths per age for 25 – 29 year-olds:  $(650 \text{ deaths per age}) \times (5 \text{ years in bin}) = 3,250 \text{ deaths}$  for 25 – 29 year-olds.

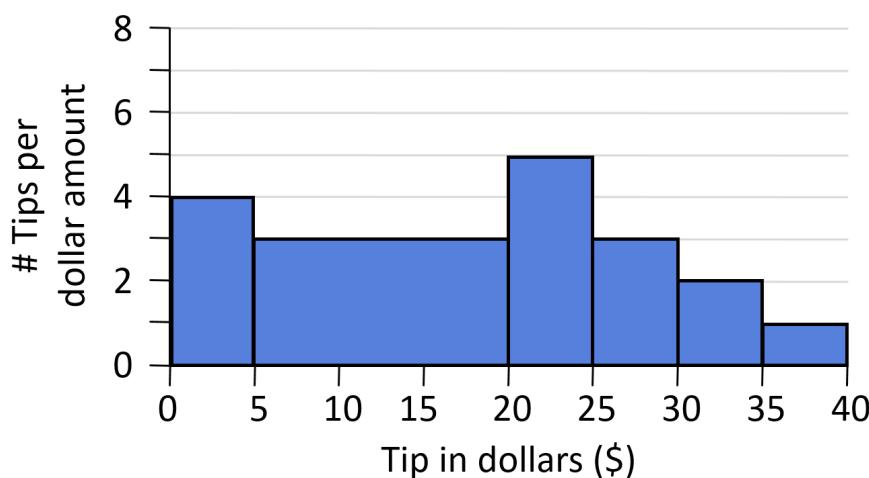
The incorrect histogram above visually suggested that a 10 year-old is more likely to die in a crash than a 15 year-old. Even though the 0 – 12 year-old bin frequency (872) is higher than the 13 – 15 year-old bin frequency (380), the unit area histogram shows that a child between the ages of 0 – 12 years-old is less likely (shorter rectangle) to die in a crash than a 13 – 15 year-old child.

Figure 1.10.3: Correct unit area histogram for IIHS crash data, showing 2014 motor vehicle crash deaths per **100,000** people.

©zyBooks 01/08/23 20:15 1267703  
 Traver Yates  
 MAT-243-J3996-OL-TRAD-UG.23EW3

**PARTICIPATION ACTIVITY****1.10.4: Unit area histogram.**

Consider the following unit area histogram showing tips earned for a waitress:



- 1) The waitress is \_\_\_\_\_ likely to earn a \$5 – \$20 tip than a \$25 – \$30 tip.



- less
- equally
- more

- 2) The waitress is \_\_\_\_\_ likely to earn a \$10 tip than a \$27 tip.



- less

equally more

- 3) The waitress earned a \_\_\_\_\_ tip more often than any other tip.

 \$5 to \$20 \$20 to \$25 \$35 to \$40

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- 4) The waitress is most likely to earn a \_\_\_\_\_ tip.

 \$5 to \$10 \$20 to \$25 \$35 to \$40

## Python-Practice 1.10.1: Histograms.

The [Height](#) dataset gives the heights, in inches, of students in a particular class.

Histograms are plotted using the `hist()` function in the `matplotlib.pyplot` library. The function takes the dataset as the first parameter and the number of bins as the optional second parameter.

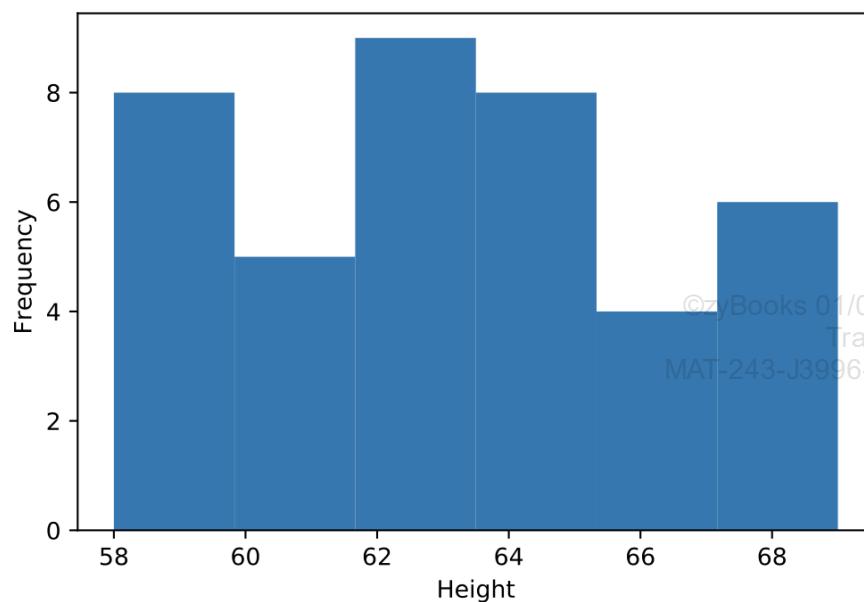
```
import pandas as pd
import matplotlib.pyplot as plt
heights = pd.read_csv('http://data-analytics.zybooks.com/height.csv')

fig, ax = plt.subplots()
plt.hist(heights['Height'], bins=6)
ax.set_xlabel('Height')
ax.set_ylabel('Frequency')
plt.savefig('histogram6.png')
plt.show()
```

©zyBooks 01/08/23 20:15 1267703

Traver Yates

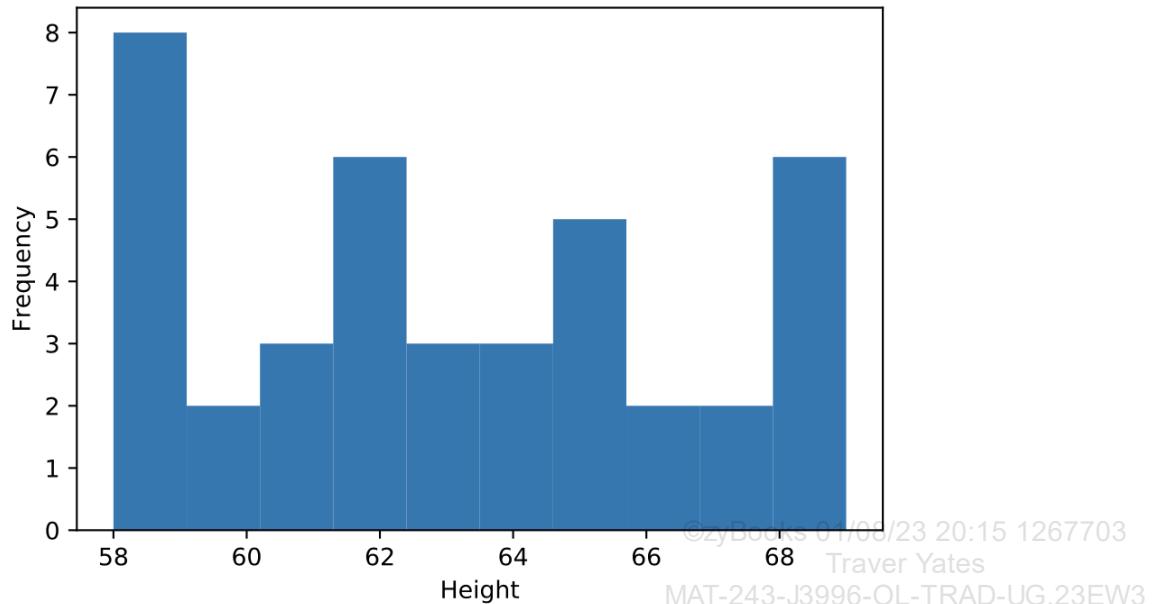
MAT-243-J3996-OL-TRAD-UG.23EW3



The same data can be represented as a histogram with more bins by modifying the bins argument.

```
import pandas as pd
import matplotlib.pyplot as plt
heights = pd.read_csv('http://data-analytics.zybooks.com/height.csv')

fig, ax = plt.subplots()
plt.hist(heights['Height'], bins=10)
ax.set_xlabel('Height')
ax.set_ylabel('Frequency')
plt.savefig('histogram10.png')
plt.show()
```



[Run example](#)

## References

(\*) "General Statistics" Insurance Institute for Highway Safety Highway Loss Data Institute 2014, <http://www.iihs.org/iihs/topics/t/general-statistics/fatalityfacts/overview-of-fatality-facts>.

## 1.11 Survey sampling

©zyBooks 01/08/23 20:15 1267703

### Descriptive and inferential statistics

Two types of statistical analysis exist to describe survey data: *descriptive statistics* and *inferential statistics*.

**Descriptive statistics** focuses on summarizing survey data about a sample drawn from a population. Summary statistics include measures of central tendency such as mean, median, and mode; and dispersion such as range and standard deviation. Descriptive statistics cannot make conclusions based on the data. Rather, descriptive statistics is a way to present data in a meaningful way.

**Inferential statistics** focuses on using information from the sample to make conclusions about the population from which the sample was drawn. The two primary methods of inferential statistics are confidence intervals, which specify the range within which a parameter falls with a given probability, and hypothesis testing, which allows differences between population parameters to be compared.

### Surveys

**Surveys** are conducted to allow statisticians to make generalizations about a population.

A **population** is any collection of objects, people, or things about which statistical inference are made.

A **parameter** of a population is a numerical characteristic of a population, such as mean, median, or standard deviation.

A **sampling unit** is an individual in the population on which a measurement can be taken.

The **sampling frame** is the subset of the population from which a sample is drawn.

The **sample** is composed of the sampling units that provide data to be collected.

A **statistic** is a numerical characteristic of a sample, rather than the population.

The following animation shows the relationship between the population, sampling unit, sampling frame, and sample.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



PARTICIPATION  
ACTIVITY

1.11.1: Sampling a population.

### Animation captions:

1. A sample is a representative subset of a population that is used to measure a parameter of the population.

2. A sampling unit is an individual in the population from which a parameter can be measured.
3. The sampling frame is the subset of the population from which samples can be drawn.
4. The sample is the subset of the sampling frame from which measurements are actually taken.

The following animation shows the relationship between a parameter and a statistic.

**PARTICIPATION ACTIVITY****1.11.2: Parameters and statistics.**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Animation captions:**

1. A preschool's population has **15** aged between **1** and **5**. Sally, aged **5**, is one such child, as is Joey, aged **4**. The dot plot shows the number of children of each age.
2. Sample A consists of **5** children, including Sally but not Joey, selected from the population of **15**.
3. Sample B consists of **5** other children selected from the population, this time including Joey but not Sally.
4. The population mean is a parameter. The mean of each sample is a statistic.

**Example 1.11.1: The Liraglutide Effect and Action in Diabetes: Evaluation of Cardiovascular Outcome Results (LEADER) clinical trial.**

The LEADER clinical trial was initiated in 2010 at 410 hospitals in 32 countries to evaluate the effect of liraglutide, a drug for treatment of type 2 diabetes, on the frequency of cardiovascular diseases such as heart attack, stroke, and heart failure<sup>1</sup>. The populations under study were type 2 diabetes patients with excessively high blood sugar taking either liraglutide or a placebo (an inactive drug). The parameters measured include blood sugar level, kidney function measurement, frequency of adverse effects and complications, and mortality rate. The overall goal of the trial was to determine whether liraglutide treatment was effective for treating type 2 diabetes without increasing the danger of cardiovascular complications.

Identify the sampling unit, sampling frame, and surveys conducted.

**Solution**

- The sampling unit was a type 2 diabetes patient at any of the hospitals at which the trial was conducted.
- The sampling frame was the subset of type 2 diabetes patients who fulfilled the criteria for inclusion into the clinical trial, such as age, cardiovascular disease status, other drugs taken, and whether informed consent to participate in the study was given.

- The surveys included the medical tests that were performed to measure blood sugar level and other health information, as well as observations of the frequency and severity of adverse effects, complications, and deaths.

**PARTICIPATION ACTIVITY**

1.11.3: Distinguishing populations and samples.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



1) An analyst obtains the salaries of all **167** federal appeals court judges in 2014 and computes the mean salary to be **\$211, 200**. Are the **167** judges a population or a sample?

- Population
- Sample

2) An analyst obtains the salaries of all **167** federal appeals court judges in 2014 and computes the mean salary to be **\$211, 200**. Is **\$211, 200** a parameter or a statistic?

- Parameter
- Statistic

3) An analyst surveys **1, 000** registered nurses across the U.S. and computes their mean earnings to be **\$27/hr**. The analyst reports that U.S. nurses have mean earnings of **\$27/hr**<sup>2 3</sup>. Are the **1, 000** nurses a population or a sample?

- Population
- Sample

4) An analyst surveys **1, 000** registered nurses across the U.S. and computes their mean earnings to be **\$27/hr**. The analyst reports that U.S. nurses have mean earnings of **\$27/hr**. Is **\$27/hr** a parameter or a statistic?

- Parameter
- Statistic

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3





5) An analyst is asked to determine how many miles each employee commutes at a **20**-person company. Should the analyst collect data from the population or from a sample?

- Population
- Sample

6) An analyst is asked to determine how many miles each employee commutes at Microsoft, which has over **100, 000** employees<sup>4</sup>. Should the analyst collect data from the population or on a sample?

- Population
- Sample

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



## Bias

In statistics, a **bias** is a difference between the parameter predicted from a survey from the true value of the parameter in the population. Two broad categories of statistical bias include selection bias and response bias.

**Selection bias** exists when the sampling units selected from a population are not representative of the entire population, and are instead biased toward certain subsets of the population. A population should be surveyed in such a way to minimize sampling bias. Several types of selection bias follow.

**Undercoverage** occurs when certain members of a population are inadequately represented in a sample.

**Nonresponse bias** occurs when a sample is biased toward members of a population that participate in a survey.

**Voluntary response bias** occurs when a sample is biased toward members that self-select for participation in a survey.

**Response bias** can result if the responses of survey participants are affected by how a question is asked or the behaviors or attitudes of the participant. Several types of response bias follow.

**Acquiescence bias** occurs when respondents tend to agree with a statement in a survey.

**Extreme responding** occurs when respondents tend to select the most extreme options available.

**Social desirability bias** occurs when respondents tend to answer questions in a way that is socially accepted by others. In other words, a social desirability bias exists when respondents over-report "good" behaviors or under-report "bad" behaviors.

## Example 1.11.2: Types of bias.

For each situation below, determine the most likely type of selection bias.

- a. A website survey
- b. A survey about the frequency of alcohol consumption
- c. In-person survey conducted at a mall in an affluent neighborhood
- d. A survey conducted online by a political news site.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Solution

- a. A comments section of a website soliciting survey responses on a controversial issue in which most of the participants express extreme viewpoints for or against may exhibit voluntary response bias because members of the population who are indisposed to the issue are not adequately represented.
- b. Social desirability bias may exist in a survey about the frequency of alcohol consumption among subsets of the population with differing social attitudes or prohibitions toward alcohol.
- c. Undercoverage. An in-person survey conducted at a shopping center in an affluent neighborhood may inadequately represent members of the population without the economic or transportation means to travel to or shop at the mall.
- d. A survey conducted online by a political news site may have voluntary response bias because people who feel passionately about the topic would be more likely to respond.

PARTICIPATION  
ACTIVITY

1.11.4: Types of bias.



Match each description to the correct type of bias.

Select the definition that matches each term

1) Social desirability bias

- Students conducting a survey on a campus-wide issue only conducted the survey in front of the main engineering building, and the concerns of humanities students on the other side of the campus were not adequately addressed.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- An internet survey sent to a remote rural area with poor internet penetration had only a **3%** response rate.
- A question on a political ballot contains the question "Is Freedom Important?", which is a relatively non-controversial statement with which **98%** of voters agreed.
- The majority of reviews on a restaurant's social media site are either one-star or five-star reviews.
- A survey asking how often respondents send text messages while driving routinely underestimates the actual frequency of texting while driving.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## 2) Undercoverage

- Students conducting a survey on a campus-wide issue only conducted the survey in front of the main engineering building, and the concerns of humanities students on the other side of the campus were not adequately addressed.
- An internet survey sent to a remote rural area with poor internet penetration had only a **3%** response rate.
- A question on a political ballot contains the question "Is Freedom Important?", which is a relatively non-controversial statement with which **98%** of voters agreed.
- The majority of reviews on a restaurant's social media site are either one-star or five-star reviews.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- A survey asking how often respondents send text messages while driving routinely underestimates the actual frequency of texting while driving.

### 3) Nonresponse bias

- Students conducting a survey on a campus-wide issue only conducted the survey in front of the main engineering building, and the concerns of humanities students on the other side of the campus were not adequately addressed.
- An internet survey sent to a remote rural area with poor internet penetration had only a **3%** response rate.
- A question on a political ballot contains the question "Is Freedom Important?", which is a relatively non-controversial statement with which **98%** of voters agreed.
- The majority of reviews on a restaurant's social media site are either one-star or five-star reviews.
- A survey asking how often respondents send text messages while driving routinely underestimates the actual frequency of texting while driving.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

### 4) Extreme responding

- Students conducting a survey on a campus-wide issue only conducted the survey in front of the main engineering building, and the concerns of humanities students on the other side of the campus were not adequately addressed.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- An internet survey sent to a remote rural area with poor internet penetration had only a **3%** response rate.
- A question on a political ballot contains the question "Is Freedom Important?", which is a relatively non-controversial statement with which **98%** of voters agreed.
- The majority of reviews on a restaurant's social media site are either one-star or five-star reviews.
- A survey asking how often respondents send text messages while driving routinely underestimates the actual frequency of texting while driving.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## 5) Acquiescence bias

- Students conducting a survey on a campus-wide issue only conducted the survey in front of the main engineering building, and the concerns of humanities students on the other side of the campus were not adequately addressed.
- An internet survey sent to a remote rural area with poor internet penetration had only a **3%** response rate.
- A question on a political ballot contains the question "Is Freedom Important?", which is a relatively non-controversial statement with which **98%** of voters agreed.
- The majority of reviews on a restaurant's social media site are either one-star or five-star reviews.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- A survey asking how often respondents send text messages while driving routinely underestimates the actual frequency of texting while driving.

**Reset**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Sampling methods

Different sampling methods can help mitigate certain types of statistical bias.

In **simple random sampling**, a sample is constructed by random selection from the population. Mathematically, simple random sampling is a sampling method in which all possible samples consisting of  $n$  units selected from a population of  $N$  units are equally likely.

In **systematic sampling**, every  $k$ th unit from a population of  $N$  units is selected to be in a sample.

In **stratified sampling**, the population is first divided into groups, or strata, depending on some characteristic. Next, samples within each stratum are randomly selected in a proportional manner.

In **cluster sampling**, the population is first divided into groups, or clusters, depending on some characteristic. Next, the sample is constructed by randomly selecting one or more clusters.

In **convenience sampling**, units are drawn from a subset of the population that is readily available.

### Example 1.11.3: Sampling methods.

Determine which sampling method is used in each situation.

- a. A community college contains **3000** students in the School of Arts and Sciences, **1000** students in the School of Engineering, and **1000** students in the School of Performing Arts. From each school, **10%** of students are randomly selected for participation in the survey, for a total of **300** students from the School of Arts and Sciences, **100** students from the School of Engineering, and **100** students from the School of Performing Arts.
- b. Participants for the survey are recruited by flagging down students crossing the main quad of the community college until the necessary number of students have been recruited.
- c. The student body of the community college consists of 1st year students, 2nd year students, 3rd year students, 4th year students, and transfer students. 3rd year students were randomly selected for participation in the survey.
- d. **500** students are randomly selected from the student body of **5000** for participation in the survey.

- e. Participants for the survey are recruited by selecting every 10th name from a list of students at the community college.

## Solution

- a. Stratified sampling
- b. Convenience sampling
- c. Cluster sampling
- d. Simple random sampling
- e. Systematic sampling

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

PARTICIPATION  
ACTIVITY

1.11.5: Sampling methods.



Match each description to the correct sampling method.

Select the definition that matches each term

1) Cluster sampling

- A major polling company in the United States randomly selects Alaska, Illinois, Texas, Florida, and Pennsylvania as the states from which to select households for a survey.
- To select households for the survey, the company consults the property tax rolls in a specific area, and selects households at random for the survey.
- To select households for the survey, the company consults the property tax rolls in a specific area, and subdivides households into several property value brackets. Households are then randomly and proportionally selected from each bracket.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- To select households for the survey, the company consults the property tax rolls in a specific area, and selects every household with an odd street number for the survey.

## 2) Systematic sampling

- A major polling company in the United States randomly selects Alaska, Illinois, Texas, Florida, and Pennsylvania as the states from which to select households for a survey.
- To select households for the survey, the company consults the property tax rolls in a specific area, and selects households at random for the survey.
- To select households for the survey, the company consults the property tax rolls in a specific area, and subdivides households into several property value brackets. Households are then randomly and proportionally selected from each bracket.
- To select households for the survey, the company consults the property tax rolls in a specific area, and selects every household with an odd street number for the survey.

## 3) Simple random sampling

- A major polling company in the United States randomly selects Alaska, Illinois, Texas, Florida, and Pennsylvania as the states from which to select households for a survey.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- To select households for the survey, the company consults the property tax rolls in a specific area, and selects households at random for the survey.
- To select households for the survey, the company consults the property tax rolls in a specific area, and subdivides households into several property value brackets. Households are then randomly and proportionally selected from each bracket.
- To select households for the survey, the company consults the property tax rolls in a specific area, and selects every household with an odd street number for the survey.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

#### 4) Stratified sampling

- A major polling company in the United States randomly selects Alaska, Illinois, Texas, Florida, and Pennsylvania as the states from which to select households for a survey.
- To select households for the survey, the company consults the property tax rolls in a specific area, and selects households at random for the survey.
- To select households for the survey, the company consults the property tax rolls in a specific area, and subdivides households into several property value brackets. Households are then randomly and proportionally selected from each bracket.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- To select households for the survey, the company consults the property tax rolls in a specific area, and selects every household with an odd street number for the survey.

**Reset**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## References

- (\*) Marso, Steven P., et al. "Liraglutide and Cardiovascular Outcomes in Type 2 Diabetes." *The New England Journal of Medicine*, 375:311-322, 28 July 2016, DOI: 10.1056/NEJMoa1603827
- (\*) "Providers and Service Use Indicators NURSES AND PHYSICIAN ASSISTANTS." *Henry J Kaiser Family Foundation*, kff.org/other/state-indicator/total-registered-nurses/
- (\*) "Registered Nurse (RN) Salary." *Payscale.com*, 2016, www.payscale.com/research/US/Job=Registered\_Nurse\_(RN)/Hourly\_Rate, 2016 data
- (\*) "Facts About Microsoft." *Microsoft*, 2015, http://news.microsoft.com/facts-about-microsoft/#sm.0001d1d65h4x8e63v142iztf9potu

## 1.12 Measures of center

### The mean

Large amounts of data can be overwhelming. A single number can summarize information about a dataset, such as the central tendency or the dispersion of the dataset. A common data summary is the **arithmetic mean** or **mean**, which is the sum of the data values in a dataset divided by the number of values in the dataset.

Mathematically, the mean of a set of  $n$  data values  $x_1, x_2, \dots, x_n$  is denoted  $\bar{x}$  and is defined as follows.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \dots + x_n}{n}$$

A similar quantity called the **weighted mean** is often calculated in addition to the mean. The **weighted mean** is a measure of center where some values are counted more than once. Weights are often expressed either as positive integers or percentages. Mathematically, the weighted mean of a set  $x_1, x_2, \dots, x_n$  with corresponding weights  $a_1, a_2, \dots, a_n$  is defined as follows.

$$\bar{x} = \frac{\sum_{i=1}^n a_i x_i}{\sum_{i=1}^n a_i} = \frac{a_1 x_1 + a_2 x_2 + \dots + a_n x_n}{a_1 + a_2 + \dots + a_n}$$

The following animation shows the relationship of the mean to the values in a dataset.

**PARTICIPATION ACTIVITY**

1.12.1: The mean.



©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Animation content:**

undefined

**Animation captions:**

1. The mean summarizes data, computed as the sum divided by the number of values.
2. Graphically, the mean is a value that balances the data values.

**Example 1.12.1: Finding weighted means.**

Find the weighted mean of the following dataset.

Data	Weight
6	2
8	2
17	1

**Solution**

Data with weights that are greater than 1 are counted more than once, which means that the dataset above is the same as 6, 6, 8, 8, 17. Using the formula, the weighted mean is

$$\bar{x} = \frac{\sum_{i=1}^n a_i x_i}{\sum_{i=1}^n a_i} = \frac{2 \cdot 6 + 2 \cdot 8 + 1 \cdot 17}{2 + 2 + 1} \stackrel{\text{©zyBooks 01/08/23 20:15 1267703}}{=} \frac{45}{5} = 9 \stackrel{\text{Traver Yates}}{} \\ \stackrel{\text{MAT-243-J3996-OL-TRAD-UG.23EW3}}{}$$

**PARTICIPATION ACTIVITY**

1.12.2: The mean.





1) What is the mean of the dataset

12, 1, 2?

**Check****Show answer**

2) What is the mean of the dataset

2, 6, 4?

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



3) What is the mean of the dataset

2, 3, 4, 1? Type as: #.#

**Check****Show answer**

4) What is the mean of the dataset

-3, 15?

**Check****Show answer**

5) What is the weighted mean of the

dataset **20, 4, 1** if the weights are  
1, 2, and 2 respectively?

**Check****Show answer**

6) What is the weighted mean of the

dataset **3, 0, 37** if the weights are  
2, 4, and 2 respectively?

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



## Python-Function 1.12.1: mean().

The DataFrame.mean() function is used to find the mean.

```
import pandas as pd

scores = pd.read_csv('http://data-
analytics.zybooks.com/ExamScores.csv')

# Prints the first few lines of data
print(scores.head())

# Prints the mean for each exam
print(scores.mean())

# Prints the mean for Exam1 only
print(scores[['Exam1']].mean())

# Prints the means for Exam1 and Exam2
print(scores[['Exam1', 'Exam2']].mean())
```

	Exam1	Exam2	Exam3
Exam4	0008/23 66:15 174	7703	
68	Traver Yates		
11	76	84	100
75			
2	89	80	100
95			
3	83	76	91
74			
4	86	99	78
78			
	Exam1	82.70	
	Exam2	79.40	
	Exam3	73.34	
	Exam4	76.50	
	dtype: float64		
	Exam1	82.7	
	dtype: float64		
	Exam1	82.7	
	Exam2	79.4	
	dtype: float64		

[Run example](#)

## The median

The **median** is the middle value in a sorted dataset. To find the median of a dataset, the dataset must first be sorted in ascending or descending order. The method of finding the median depends on whether the number of data values  $n$  is even or odd.

- If  $n$  is odd, the median is the middle value of the sorted dataset. Specifically, the median is the  $\left(\frac{n+1}{2}\right)$ th value.
- If  $n$  is even, the median is the mean of the middle two values of the sorted dataset. Specifically, the median is the mean of the  $\left(\frac{n}{2}\right)$ th and  $\left(\frac{n}{2} + 1\right)$ th values.

### Example 1.12.2: The median.

Find the median of each of the following datasets.

- 10, 20, 20, 30, 60, 60, 80
- 99, 80, 60, 60, 30, 20, 20, 10
- 100, 3, 6, 9, 2

**Solution**

a. The dataset contains 7 values and is already sorted in ascending order. Thus, the median is the  $\frac{7+1}{2} = 4$ th data value, which is 30.

b. The dataset contains 8 values and is already sorted in descending order. Thus, the median is the mean of the  $\frac{8}{2} = 4$ th data value and the  $\frac{8+1}{2} = 5$ th data value, which is  $\frac{60+30}{2} = \frac{90}{2} = 45$ .

c. The dataset must first be arranged in ascending or descending order to find the median. In ascending order, the dataset is 2, 3, 6, 9, 100. Thus, the median is the  $\frac{5+1}{2} = 3$ rd data value, which is 6.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

### Python-Function 1.12.2: median().

The DataFrame.median() function is used to find the median.

```
import pandas as pd

# Loads the ExamScores dataset
scores = pd.read_csv('http://data-
analytics.zybooks.com/ExamScores.csv')

# Prints the median for each exam
print(scores.median())

# Prints the median for Exam1 only
print(scores[['Exam1']].median())
```

Exam1	83.0
Exam2	79.5
Exam3	74.5
Exam4	75.0
dtype:	float64
Exam1	83.0
dtype:	float64

[Run example](#)

**PARTICIPATION ACTIVITY**

### 1.12.3: The median.



- 1) What is the median of the dataset  
1, 4, 5, 9, 11?

- 6
- 5

- 2) What is the median of the dataset  
4, 3, 2, 6, 7?

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



2 4

3) What is the median of the dataset

2, 3, 5, 18?

 3 4 5 7

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

4) What is the median of the dataset

-1, -5, -3, 6, 7?

 -3 -5 -1

## Outliers

The dataset **100, 3, 6, 9, 2** in a previous example illustrates an advantage of the median over the mean. The data value **100** is much larger than the other data values. Such a value is an **outlier**, or a data value that is either much greater than or much less than the rest of the data and not representative of the rest of the data being considered. Compared to the median of **6**, the mean is

$$\frac{2 + 3 + 6 + 9 + 100}{5} = \frac{120}{5} = 24, \text{ which is much larger than } 6 \text{ due to influence by the outlier of } 100.$$

As a practical example, the net worth (in USD) of 5 particular individuals in Medina, Washington in 2015 was **\$300,000, \$400,000, \$250,000, \$80,000,000,000**, and **\$600,000**. The outlier of **\$80,000,000,000** is due to Bill Gates, a co-founder of Microsoft, living in Medina. The mean net worth of **\$16,000,000,000** is thus a poor data summary because the mean suggests that all 5 people are wealthy multimillionaires.

The following animation shows the relationship of the median to the mean and the values in a dataset, including outliers.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**PARTICIPATION ACTIVITY**

1.12.4: The median.

### Animation captions:

1. To find the median, the data must first be sorted.
2. The median is the middle value among the sorted values.

3. The value of an outlier affects the mean but not the median.

### Example 1.12.3: Pensions in San Diego.

In recent years, various scandals have been reported relating to exorbitant pensions that city employees approved for themselves, sometimes resulting in city staff later being found guilty of crimes and imprisoned. The following table summarizes data for a major city.

In 2011, the mean pension among the 10 highest earning city employees was **\$239, 940** (the median was **\$231, 922**).

Pensions are paid for the pensioners' remaining life, often 30 years or more. The following table lists the pension amounts and the person's last job position, which the above mean and median summarize.

Last job position	Pension amount
Assistant City Attorney	\$307, 758
Investment Officer	\$255, 509
Fire Battalion Chief	\$244, 435
Assistant Police Chief	\$242, 947
City Librarian	\$234, 091
Fire Chief	\$229, 753
Fire Battalion Chief	\$228, 392
Deputy City Attorney	\$224, 863
Fire Battalion Chief	\$217, 649
Assistant Water Department Director	\$214, 007

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

### Example 1.12.4: Misleading lawyer ads.

Law firms sometimes report mean settlement amounts rather than the more appropriate median. The idea is to lure clients into believing they may win a larger settlement than is

actually likely, since the mean is influenced by a few big wins. The following provides two summaries of data for a law firm whose win data is shown below.

Let us represent you: Our mean settlement amount is \$1,529,000.

Let us represent you: Our median settlement amount is \$675,000.

Examining the table below, one can see that the above mean is misleading, due to being unduly influenced by the top 2-3 settlements. The data summary using the median more appropriately reflects what a client might expect to win.

\$7,500,000	Group settlement for motor vehicle accident
\$5,700,000	Federal verdict against Veteran's Administration for not detecting neurological condition and performing surgery before woman became paralyzed
\$4,750,000	Medical malpractice involving undiagnosed kidney failure prior to baby delivery
\$3,000,000	Nursing home neglect leading to bedsores
\$1,400,000	Nursing home neglect
\$1,000,000	Nursing home neglect leading to a fall
\$950,000	Nursing home neglect leading to death
\$830,000	Motor vehicle accident resulting in death
\$750,000	Nursing home neglect leading to bedsores
\$700,000	Medical malpractice involving death due to overdose
\$650,000	Nursing home neglect leading to infection
\$600,000	Product liability case involving leg injury from an ATV
\$550,000	Client struck by a limousine
\$500,000	Nursing home neglect resulting in death
\$400,000	Medical malpractice involving pressure sores
\$375,000	Medical malpractice involving pressure sores in hospital
\$300,000	Motor vehicle collision

\$230,000	Client struck in the arm by a stray bullet while driving
\$215,000	Client suffered neck injury from a motor vehicle accident
\$194,000	Client injured knee falling through a restaurant trap door

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Example 1.12.5: Mean and median age of marriage.

The median age at first marriage in a certain country increased from **26.8** for men and **25.1** for women in 2000 to **28.2** for men and **26.1** for women in 2010. In this case, the data is likely skewed to the right. People do not usually marry much earlier than **18**, but can marry as old as **50** or later. The outliers here are not a few high numbers, but rather a skewing of data that increases the mean. Thus, the median is more often reported compared to the mean.

#### PARTICIPATION ACTIVITY

##### 1.12.5: Mean and median.



- 1) The pension data above can be summarized nearly equally well using either the mean or the median.

- True
- False



- 2) The settlement data above can be summarized nearly equally well using either the mean or the median.

- True
- False



- 3) In the marriage example above, the data presenter likely chose median because the presenter is trying to lure people into believing marriage age is younger than justified by the data.

- True
- False

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3





4) Law firms are not the only companies that report the mean when the median would be more appropriate.

- True
- False

5) Home prices for a given city or state are commonly summarized in news articles using the median price.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- True
- False

## The mode

The **mode** is the most frequently-occurring value in a dataset and is another measure of center. A dataset may have multiple modes if multiple values have the same maximum frequency. A dataset with only unique values does not have a mode.

### Example 1.12.6: The mode.

Find the mode or modes of each dataset.

- a. 1, 2, 2, 2, 3, 3, 4
- b. 1, 2, 2, 3, 3, 4
- c. 1, 2, 3, 4

### Solution

- a. Since **2** is the value with the highest frequency (**3**), the mode is **2**.
- b. Since both **2** and **3** have the highest frequency (**2**), the dataset has two modes, **2** and **3**.
- c. Since every value in the dataset is unique, the dataset does not have a mode.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

### Python-Function 1.12.3: mode().

The DataFrame.mode() function is used to find the mode. Multiple modes are shown on separate lines. "NaN" (not a number) is used to fill in blank spaces for columns with fewer modes.

```

import pandas as pd

# Loads the ExamScores dataset
scores = pd.read_csv('http://data-
analytics.zybooks.com/ExamScores.csv')

# Prints the mode(s) for each exam
print(scores.mode())

# Prints the mode(s) for Exam1 only
print(scores[['Exam1']].mode())

```

	Exam1	Exam2	Exam3
Exam4			
0	78	84.0	100.0
74.0			
1	83	100.0	NaN
75.0			
2	89	NaN	NaN
NaN			
Exam1			
0	78		
1	83		
2	89		

©zyBooks 01/08/23 20:15 1267703

MAT-243-J3996-OL-TRAD-UG.23EW3

[Run example](#)**PARTICIPATION ACTIVITY**

1.12.6: The mode.



1) What is the mode of the dataset



1, 4, 4, 5, 5, 9, 9, 9?

- 1
- 5
- 9

2) Which of the following statements is true about the dataset 2, 5, 6, 7?



- No mode exists
- All values in the dataset are modes.

3) After a baseball tournament, the number of runs scored by each player is 0, 0, 0, 0, 0, 0, 0, 3, 3, 6, 10. What is the mode?



- 0
- 7

©zyBooks 01/08/23 20:15 1267703

Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## Python-Practice 1.12.1: Mean and median.

The [rent](#) dataset lists the monthly rent (in USD) of 10 apartments with area less than 1000 ft<sup>2</sup> in 6 cities ([Source](#)). The columns of the dataset give the rents in Santa Monica, CA, Boise, ID, Tucson, AZ, Detroit, MI, Pittsburgh, PA, and Orlando, FL.

To find the mean and median, the CSV file is first imported into Python as a pandas DataFrame.

```
import pandas as pd
rent = pd.read_csv('http://data-analytics.zybooks.com/rent.csv')
print(rent)
```

	Santa Monica CA	Boise ID	Tucson AZ	Detroit MI	Pittsburgh PA	Orlando FL	Traver Yate	MAT-243-J3996-OL-TRAD-UG.23EW3
0	10230	1600	2495	3195	2480	2242		
1	10000	1500	2200	2695	2435	2000		
2	9000	1029	2150	2595	2405	1912		
3	8500	1025	1800	2495	2350	1895		
4	8250	980	1650	2495	2320	1800		
5	8000	950	1600	1675	2316	1765		
6	7000	950	1500	1525	2305	1685		
7	6500	925	1500	1480	2290	1670		
8	6000	925	1500	1410	2275	1665		
9	5815	900	1500	1400	2265	1625		

The mean rent for Santa Monica is found as follows.

```
print(rent['Santa Monica CA'].mean())
```

7929.5

The mean rent for all cities is found as follows.

```
print(rent.mean())
```

Santa Monica CA	7929.5
Boise ID	1078.4
Tucson AZ	1789.5
Detroit MI	2096.5
Pittsburgh PA	2344.1
Orlando FL	1825.9
dtype: float64	

The median rent for all cities is found as follows.

```
print(rent.median())
```

Santa Monica CA	8125.0
Boise ID	965.0
Tucson AZ	1625.0
Detroit MI	2085.0
Pittsburgh PA	2318.0
Orlando FL	1782.5
dtype: float64	

[Run example](#)

## 1.13 Measures of variability

### Variance and standard deviation

**Variability** is the difference between values in a dataset and the center of the dataset. A measure of center alone does not indicate the extent of variability. Ex: the data set **1, 2, 8, and 9** and the data set **4, 5, 5, and 6** both have a mean of **5** and a median of **5**. However, the values in the first dataset have a larger variability. Two common measures of variability are **variance** and **standard deviation**.

**Variance**, is the average of the square difference from the mean. **Standard deviation** is the square root of the variance. By definition, the variance is the square of the standard deviation.

The formula for variance and standard deviation depends on whether the dataset contains the whole population or or a subset of the population. The sample standard deviation is denoted by  $s$ , while the population standard deviation is denoted by  $\sigma$ . The formulas are given below.

$$\sigma^2 = \frac{\sum_{i=1}^n (a_i - \mu)^2}{n} \quad (\text{Population variance})$$

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (a_i - \mu)^2}{n}} \quad (\text{Population standard deviation})$$

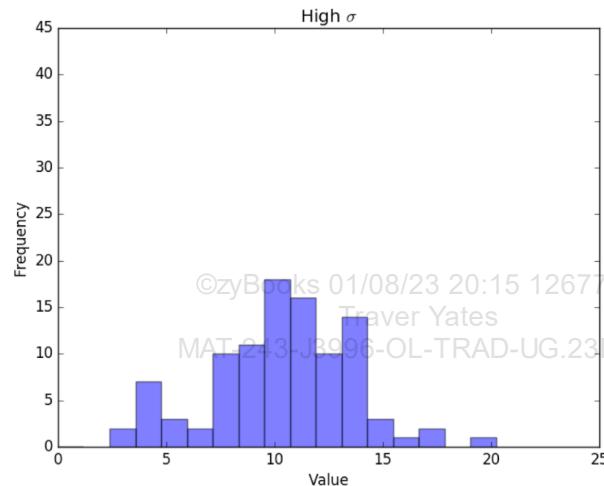
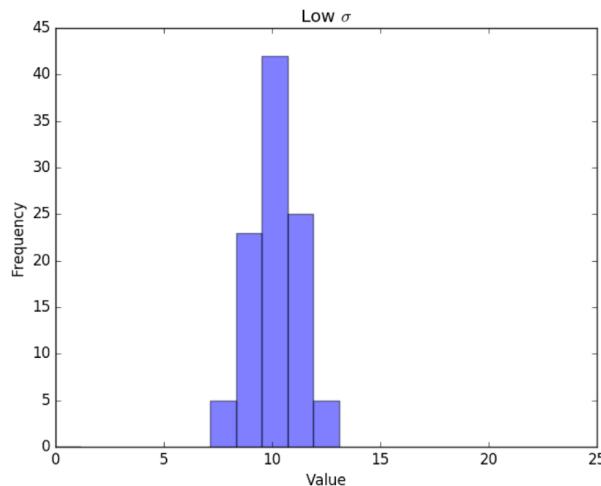
$$s^2 = \frac{\sum_{i=1}^n (a_i - \bar{x})^2}{n-1} \quad (\text{Sample variance})$$

$$s = \sqrt{\frac{\sum_{i=1}^n (a_i - \bar{x})^2}{n-1}} \quad (\text{Sample standard deviation})$$

In the formulas above,  $n$  is the number of data values,  $a_1, a_2, \dots, a_n$  are the  $n$  data values,  $\bar{x}$  is the sample mean, and  $\mu$  is the population mean. The numerator of the fraction in both variance formulas is often referred to as the sum of the square differences. A large standard deviation indicates a more spread-out data set. A small standard deviation indicates a more tightly clustered data set.

The following figure shows histograms of data with a low standard deviation and a high standard deviation.

Figure 1.13.1: Data sets with low and high standard deviations.



### Example 1.13.1: Finding the variance and standard deviation.

Find the variance and standard deviation of the dataset: 4, 5, 6, 13

- assuming that the dataset is a subset of the population
- assuming that the dataset represents the whole population

#### Solution

- Since the dataset is a subset of measurements from a population, the sample variance and sample standard deviation are obtained.

First, the sample mean is found.

$$\bar{x} = \frac{4 + 5 + 6 + 13}{4} = \frac{28}{4} = 7$$

Using a table, the sum of the square differences can be obtained.

$a_i$	$\bar{x}$	$a_i - \bar{x}$	$(a_i - \bar{x})^2$
4	7	-3	9
5	7	-2	4
6	7	-1	1
13	7	6	36

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

From the table above, the sum of the square differences is

$$\sum_{i=1}^4 (a_i - \bar{x})^2 = 9 + 4 + 1 + 36 = 50$$

Thus, the sample variance and sample standard deviation are

$$s^2 = \frac{\sum_{i=1}^4 (a_i - \bar{x})^2}{4 - 1} = \frac{50}{3} \approx 16.667$$

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

$$s = \sqrt{\frac{50}{3}} \approx 4.082$$

- b. Since the dataset represents the entire population, the population variance and standard deviation are obtained. The population mean is  $\mu = 7$  and the sum of the square differences is  $\sum_{i=1}^4 (a_i - \mu)^2 = 50$ , using the same set of calculations as shown above.

Thus, the population variance and population standard deviation are

$$\sigma^2 = \frac{\sum_{i=1}^4 (a_i - \mu)^2}{4} = \frac{50}{4} = 12.5$$

$$\sigma = \sqrt{\frac{\sum_{i=1}^4 (a_i - \mu)^2}{4}} \approx 3.536$$

## Analysis

Although the population mean and sample mean are the same in this example,  $\bar{x}$  and  $\mu$  are generally different. In most cases,  $\mu$  is unknown or difficult to calculate. The sample mean can be used to estimate the population mean, but the sample mean is strongly susceptible to the presence of extreme values.

The sample variance and standard deviation are always greater than the population variance and standard deviation, because the value depends strongly on the elements of the subset taken during sampling. Thus, samples display greater variability than the entire population.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Python-Function 1.13.1: std() and var().

The DataFrame.std() function is used to find the standard deviation, and the DataFrame.var() is used to find the variance.

```

import pandas as pd

# Loads the ExamScores dataset
scores = pd.read_csv('http://data-
analytics.zybooks.com/ExamScores.csv')

# Prints the standard deviation for each exam
print(scores.std())

# Prints the standard deviation for Exam1 only
print(scores[['Exam1']].std())

# Prints the variance for each exam
print(scores.var())

# Prints the variance for Exam1 only
print(scores[['Exam1']].var())

```

©zyBooks 01/08/23 20:15 1267703  
MAT-243-J3996-OL-TRAD-UG.23EW3

Exam1	9.291756
Exam2	14.332780
Exam3	21.754296
Exam4	8.056560
	dtype: float64
Exam1	9.291756
	dtype: float64
Exam1	86.336735
Exam2	205.428571
Exam3	473.249388
Exam4	64.908163
	dtype: float64
Exam1	86.336735
	dtype: float64

[Run example](#)

PARTICIPATION ACTIVITY

1.13.1: Variance and standard deviation.



Consider the dataset 1, 2, 4, 5 taken from a subset of a population.

- 1) What is the sample mean?




**Check**

**Show answer**

- 2) What is the sum of the squares of the differences between each data value and the sample mean?




**Check**

**Show answer**

- 3) What is the sample variance? Type as: #.###




**Check**

**Show answer**



- 4) What is the sample standard deviation? Type as: #.###

**Check****Show answer**

- 5) Suppose the data represents measures from the entire population. What is the population variance? Type as: #.#

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



- 6) Suppose the data represents measures from the entire population. What is the population standard deviation? Type as: #.###

**Check****Show answer**

### Example 1.13.2: Course student evaluation data summary.

Below is a data summary of student evaluations for a particular course at a major university. The summary includes and compares three sets of data: the course (and professor), all courses within that course's department, and all courses at the university. For each set, the summary provides the mean (Mean), median (Med), and standard deviation (SD). Items to note:

- The median conveys little useful information for this data (e.g., being 5.0 for nearly all questions under Course); the mean is clearly superior. ©zyBooks 01/08/23 20:15 1267703 Traver Yates
- Per the standard deviation, this course professor's scores for question 13 (the main question of interest for professors) have less variation (0.5) than for the department (0.8) or university (0.9), indicating students were more consistent in rating this professor (highly).
- The counts per rating category are provided, to provide further insight into how the ratings were distributed.

- "Percentiles" (shown as "% tile") are also indicated showing how the course ranked compared to department or university courses. Ex: For question 19, this course was rated higher than 83% of all courses at the university (in other words, in the top 17% of all courses).

Questions	Course										Department				Campus							
	High					Low					% tile		Mean		SD		% tile		Mean		SD	
	5	4	3	2	1	N/A	Mean	Med	SD													
1 I had a strong desire to take this course	38	22	7	1	2	-	4.3	5.0	0.9	73	4.0	4.0	1.0	72	4.0	4.0	1.0					
2 I attended class regularly	54	13	1	2	-	-	4.7	5.0	0.6	83	4.2	5.0	1.0	83	4.4	5.0	0.9					
3 I put considerable effort into this course	46	17	2	5	-	-	4.5	5.0	0.9	73	4.2	4.0	0.8	75	4.3	4.0	0.8					
4 I gained a good understanding of the course content	45	21	2	1	-	-	4.6	5.0	0.6	91	4.1	4.0	0.8	84	4.2	4.0	0.9					
5 I normally spent at least two hours preparing for each hour of class	34	23	11	2	-	-	4.3	4.0	0.8	79	3.8	4.0	1.1	75	3.9	4.0	1.1					
6 Instructor was prepared and organized	54	15	1	-	-	-	4.8	5.0	0.5	92	4.4	4.0	0.8	92	4.3	5.0	0.9					
7 Instructor used class time effectively	45	22	3	-	-	-	4.6	5.0	0.6	75	4.3	4.0	0.8	83	4.3	4.0	0.9					
8 Instructor was clear and understandable	55	14	1	-	-	-	4.8	5.0	0.5	100	4.3	4.0	0.8	93	4.2	4.0	1.0					
9 Instructor exhibited enthusiasm for subject and teaching	59	10	1	-	-	-	4.8	5.0	0.4	91	4.4	5.0	0.8	90	4.4	5.0	0.8					
10 Instructor respected students; sensitive to and concerned with their progress	55	14	1	-	-	-	4.8	5.0	0.5	93	4.3	4.0	0.9	92	4.3	5.0	0.9					
11 Instructor was available and helpful	47	20	3	-	-	-	4.6	5.0	0.6	77	4.3	4.0	0.8	83	4.3	4.0	0.9					
12 Instructor was fair in evaluating students	51	16	3	-	-	-	4.7	5.0	0.6	100	4.3	4.0	0.8	88	4.3	4.0	0.9					
13 Instructor was effective as a teacher overall	51	17	1	-	-	-	4.7	5.0	0.5	100	4.3	4.0	0.8	88	4.3	4.0	0.9					
14 The syllabus clearly explained the structure of the courses	49	19	2	-	-	-	4.7	5.0	0.5	85	4.3	4.0	0.8	87	4.4	5.0	0.8					
15 The examinations reflected the materials covered during the course	50	17	3	-	-	-	4.7	5.0	0.6	100	4.3	4.0	0.8	87	4.3	4.0	0.9					
16 The required readings contributed to my learning	44	22	4	-	-	-	4.6	5.0	0.6	91	4.2	4.0	0.9	80	4.2	4.0	0.9					
17 The assignments contributed to my learning	48	21	1	-	-	-	4.7	5.0	0.5	100	4.3	4.0	0.9	86	4.3	4.0	0.9					
18 Supplementary materials (e.g. films, slides, videos, demonstrations, guest lectures, iLearn, web pages, etc) were informative	41	26	2	-	-	-	4.6	5.0	0.6	100	4.2	4.0	0.9	83	4.2	4.0	0.9					
19 The course overall as a learning experience was excellent	48	19	2	1	-	-	4.6	5.0	0.6	85	4.2	4.0	0.9	83	4.2	4.0	0.9					

## Python-Practice 1.13.1: Standard deviation and variance.

The [rent](#) dataset is first imported.

```
import pandas as pd
rent = pd.read_csv('http://data-analytics.zybooks.com/rent.csv')
print(rent)
```

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

	Santa Monica CA	Boise ID	Tucson AZ	Detroit MI	Pittsburgh PA
Orlando FL	2242	10230	1600	2495	3195
1	2000	10000	1500	2200	2695
2	1912	9000	1029	2150	2595
3	1895	8500	1025	1800	2495
4	1800	8250	980	1650	2495
5	1765	8000	950	1600	1675
6	1685	7000	950	1500	1525
7	1670	6500	925	1500	1480
8	1665	6000	925	1500	1410
9	1625	5815	900	1500	1400
					2265

The standard deviation in rent for all cities is found as follows.

```
print(rent.std())
```

Santa Monica CA	1572.680744
Boise ID	253.125002
Tucson AZ	363.153243
Detroit MI	664.308538
Pittsburgh PA	72.520419
Orlando FL	191.608658
<b>dtype:</b>	<b>float64</b>

The variance in rent for all cities is found as follows.

```
print(rent.var())
```

Santa Monica CA	2.473325e+06
Boise ID	6.407227e+04
Tucson AZ	1.318803e+05
Detroit MI	4.413058e+05
Pittsburgh PA	5.259211e+03
Orlando FL	3.671388e+04
<b>dtype:</b>	<b>float64</b>

[Run example](#)

## Mean absolute deviation

**Mean absolute deviation (MAD)** is the mean of the absolute difference between each value and the mean of the values. The MAD uses the absolute value instead of the square root of a sum of squares to avoid negative distances. The formula for the MAD is given below.

$$\frac{\sum_{i=1}^n |a_i - \bar{x}|}{n} \quad (\text{Mean absolute deviation})$$



## Animation captions:

1. A measure of center alone does not indicate the extent of the variability of the data values.
2. The mean absolute deviation is the mean of the distances of the data values from the mean of the data values.
3. Data has the same mean but a smaller mean absolute deviation if the data points have a smaller variability.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates

## Python-Function 1.13.2: mad().

The DataFrame.mad() function is used to find the mean absolute deviation.

```
import pandas as pd

# Loads the ExamScores dataset
scores = pd.read_csv('http://data-
analytics.zybooks.com/ExamScores.csv')

# Prints the mean absolute deviation of all scores
print(scores.mad())

# Prints the mean absolute deviation of Exam 1
print(scores['Exam1'].mad())
```

```
Exam1
7.1360
Exam2
12.1200
Exam3
17.7928
Exam4
5.7000
dtype: float64
7.136
```

[Run example](#)



Use the dataset 1, 2, 4, 5 to answer the following.

- 1) What is the mean?

  
//

**Check**

**Show answer**

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- 2) What is the sum of the absolute differences of each data value and the mean?

  
//

**Check****Show answer**

- 3) What is the mean absolute deviation of the data values? Type as: #.#

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Check****Show answer**

## 1.14 Introduction to probability

### Experiments, outcomes, and sample spaces

An **experiment** is a procedure that results in one out of a number of possible outcomes. An **outcome** is the result of an experiment. Ex: Rolling a six-sided die is an experiment, and the number of dots displayed is an outcome. The set of all possible outcomes is called the **sample space** of the experiment and is denoted  $S$ . Ex: The sample space of a six-sided die roll is  $S = \{1, 2, 3, 4, 5, 6\}$ .

**PARTICIPATION ACTIVITY**

1.14.1: Experiments, outcomes, and sample spaces.



#### Animation captions:

1. Rolling a six-sided die is an experiment that results in 1, 2, 3, 4, 5, or 6.
2. A roll of 1 is an experiment and the result 1 is an outcome. A different roll of 3 is an experiment, and the result 3 is an outcome.
3. The sample space is the set of all possible outcomes.

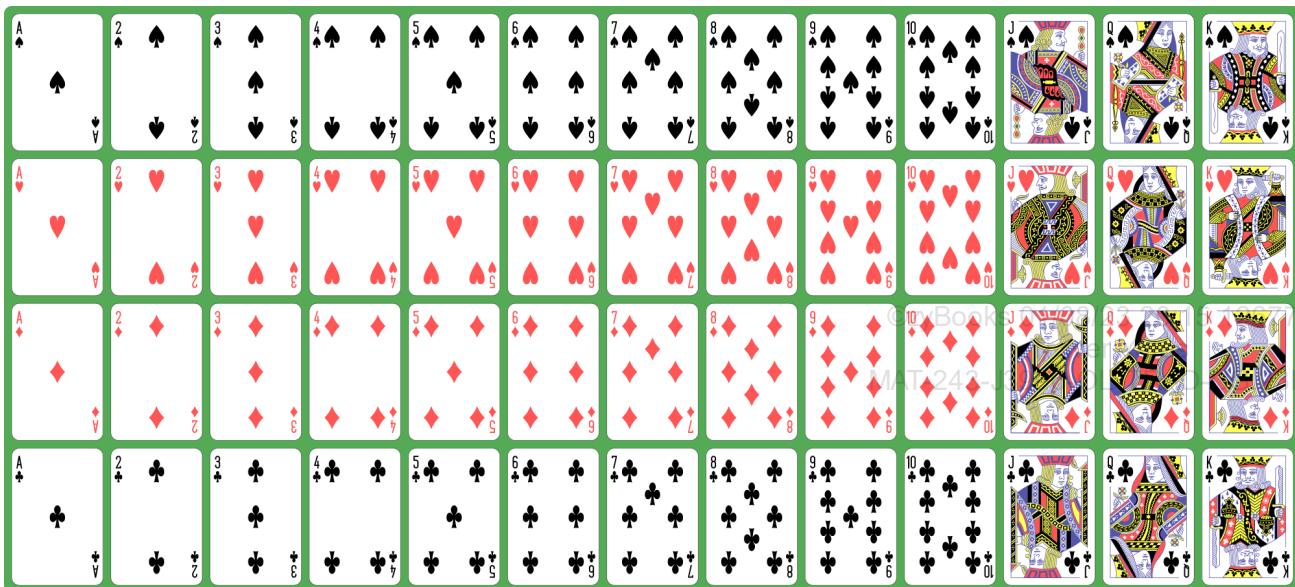
#### Example 1.14.1: Experiments, outcomes, and sample space.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

A standard deck of playing cards contains 52 cards. Each card has a rank of 2, 3, 4, 5, 6, 7, 8, 9, 10, J (for jack), Q (for queen), K (for king), or A (for ace). Each card also has a suit of spades (♠), clubs (♣), hearts (♥), or diamonds (♦). Spades and clubs are black cards, and hearts and diamonds are red cards. The cards in a standard deck are shown below.



Source: English pattern playing cards deck ([Dmitry Fomin / CC0 1.0](#) via Wikimedia Commons)<sup>1</sup>

A person draws a random card from the deck. The card drawn is the **3** of spades.

- Is drawing a random card from the deck an experiment or an outcome?
- Is drawing the **3** of spades an experiment or an outcome?
- What is the size of the sample space of this experiment?

### Solution

- The act of drawing a random card is an experiment. The experiment can be repeated many times, each time with a potentially different outcome from drawing the **3** of spades.
- Drawing the **3** of spades is one of many possible outcomes of the experiment. Repeating the experiment may result in a different card being drawn, which would be a different outcome.
- Since one of **52** different cards could be drawn each time the experiment is performed, the sample space contains **52** outcomes, one for each card in the deck.

#### PARTICIPATION ACTIVITY

1.14.2: Experiments, outcomes, and sample spaces.



- 1) A 6-sided die is rolled and the number of dots displayed is recorded. Is the roll of the die an experiment or an outcome?

- an experiment
- an outcome

©zyBooks 01/08/23 20:15 1267700  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



2) A 6-sided die is rolled and the number of dots displayed is recorded. 4 dots are displayed. Is the 4 dots being displayed an experiment or an outcome?

- an experiment
- an outcome

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



3) An experiment of rolling a die is repeated 3 times. The experiments result in the 3 outcomes 2, 4, and 5. Do the three outcomes represent the sample space?

- Yes
- No



4) A 6-sided die is rolled. Is the sample space  $\{1, 2, 3, 4, 5, 6\}$ ?

- Yes
- No



5) An experiment involves flipping two coins. Two coins are flipped, and heads is observed on the first coin and tails on the second coin, represented as  $HT$ . Is  $HT$  an outcome of this experiment?

- Yes
- No



6) An experiment involves flipping two coins. Two coins are flipped, and the outcome  $HT$  is observed. Is the sample space  $\{H, T\}$ ?

- Yes
- No

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## Events

A subset of the sample space is called an **event**. Ex: For a die roll, the event  $A$  is rolling an even number. The event  $A = \{2, 4, 6\}$  is a subset of the sample space  $\{1, 2, 3, 4, 5, 6\}$ .

A **compound event** is a subset of the sample space consisting of more than one outcome. Ex: The event  $A = \{2, 4, 6\}$  is a compound event since rolling an even number consists of three outcomes.

A **simple event** is a subset with a single outcome. Ex: The event  $C$  is rolling a **5** on the die. Thus,  $C = \{5\}$  is a simple event because  $C$  contains only one outcome.

### Example 1.14.2: Simple and compound events.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

A 6-sided die is rolled. For each of the following events, list the possible outcomes, and whether the event is simple or compound.

- a. A **5** is rolled.
  - b. A number less than **5** is rolled.
  - c. An even number is rolled.
  - d. An even prime number is rolled.
  - e. A number is rolled.
- 
- a. The event contains only one outcome and is  $\{5\}$ . Since the event has only one outcome, the event is a simple event.
  - b. The event contains four outcomes and is  $\{1, 2, 3, 4\}$ . Since the event has more than one outcome, the event is a compound event.
  - c. The event is  $\{2, 4, 6\}$ . Since the event has more than one outcome, the event is a compound event.
  - d. The only even prime number is **2**, so the event is  $\{2\}$ . Since the event has only one outcome, the event is a simple event.
  - e. A number will always be rolled, so the event contains all outcomes in the sample space, or  $\{1, 2, 3, 4, 5, 6\}$ . Since the event has more than one outcome, the event is a compound event.

#### PARTICIPATION ACTIVITY

#### 1.14.3: Events.



- 1) The event  $B$  is rolling less than a **3** on a six-sided die.  $B$  is best described as

\_\_\_\_\_.



- a simple event
- a compound event
- not an event

- 2) Noah conducts an experiment by flipping a coin twice. The sample

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



space for this experiment is  $\{HH, HT, TH, TT\}$ . Which of the following is a simple event?

- Noah flips two heads.
  - Noah flips one heads and one tails.
  - Since Noah is flipping the coin twice, no simple events can exist.
- 3) A gym teacher defines an experiment to be running a mile race three times in gym class. Sue runs the mile three times. Sue's fastest time was 6:10. This time can best be described as \_\_\_\_\_.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- a simple event
- a compound event
- not an event

PARTICIPATION  
ACTIVITY

1.14.4: Experiments.



Researchers sought to learn how many calories were expended by students in a typical school day. The researchers placed calorie measuring devices on children during school hours for a school day and then read the displayed value for each child. Possible calorie expenditure readings for one measuring device were integers from **1** to **5000** calories.

Select the definition that matches each term

- 1) Experiment
- 1, 2, 3, ..., 4999, 5000
  - 1200
  - A "low" calorie expenditure defined as fewer than 800 calories.
  - Placing a calorie measuring device on a child for a day and then reading the device value.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- 2) Outcome

- 1, 2, 3, . . . , 4999, 5000
- 1200
- A "low" calorie expenditure defined as fewer than 800 calories.
- Placing a calorie measuring device on a child for a day and then reading the device value.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## 3) Sample space

- 1, 2, 3, . . . , 4999, 5000
- 1200
- A "low" calorie expenditure defined as fewer than 800 calories.
- Placing a calorie measuring device on a child for a day and then reading the device value.

## 4) Event

- 1, 2, 3, . . . , 4999, 5000
- 1200
- A "low" calorie expenditure defined as fewer than 800 calories.
- Placing a calorie measuring device on a child for a day and then reading the device value.

Reset

## Probability

**Probability** is a measure of how likely an event is to occur. The probability of an event  $A$  is denoted  $P(A)$ , and is the sum of the probabilities of each outcome in the event.

One definition of probability is the number of desired outcomes divided by the total number of outcomes in the sample space, assuming that all outcomes are equally likely. Ex: If a (fair) coin is flipped, the probability that the coin turns up heads is 1 desired outcome of heads divided by 2 outcomes in the sample space (heads and tails), or  $\frac{1}{2}$ .

However, the size of the sample space often cannot be counted or determined in a practical way, or the different outcomes are not equally likely. Thus, another definition of probability is the relative

frequency of the desired outcome, or the proportion of times the outcome will occur when an experiment is repeated an infinite number of times. Ex: If a coin is flipped many times, the frequency of heads will approach  $\frac{1}{2}$  over time. Thus, the probability of heads is  $\frac{1}{2}$ . Ex: If the frequency of red hair in a population is 0.17, and the experiment of selecting a random person from the population is performed many times, the probability that a person with red hair is selected will approach 0.17.

**PARTICIPATION ACTIVITY**

1.14.5: Calculating the probability of a compound event

01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Animation content:**

undefined

**Animation captions:**

1. To find the probability of an event, the outcomes in the event should first be identified.
2. Since the die is fair, each of the six numbers is equally likely, and the probability of each number is  $\frac{1}{6}$ .
3. The probability of an event is the sum of the probabilities of the outcomes in the event.

**Example 1.14.3: Probability with equally likely outcomes.**

In many pen-and-paper role-playing games, a 20-sided die (or d20), shown below, is used to determine the success or failure of actions in the game. Each side of the d20 is labeled with a number from 1 to 20.



zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Source: D20 (4510874033) ([janet galore](#) / [CC-BY-2.0](#) via Wikimedia Commons)<sup>2</sup>

- a. A player critically fails a certain action if a 1, 2, or 3 is rolled. Assuming a fair d20, what is the probability of critical failure?
- b. What is the probability that the player avoids a critical failure?

**Solution**

a. Since the d20 is fair, the probability of each number is equally likely. Thus,

$P(1) = P(2) = P(3) = \frac{1}{20}$ . The probability of critical failure is the sum of the probabilities of the outcomes 1, 2, and 3, or

$$P(1) + P(2) + P(3) = \frac{1}{20} + \frac{1}{20} + \frac{1}{20} = \frac{3}{20}.$$

b. Rolling any number 4 or greater avoids critical failure. 17 such numbers can be rolled (4, 5, 6, ..., 20). Thus, the probability of avoiding critical failure is

$$P(4) + P(5) + \dots + P(20) = 17 \left(\frac{1}{20}\right) = \frac{17}{20}.$$

**PARTICIPATION ACTIVITY****1.14.6: Probability with equally likely outcomes.**

The same role-playing game also uses 10-sided dice, or d10. Each side of the d10 is labeled with a number from 1 to 10.



Source: White RPG dice including a d4, d6, d8, d10, d12, and a d20 (modified) ([Veikk0.ma](#) / [CC-BY-SA-4.0](#) via Wikimedia Commons)<sup>3</sup>

- 1) A player succeeds at an action if a 9 or 10 is rolled. Assuming a fair d10, what is the probability of success? Type as: #.#

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- 2) What is the probability that the player fails? Type as: #.#



**Check****Show answer**

Unlike flipping a coin or rolling a die, the probability of each outcome need not be the same.

### Example 1.14.4: Probability with outcomes that are not equally likely.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

The relative frequencies of the eight major human blood types in the United States are given below.

Blood type	Frequency
O+	0.38
O-	0.07
A+	0.34
A-	0.06
B+	0.09
B-	0.02
AB+	0.03
AB-	0.01

Source: [San Diego Blood Bank](#)<sup>4</sup>

- a. What is the probability that a randomly selected American has a blood type of O+?
- b. What is the probability that a randomly selected American has a blood type of either O+ or O-?
- c. What is the probability that a randomly selected American does not have blood type O?

### Solution

©zyBooks 01/08/23 20:15 1267703

- a. The event of randomly selecting an O+ American has only one outcome, O+. The probability is  $P(O+) = 0.38$ .
- b. The event of randomly selecting an O+ or O- American has two outcomes, O+ and O-. The probability is  $P(O+) + P(O-) = 0.38 + 0.07 = 0.45$ .
- c. The event of randomly selecting an American who is not type O has six outcomes, A+, A-, B+, B-, AB+, and AB-. The probability is  

$$P(A+) + P(A-) + P(B+) + P(B-) + P(AB+) + P(AB-) \\ = 0.34 + 0.06 + 0.09 + 0.02 + 0.03 + 0.01 = 0.55$$

**PARTICIPATION ACTIVITY**

1.14.7: Probability with outcomes that are not equally likely.



Flower color in a certain plant is controlled by a gene with two alleles, or variants, called A and a. Every plant has two alleles that make up the genotype, or genetic composition. In turn, the plant's genotype determines the plant's phenotype, or outward appearance.

In a certain garden, the genotype and phenotype frequencies of this plant are as follows.

Genotype	Phenotype	Frequency
AA	Purple flowers	0.49
Aa	Purple flowers	0.42
aa	White flowers	0.09

- 1) A flower is chosen at random from the garden. What is the probability that the flower is purple? Type as:

#.##

**Check****Show answer**

- 2) A flower is chosen at random from the garden. What is the probability that the flower carries at least one a allele? Type as: #.##

**Check****Show answer**

## References

- (\*1) Fomin, Dmitry. "English pattern playing cards deck" *Wikimedia Commons*. 13 June 2019, [https://upload.wikimedia.org/wikipedia/commons/thumb/8/81/English\\_pattern\\_playing\\_cards\\_deck.svg.png](https://upload.wikimedia.org/wikipedia/commons/thumb/8/81/English_pattern_playing_cards_deck.svg.png).

(\*2) janet galore. "D20 (4510874033)" *Wikimedia Commons*. 13 June 2019,  
[https://upload.wikimedia.org/wikipedia/commons/d/d8/D20\\_%284510874033%29.jpg](https://upload.wikimedia.org/wikipedia/commons/d/d8/D20_%284510874033%29.jpg).

(\*3) Veikk0.ma. "White RPG dice including a d4, d6, d8, d10, d12, and a d20" (modified) *Wikimedia Commons*. 13 June 2019,  
[https://upload.wikimedia.org/wikipedia/commons/b/b8/White\\_RPG\\_dice\\_including\\_a\\_d4%2C\\_d6%2C\\_d8%2C\\_d10%2C\\_d12%2C\\_d20%28modified%29.jpg](https://upload.wikimedia.org/wikipedia/commons/b/b8/White_RPG_dice_including_a_d4%2C_d6%2C_d8%2C_d10%2C_d12%2C_d20%28modified%29.jpg).

(\*4) San Diego Blood Bank. "What is the Most Common Blood Type?" *San Diego Blood Bank*. 13 June 2019. <https://www.sandielobloodbank.org/what-most-common-blood-type>

©zyBooks 01/08/23 20:15 1267703  
Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## 1.15 Addition rule and complements

### Events as sets

Events are sets of outcomes and subsets of the sample space. Thus, the terminology and notation of sets can be used to describe events and outcomes.

A sample space and sets of events in the sample space can be represented visually with a **Venn diagram**. In a Venn diagram, the sample space is represented by a rectangle. Every outcome appears inside the rectangle. An event is represented by a circle enclosing the outcomes in that event. Events with outcomes in common are represented as overlapping circles.

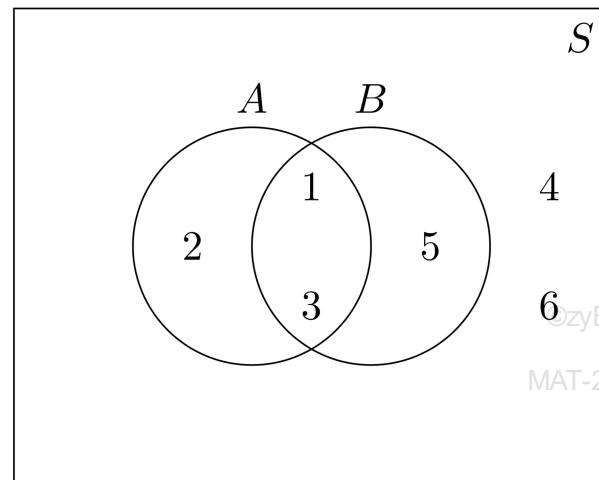
#### Example 1.15.1: Constructing a Venn diagram.

A 6-sided die is rolled. Let  $A$  be the event a number less than 4 is rolled. Let  $B$  be the event an odd number is rolled. Construct the Venn diagram for these events.

#### Solution

The sample space  $S = \{1, 2, 3, 4, 5, 6\}$  contains all 6 outcomes of rolling the die.  $A$  contains the three outcomes 1, 2, and 3 and  $B$  contains the three outcomes 1, 3, and 5. Outcomes 1 and 3 are common to both  $A$  and  $B$ . Outcome 2 is in  $A$  but not  $B$ , and outcome 5 is in  $B$  but not  $A$ . Outcomes 4 and 6 are in the sample space but in neither  $A$  nor  $B$ . Thus, the Venn diagram is as follows.

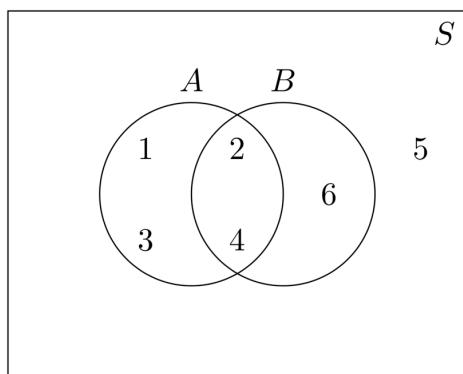
©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

**PARTICIPATION ACTIVITY**

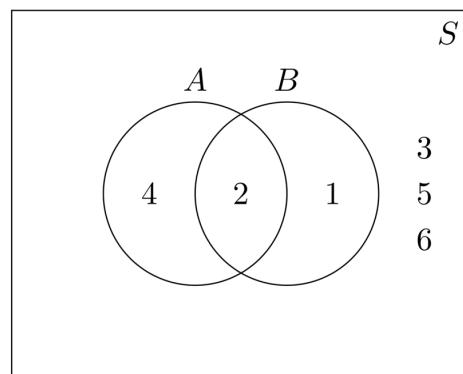
1.15.1: Using a Venn diagram to represent events.



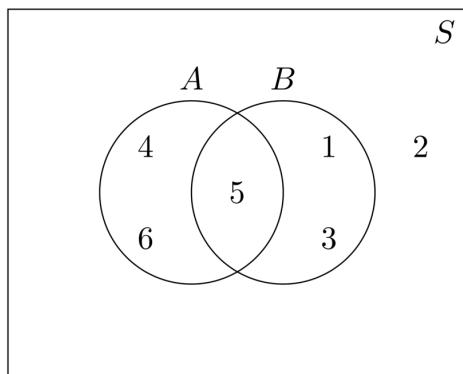
Each of the following Venn diagrams represent events when rolling a 6-sided die. Match each Venn diagram with the correct event descriptions.



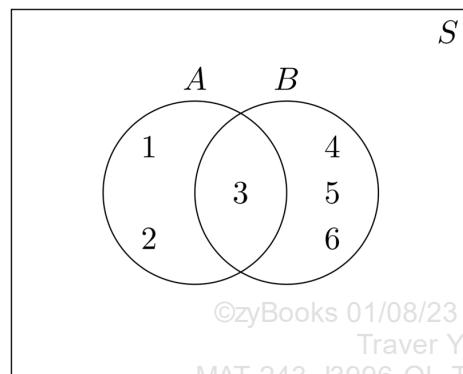
(a)



(b)



(c)



(d)

Select the definition that matches each term

- 1) (d)

- $A$  is the event of rolling a number less than 4.  
 $B$  is the event of rolling a number greater than 2.
- $A$  is the event of rolling a number less than 5.  
 $B$  is the event of rolling an even number.
- $A$  is the event of rolling an even number less than 5.  
 $B$  is the event of rolling a number less than 3.
- $A$  is the event of rolling a number greater than 3  
 $B$  is the event of rolling an odd number.

2) (c)

- $A$  is the event of rolling a number less than 4.  
 $B$  is the event of rolling a number greater than 2.
- $A$  is the event of rolling a number less than 5.  
 $B$  is the event of rolling an even number.
- $A$  is the event of rolling an even number less than 5.  
 $B$  is the event of rolling a number less than 3.
- $A$  is the event of rolling a number greater than 3  
 $B$  is the event of rolling an odd number.

3) (a)

- $A$  is the event of rolling a number less than 4.  
 $B$  is the event of rolling a number greater than 2.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- $A$  is the event of rolling a number less than 5.  
 $B$  is the event of rolling an even number.
- $A$  is the event of rolling an even number less than 5.  
 $B$  is the event of rolling a number less than 3.
- $A$  is the event of rolling a number greater than 3  
 $B$  is the event of rolling an odd number.

4) (b)

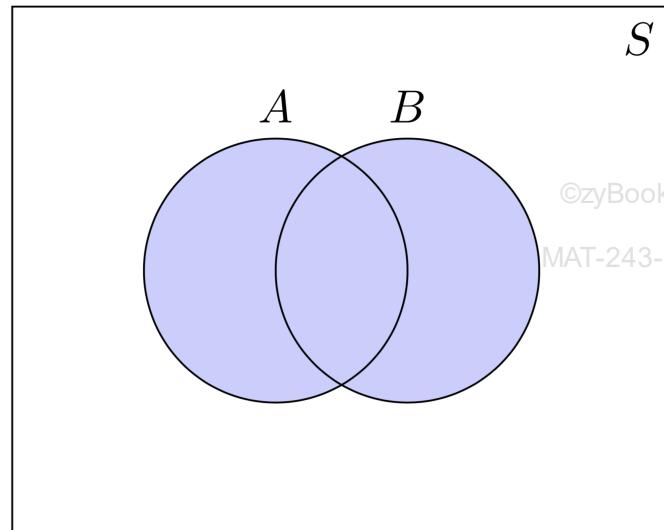
- $A$  is the event of rolling a number less than 4.  
 $B$  is the event of rolling a number greater than 2.
- $A$  is the event of rolling a number less than 5.  
 $B$  is the event of rolling an even number.
- $A$  is the event of rolling an even number less than 5.  
 $B$  is the event of rolling a number less than 3.
- $A$  is the event of rolling a number greater than 3  
 $B$  is the event of rolling an odd number.

Reset

## Union, intersection, and complement

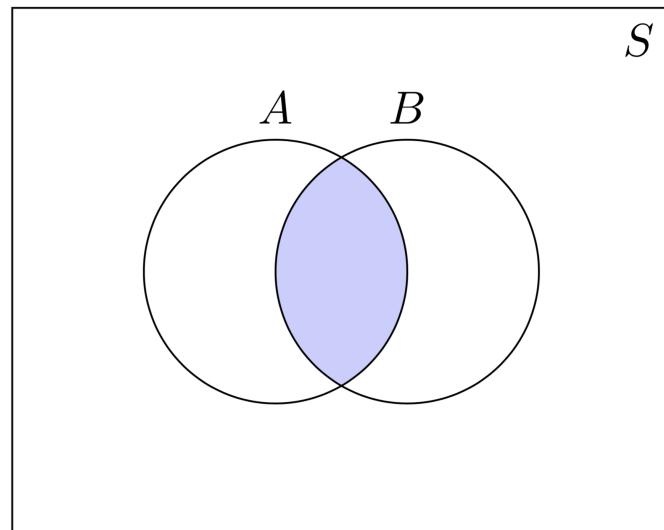
Several set operations are used frequently in probability.

The **union** of two events  $A$  and  $B$  is denoted  $A \cup B$  and is the event that includes outcomes in  $A$  or  $B$  or both.

Figure 1.15.1: Venn diagram of  $A \cup B$ .

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

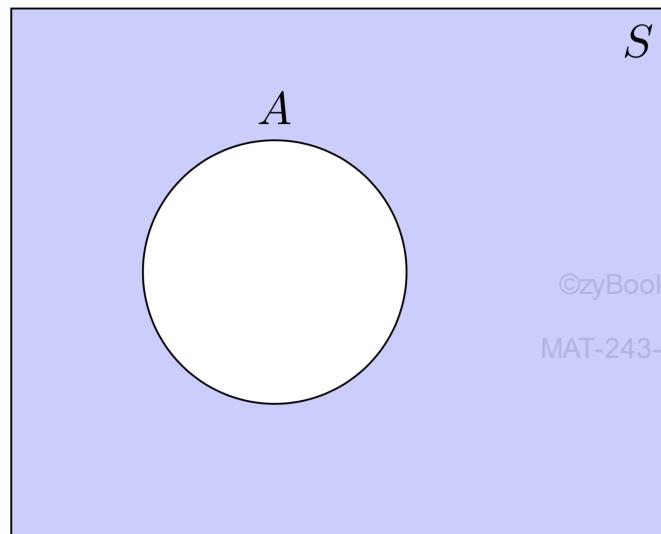
The **intersection** of two events  $A$  and  $B$  is denoted  $A \cap B$  and is the event consisting of outcomes that are in both  $A$  and  $B$ .

Figure 1.15.2: Venn diagram of  $A \cap B$ .

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The **complement** of an event  $A$  is denoted  $\bar{A}$  and is the event consisting of outcomes that are not in  $A$ .

Figure 1.15.3: Venn diagram of  $\bar{A}$ .



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

The **empty set** is the event consisting of no outcomes. Ex: A coin is tossed. The intersection of the event of the coin coming up heads and the event of the coin coming up tails is the empty set because no outcomes exist where the coin comes up both heads and tails. In other words, a coin cannot come up both heads and tails at the same time.

### Example 1.15.2: Union, intersection, and complement.

An integer between **1** and **10** (inclusive) is selected at random. Let  $A$  be the event a number greater than **5** is selected,  $B$  be the event that a number smaller than **5** is selected, and  $C$  be the event an even number is selected.

- What is the sample space?
- What is  $A \cap C$ ?
- What is  $A \cup C$ ?
- What is  $A \cap B$ ?
- What is  $A \cup B$ ?
- What is  $\overline{A}$ ?

### Solution

- ©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3
- The sample space is all of the integers between **1** and **10** inclusive,  $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$ .
  - $A \cap C$  is the set of outcomes in the sample space that are both greater than **5** and even, or  $\{6, 8, 10\}$ .
  - $A \cup C$  is the set of outcomes in the sample space that are either greater than **5** or even or both, or  $\{2, 4, 6, 7, 8, 9, 10\}$ .
  - Since a number cannot be both greater than **5** and less than **5**,  $A \cap B$  has no outcomes and is the empty set  $\emptyset$ .

- e.  $A \cup B$  is the set of outcomes in the sample space that are either greater than 5 or less than 5, or  $\{1, 2, 3, 4, 6, 7, 8, 9, 10\}$ .
- f.  $\overline{A}$  is the set of outcomes in the sample space that does not include the outcomes of  $A$ . Since  $A = \{6, 7, 8, 9, 10\}$ ,  $\overline{A} = \{1, 2, 3, 4, 5\}$ .

**PARTICIPATION ACTIVITY**

## 1.15.2: Union, intersection, and complement.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



A 6-sided die is rolled. Let  $A$  be the event that an odd number is rolled,  $B$  be the event that the numbers 1, 2, or 3 is rolled, and  $C$  be the event that an even number is rolled. Match each set of outcomes with the correct expression.

Select the definition that matches each term

1)  $A \cap B$

- $\emptyset$
- $\{1, 3\}$
- $\{1, 2, 3, 5\}$
- $\{1, 2, 3, 4, 6\}$
- $\{4, 5, 6\}$

2)  $A \cup B$

- $\emptyset$
- $\{1, 3\}$
- $\{1, 2, 3, 5\}$
- $\{1, 2, 3, 4, 6\}$
- $\{4, 5, 6\}$

3)  $B \cup C$

- $\emptyset$
- $\{1, 3\}$
- $\{1, 2, 3, 5\}$
- $\{1, 2, 3, 4, 6\}$
- $\{4, 5, 6\}$

4)  $\overline{B}$

- $\emptyset$

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- $\{1, 3\}$
- $\{1, 2, 3, 5\}$
- $\{1, 2, 3, 4, 6\}$
- $\{4, 5, 6\}$

5)  $A \cap C$ 

- $\emptyset$
- $\{1, 3\}$
- $\{1, 2, 3, 5\}$
- $\{1, 2, 3, 4, 6\}$
- $\{4, 5, 6\}$

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

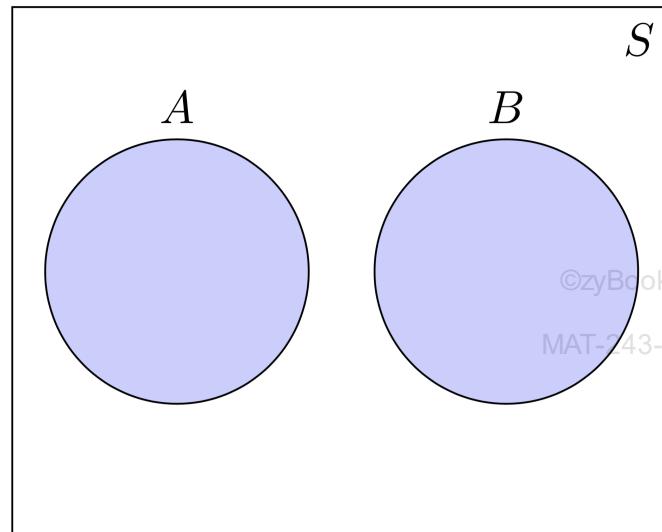
**Reset**

## Mutually exclusive events

Two events  $A$  and  $B$  are **mutually exclusive** if  $A \cap B = \emptyset$ , that is,  $A$  and  $B$  have no outcomes in common. Ex: A coin is tossed. The event of the coin coming up heads and the event of the coin coming up tails are mutually exclusive because the coin cannot come up both heads and tails at the same time.

Mutually exclusive events are represented by non-overlapping circles in a Venn diagram.

Figure 1.15.4: Venn diagram of mutually exclusive events  $A$  and  $B$ .

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

### Example 1.15.3: Mutually exclusive events.

An integer between **1** and **10** (inclusive) is selected at random. Let  $A$  be the event a number greater than **5** is selected,  $B$  be the event that a number smaller than **5** is selected, and  $C$  be the event an even number is selected.

- Are events  $A$  and  $B$  mutually exclusive?
- Are events  $B$  and  $C$  mutually exclusive?
- Are events  $A$  and  $\bar{A}$  mutually exclusive?

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

### Solution

- A number cannot be both greater than **5** and less than **5**. Therefore,  $A \cap B = \emptyset$ , and  $A$  and  $B$  are mutually exclusive.
- Since  $B = \{1, 2, 3, 4\}$  and  $C = \{2, 4, 6, 8, 10\}$ ,  $B \cap C = \{2, 4\}$ , and  $B$  and  $C$  are not mutually exclusive.
- Since  $A = \{6, 7, 8, 9, 10\}$  and  $\bar{A} = \{1, 2, 3, 4, 5\}$ ,  $A$  and  $\bar{A}$  have no outcomes in common and are mutually exclusive. In general, an event and the complement of that event are mutually exclusive.

#### PARTICIPATION ACTIVITY

1.15.3: Identifying mutually exclusive events.



Given each description, select the pair of events  $A$  and  $B$  that are mutually exclusive.

1) A **6**-sided die is rolled.



- $A$  is the event of rolling a number less than **5**.  
 $B$  is the event of rolling a number greater than **3**.
- $A$  is the event of rolling a number greater than **5**.  
 $B$  is the event of rolling a number less than **3**.
- $A$  is the event of rolling an even number.  
 $B$  is the event of rolling a prime number.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

2) A card from a standard **52**-card deck of playing cards is drawn.



- $A$  is the event of drawing a 5.  
 $B$  is the event of drawing a club.
- $A$  is the event of drawing a club.  
 $B$  is the event of drawing a black card.
- $A$  is the event of drawing a club.  
 $B$  is the event of drawing a spade.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

- 3) A coin is flipped five times. Ex: HTHTH means that the five flips are heads, tails, heads, tails, and heads, in that order.

- $A$  is the event that the first flip comes up heads.  
 $B$  is the event that the last flip comes up tails.
- $A$  is the event that the first two flips come up heads.  
 $B$  is the event that the last two flips come up tails.
- $A$  is the event that the first three flips come up heads.  
 $B$  is the event that the last three flips come up tails.



## The axioms of probability

Probability has three fundamental properties, or axioms.

### The axioms of probability

Let  $A$  and  $B$  be mutually exclusive events in the sample space  $S$ .

01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

1.  $0 \leq P(A) \leq 1$
2.  $P(S) = 1$
3.  $P(A \cup B) = P(A) + P(B)$

The smallest possible probability is **0** when an outcome never occurs, **1** when an outcome always occurs. Thus, axiom 1 states the probability of an event must be between **0** and **1**, inclusive.

An outcome in the sample space must always occur in the long run. Thus, axiom 2 states that the probability of an outcome in the sample space is **1**.

Individual outcomes are mutually exclusive events, so all events are unions of mutually exclusive events. Thus, axiom 3 implies that the long-term proportion that each of the outcomes in an event occurs is the sum of the probabilities of each outcome.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

The axioms of probability can be used to prove some properties of probability.

#### Example 1.15.4: Probability of the empty set.

Use the axioms of probability to show that  $P(\emptyset) = 0$ .

#### Solution

The empty set  $\emptyset$  is the event of no outcomes and the sample space  $S$  is the event of every outcome. Thus,  $\emptyset$  and  $S$  are mutually exclusive. Also, since  $S$  includes every outcome,  $S \cup \emptyset = S$ . By axiom 3,  $P(S \cup \emptyset) = P(S) = P(S) + P(\emptyset)$ . By axiom 2,  $P(S) = 1$ . Thus,  $1 = 1 + P(\emptyset)$ , and  $P(\emptyset) = 0$ .

#### PARTICIPATION ACTIVITY

1.15.4: The axioms of probability.



Let  $A$  be an event in the sample space  $S$ .

1)  $P(A) > 1$



True

False

2)  $P(A) < 0$



True

False

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

#### The complement rule

The axioms of probability have two important consequences, which are the complement rule and the addition rule of probability. The **complement rule** relates the probability of an event to the probability of the complement of the event.

## Theorem 1.15.1: The complement rule.

For any event  $A$ , if the probability of  $A$  is  $P(A)$ , the probability of  $\bar{A}$  is  $P(\bar{A}) = 1 - P(A)$ .

### Proof.

By definition of complement,  $A$  and  $\bar{A}$  are mutually exclusive, and  $A \cup \bar{A} = S$ . By axiom 3,  $P(A) + P(\bar{A}) = P(S)$ . By axiom 2,  $P(S) = 1$ . Thus,  $P(A) + P(\bar{A}) = 1$ . Rearranging gives  $P(\bar{A}) = 1 - P(A)$ . ■

### PARTICIPATION ACTIVITY

1.15.5: The complement rule.



### Animation content:

undefined

### Animation captions:

1. The union of  $A$  and  $\bar{A}$  is the sample space  $S$ .
2. Since  $A$  and  $\bar{A}$  are mutually exclusive,  $P(A) + P(\bar{A}) = P(S)$ , by axiom 3.
3. By axiom 2,  $P(S) = 1$ . Thus,  $P(A) + P(\bar{A}) = 1$ .
4. Rearranging gives the complement rule  $P(\bar{A}) = 1 - P(A)$ .

If the complement of an event can be inferred from the description of the event, the probability of the complement can be calculated using the complement rule, as the following example illustrates.

### Example 1.15.5: Applying the complement rule.

The San Francisco Bridges and the Mississauga Dinosaurs are competing for a basketball championship decided by a best-of-7 series, meaning that the first team to win 4 games wins the championship. Sports experts predict a probability of 0.744 for the Bridges to sweep the series, meaning that the Bridges win 4 games in a row and the Dinosaurs do not win any games.

Let  $A$  be the event that the Bridges sweep the series.

- a. What is  $P(A)$ ?
- b. What is  $\bar{A}$ ?
- c. What is  $P(\bar{A})$ , and what is the interpretation of this probability?

## Solution

- a. Since the probability of the Bridges sweeping the series is  $0.744$ ,  $P(A) = 0.744$ .
- b. The sample space consists of outcomes in which the Dinosaurs win 0, 1, 2, 3, or 4 games. Since  $A$  is the event that the Dinosaurs win 0 games, the complement of  $A$  is the event that the Dinosaurs win 1, 2, 3, or 4 games. In other words,  $\bar{A}$  is the event that the Dinosaurs win at least one game, that is, the Bridges do not sweep the series.
- c. By the complement rule,  $P(\bar{A}) = 1 - P(A) = 1 - 0.744 = 0.256$ .  $P(\bar{A})$  is the probability that the Dinosaurs win at least one game in the series.

### PARTICIPATION ACTIVITY

1.15.6: Applying the complement rule.



In a certain city, 17% of residents own a house, 48% of residents own a car, and 42% of residents own neither a house nor a car. Let  $A$  be the event that a randomly selected resident from the city owns a house and  $B$  be the event that a randomly selected resident from the city owns a car.

- 1) What is the probability that a resident does not own a house?

**Check****Show answer**

- 2) What is the probability that a resident does not own a car?

**Check****Show answer**

- 3) What is the probability that a resident owns a car, a house, or both?

**Check****Show answer**

## The addition rule

The **addition rule** generalizes axiom 3 to events that are not mutually exclusive.

### Theorem 1.15.2: The addition rule.

For any events  $A$  and  $B$ ,  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

**Proof.**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

$$\begin{aligned} P(A \cup B) &= P(A \cap \bar{B}) + P(A \cap B) + P(\bar{A} \cap B) \\ &= P(A \cap \bar{B}) + P(A \cap B) + P(\bar{A} \cap B) + P(A \cap B) - P(A \cap B) \\ &= [P(A \cap \bar{B}) + P(A \cap B)] + [P(\bar{A} \cap B) + P(A \cap B)] - P(A \cap B) \\ &= P(A) + P(B) - P(A \cap B) \end{aligned}$$

■

If  $A$  and  $B$  are mutually exclusive, then  $A \cap B = \emptyset$ ,  $P(A \cap B) = 0$ , and the addition rule reduces to axiom 3.

PARTICIPATION ACTIVITY

1.15.7: The addition rule.



### Animation content:

undefined

### Animation captions:

1. Let events  $A$  and  $B$  be non-mutually exclusive events.
2. Event  $A$  has probability  $P(A)$  and event  $B$  has probability  $P(B)$ .
3. Adding  $P(A)$  and  $P(B)$  results in  $P(A \cap B)$  being counted twice.
4. Subtracting  $P(A \cap B)$  gives  $P(A \cup B)$  for non-mutually exclusive events  $A$  and  $B$ .

### Example 1.15.6: Applying the addition rule.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

A card is drawn from a standard deck of 52 playing cards. Let  $A$  be the event that a 5 is drawn, and let  $B$  be the event that a heart is drawn.

- a. What are  $P(A)$  and  $P(B)$ ?
- b. Are  $A$  and  $B$  mutually exclusive events?
- c. What is  $P(A \cap B)$ ?
- d. What is  $P(A \cup B)$ ?

**Solution**

- a. Of the 52 cards in the deck, 4 cards are 5s. Thus, the probability of drawing a 5 is  $P(A) = \frac{4}{52} = \frac{1}{13}$ . Similarly, 13 cards are hearts, so the probability of drawing a heart is  $P(B) = \frac{13}{52} = \frac{1}{4}$ .

b. Since drawing the 5 of hearts is an outcome in both  $A$  and  $B$ ,  $A$  and  $B$  are not mutually exclusive events.

c.  $A \cap B$  is the event of drawing a card that is both a 5 and a heart, that is, drawing the single 5 of hearts in the deck. Thus,  $P(A \cap B) = \frac{1}{52}$ .

d.  $P(A \cup B)$  is the event of drawing either a 5 or a heart, or both. Since  $A$  and  $B$  are not mutually exclusive, the addition rule is used.

©zyBooks 01/08/23 20:15 1267703

MAT-243-J3996-OL-TRAD-UG.23EW3

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ &= \frac{1}{13} + \frac{1}{4} - \frac{1}{52} \\ &= \frac{16}{52} \\ &= \frac{4}{13} \end{aligned}$$

**PARTICIPATION ACTIVITY**

1.15.8: Applying the addition rule.



In a certain city, 46% of residents own a house, 67% of residents own a car, and 21% of residents own both a house and a car. Let  $A$  be the event that a randomly selected resident from the city owns a house and  $B$  be the event that a randomly selected resident from the city owns a car.

- 1) What is the probability that a resident owns a house, a car, or both?

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- 2) What is the probability that a resident owns neither a house nor a car?

**Check****Show answer**

- 3) What is the probability that a resident owns a house but not a car?

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## 1.16 Multiplication rule and independence

### Independent events

Two events are **independent** if the probability of one event does not affect the probability of the other. Ex: A nickel and a dime are flipped. Whether the nickel comes up heads or tails is not affected by whether the dime comes up heads or tails. Thus, flipping the nickel and flipping the dime are independent events.

#### Example 1.16.1: Identifying independent events.

- a. A coin is flipped and a 6-sided die is rolled, separately. Let  $A$  be the event of getting tails with the coin and  $B$  be the event of getting a number less than 3 with the die. Are  $A$  and  $B$  independent events?
- b. A coin is flipped twice. Let  $A$  be the event of getting tails on the first flip and  $B$  be the event of getting tails on the second flip. Are  $A$  and  $B$  independent events?
- c. A coin is flipped twice. Let  $A$  be the event of getting heads on the second flip and  $B$  be the event of getting tails on the second flip. Are  $A$  and  $B$  independent events?
- d. A bag contains 5 red marbles and 5 black marbles. Two marbles are picked at random from the bag without replacement. Let  $A$  be the event of picking a red marble on the first pick, and  $B$  be the event of picking a red marble on the second pick. Are  $A$  and  $B$  independent events?

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

### Solution

- a. Assuming that the coin and the die are not physically connected, the coin flip will not affect the die roll and vice versa. Thus,  $A$  and  $B$  are independent events.
- b. Since the outcome of the first flip does not influence the outcome of the second flip,  $A$  and  $B$  are independent events.

- c. Since the second flip cannot be both heads and tails, event  $A$  occurring prevents event  $B$  from occurring, and vice versa. Thus,  $A$  and  $B$  are dependent events.
- d. Since a marble is removed from the bag after event  $A$ , the probability of  $B$  is changed by event  $A$  occurring. Thus,  $A$  and  $B$  are dependent events.

**PARTICIPATION ACTIVITY****1.16.1: Identifying independent events.**

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



Determine whether  $A$  and  $B$  are independent or dependent events.

- 1) One card is drawn from each of two decks of cards.



$A$  is the event of drawing a red card from the first deck.

$B$  is the event of drawing a red card from the second deck.

- Independent events  
 Dependent events

- 2) Two cards are drawn from the same deck of cards. The first card is not put back into the deck after being drawn.



$A$  is the event of drawing a red card on the first draw.

$B$  is the event of drawing a red card on the second draw.

- Independent events  
 Dependent events

- 3) Two cards are drawn from the same deck of cards. The first card is put back into the deck after being drawn.



$A$  is the event of drawing a red card on the first draw.

$B$  is the event of drawing a red card on the second draw.

- Independent events  
 Dependent events

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## The multiplication rule

The **multiplication rule** gives the probability of 2 independent events happening together.

### Theorem 1.16.1: The multiplication rule.

Let  $A$  and  $B$  be independent events. The probability that both  $A$  and  $B$  occur is

©zyBooks 01/08/23 20:15 1267703

Traver Yates

$$P(A \cap B) = P(A)P(B)$$

MAT-243-J3996-OL-TRAD-UG.23EW3

Independent events are different from mutually exclusive events. Mutually exclusive events are dependent events by definition, because if one event occurs, the other event cannot occur. However, the addition rule can be used to find the probability of the union of two events regardless of whether the events are independent or mutually exclusive.

### Example 1.16.2: Applying the multiplication rule.

- a. A coin is flipped and a 6-sided die is rolled, separately. Let  $A$  be the event of getting tails with the coin and  $B$  be the event of getting a number less than 3 with the die. What is  $P(A \cap B)$ ?
- b. A coin is flipped twice. Let  $A$  be the event of getting tails on the first flip and  $B$  be the event of getting tails on the second flip. What is  $P(A \cap B)$ ?

#### Solution

- a. Since  $A$  and  $B$  are independent events, the multiplication rule applies.  $P(A) = \frac{1}{2}$  and  $P(B) = \frac{1}{3}$ . Thus,  $P(A \cap B) = P(A)P(B) = \left(\frac{1}{2}\right)\left(\frac{1}{3}\right) = \frac{1}{6}$ .
- b. Since  $A$  and  $B$  are independent events, the multiplication rule applies.  $P(A) = P(B) = \frac{1}{2}$ . Thus,  $P(A \cap B) = P(A)P(B) = \left(\frac{1}{2}\right)\left(\frac{1}{2}\right) = \frac{1}{4}$ .

#### PARTICIPATION ACTIVITY

1.16.2: Probabilities of independent and mutually exclusive events  
©zyBooks 01/08/23 20:15 1267703  
Traver Yates

Let events  $A$  and  $B$  have probabilities  $P(A) = 0.3$  and  $P(B) = 0.5$ , respectively.

- 1) If  $A$  and  $B$  are mutually exclusive, what is  $P(A \cap B)$ ?



**Check****Show answer**

- 2) If  $A$  and  $B$  are mutually exclusive, what is  $P(A \cup B)$ ?

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- 3) If  $A$  and  $B$  are independent, what is  $P(A \cap B)$ ?

**Check****Show answer**

- 4) If  $A$  and  $B$  are independent, what is  $P(A \cup B)$ ?

**Check****Show answer**

## The generalized multiplication rule

The multiplication rule can also be generalized to more than 2 independent events. Given independent events  $A_1, A_2, \dots, A_n$ ,

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1)P(A_2)\dots P(A_n)$$

Example 1.16.3: Applying the multiplication rule to more than two events.

Three 6-sided dice are rolled.

- What is the probability that three 1s are rolled?
- What is the probability that three 1s or three 6s are rolled?
- What is the probability that at least one 1 is rolled?

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Solution

Each die roll is an independent event. Thus, the multiplication rule can be applied.

a. For each die, the probability of rolling a 1 is  $\frac{1}{6}$ .

$$\begin{aligned} P(\text{three 1s}) &= P(1 \text{ on die 1} \cap 1 \text{ on die 2} \cap 1 \text{ on die 3}) \\ &= P(1 \text{ on die 1}) \cdot P(1 \text{ on die 2}) \cdot P(1 \text{ on die 3}) \\ &= \left(\frac{1}{6}\right) \left(\frac{1}{6}\right) \left(\frac{1}{6}\right) \\ &= \frac{1}{216} \end{aligned}$$

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

b. The event of rolling three 1s and the event of rolling three 6s are mutually exclusive events. Thus, the addition rule is used.

$$\begin{aligned} P(\text{three 1s} \cup \text{three 6s}) &= P(\text{three 1s}) + P(\text{three 6s}) \\ &= \frac{1}{216} + \frac{1}{216} \\ &= \frac{2}{216} \\ &= \frac{1}{108} \end{aligned}$$

c. The complement rule is used with the multiplication rule. The complement of the event that at least one 1 is rolled is the event that no 1s are rolled. For each die, the probability of not rolling a 1 is  $P(\text{not 1}) = 1 - P(1) = 1 - \frac{1}{6} = \frac{5}{6}$ .

$$\begin{aligned} P(\text{at least one 1}) &= 1 - P(\text{no 1s}) \\ &= 1 - P(\text{not 1 on die 1}) \cdot P(\text{not 1 on die 2}) \cdot P(\text{not 1 on die 3}) \\ &= 1 - \left(\frac{5}{6}\right) \left(\frac{5}{6}\right) \left(\frac{5}{6}\right) \\ &= 1 - \frac{125}{216} \\ &= \frac{91}{216} \end{aligned}$$

Thus, the probability that at least one 1 is rolled is  $\frac{91}{216}$ .

**PARTICIPATION ACTIVITY**
**1.16.3: Applying the multiplication rule.**

©zyBooks 01/08/23 20:15 1267703

Traver Yates



MAT-243-J3996-OL-TRAD-UG.23EW3

A basketball player is fouled in the act of shooting a three-point shot and is awarded three free throws. The player makes free throws 80% of the time. Assume that each free throw is an independent event.

- 1) What is the probability that the player makes all three free throws?



Type as: #.###

**Check****Show answer**

- 2) What is the probability that the player misses all three free throws?

Type as: #.###

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Check****Show answer**

- 3) What is the probability that the player misses at least one free throw? Type as: #.###

**Check****Show answer**

- 4) What is the probability that the player makes at least one free throw? Type as: #.###

**Check****Show answer**

## 1.17 Conditional probability



This section has been set as optional by your instructor.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Conditional probability

Conditional probability is a measure of the probability of one event given that another event has occurred. The **conditional probability** of an event  $A$  given event  $B$  has occurred is denoted as

$$P(A|B)$$
 and is equal to  $\frac{P(A \cap B)}{P(B)}$ .

If event  $B$  has occurred, then the sample space in which  $A$  can occur is reduced from the entire sample space  $S$  to only  $B$ . The subset of  $A$  that occurs with  $B$  as the sample space is  $P(A \cap B)$ . Thus,  $P(A|B)$  is the proportion of  $P(B)$  represented by  $P(A \cap B)$ .

**PARTICIPATION ACTIVITY**

1.17.1: Conditional probability with Venn diagrams.

**Animation content:**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

undefined

**Animation captions:**

1. The conditional probability  $P(A|B)$  is the probability that  $A$  occurs given that  $B$  has occurred.
2. If event  $B$  has already occurred, then the sample size is reduced from  $S$  to  $B$ .  $A \cap B$  is the subset of outcomes in  $A$  with  $B$  as the sample space.
3. Thus, the probability of  $A$  occurring given that  $B$  has occurred is  $P(A|B) = \frac{P(A \cap B)}{P(B)}$ .

**Example 1.17.1: Conditional probability.**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

A 6-sided die is rolled. Let  $A$  be the event a number greater than 4 is rolled and  $B$  be the event a prime number is rolled.

- What is  $P(A|B)$ ?
- What is  $P(B|A)$ ?

### Solution

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Since  $A = \{5, 6\}$ ,  $B = \{2, 3, 5\}$ , and  $A \cap B = \{5\}$ ,  $P(A) = \frac{2}{6} = \frac{1}{3}$ ,  $P(B) = \frac{3}{6} = \frac{1}{2}$ , and  $P(A \cap B) = \frac{1}{6}$ .

- Given that a prime number was rolled, the sample space is reduced to  $B = \{2, 3, 5\}$ .

With  $B$  as the sample space, only the outcome 5 is also in  $A$ . Therefore, the probability of rolling a number greater than 4 given that a prime number was rolled is

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{1/6}{1/2} = \frac{1}{3}.$$

- Given that a number greater than 4 was rolled, the sample space is reduced to

$A = \{5, 6\}$ . With  $A$  as the sample space, only the outcome 5 is also in  $B$ . Therefore, the probability of rolling a prime number given that a number greater than 4 was rolled is  $P(B|A) = \frac{P(B \cap A)}{P(A)} = \frac{1/6}{1/3} = \frac{1}{2}$ .

### Example 1.17.2: Conditional probability in cat tail genetics.

The Manx cat is a breed of cat with a shortened or missing tail.



©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

Source: A Rumpy Manx Cat (Michelle Weigold / [CC-BY-SA-4.0](#) via Wikimedia Commons)<sup>1</sup>

When two Manx cats are bred together, three outcomes are possible. A non-Manx kitten with a normal tail is born  $\frac{1}{4}$  of the time. A Manx kitten with a missing tail is born  $\frac{1}{2}$  of the time. However,  $\frac{1}{4}$  of the time, the mutation controlling tail development has a lethal effect on spinal cord development, and a kitten is not born.

- Two Manx cats are bred together. What is the probability of a kitten being born?
- Two Manx cats are bred together. What fraction of kittens are expected to be Manx?

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Solution**

- Let  $B$  be the event that a kitten is born.  $B$  is the complement of the event that a kitten is not born. Thus,  $P(B) = 1 - \frac{1}{4} = \frac{3}{4}$ .
- Let  $M$  be the event a Manx kitten is born. The probability of the event  $M \cap B$  is the same as the probability of  $M$ , that is, the probability that a Manx kitten is born and that a kitten is born is the same as the probability that a Manx kitten is born. Thus,  $P(M \cap B) = P(M) = \frac{1}{2}$ .

The fact that a kitten is born means that  $B$  has occurred. Thus, the probability of a Manx kitten among the kittens born is  $P(M|B) = \frac{P(M \cap B)}{P(B)} = \frac{1/2}{3/4} = \frac{2}{3}$ .

**PARTICIPATION ACTIVITY**

1.17.2: Conditional probability in flower genetics.



Flower color in a certain plant is controlled by a gene with two alleles, or variants, called A and a. Every plant has two alleles that make up the genotype, or genetic composition. In turn, the plant's genotype determines the plant's phenotype, or outward appearance.

In a certain garden, the genotype and phenotype frequencies of this plant are as follows.

Genotype	Phenotype	Frequency
AA	Violet flowers	0.35
Aa	Violet flowers	0.58
aa	White flowers	0.07

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- A flower is chosen at random from the garden. What is the probability that the flower is violet? Type as: #.##

**Check****Show answer**

- 2) A violet flower is chosen at random from the garden. What is the probability that the flower carries at least one a allele? Type as: #.###

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**Check****Show answer**

- 3) A flower with at least one a allele is chosen at random from the garden. What is the probability that the flower is white? Type as: #.###

**Check****Show answer**

## Properties of conditional probability

If  $A$  and  $B$  are mutually exclusive events, then  $P(A \cap B) = 0$  and  $P(A|B) = P(B|A) = 0$ , that is, if one event is known to have occurred, the other mutually exclusive event cannot occur.

Rearranging the multiplication rule for independent events  $P(A \cap B) = P(A)P(B)$  gives  $P(A) = \frac{P(A \cap B)}{P(B)}$ . Thus,  $P(A) = P(A|B)$  if and only if  $A$  and  $B$  are independent, that is, the probability of  $A$  occurring is the same whether or not  $B$  has occurred. Similarly,  $P(B) = P(B|A)$  if and only if  $A$  and  $B$  are independent. This property can be used to test whether two events are independent or dependent.

Example 1.17.3: Determining whether two events are independent or dependent.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- If  $P(A|B) = 0.8$ ,  $P(B) = 0.8$ , and  $P(A) = 0.6$ , are  $A$  and  $B$  independent events?
- If  $P(A|B) = 0.8$ ,  $P(B) = 0.6$ , and  $P(A) = 0.8$ , are  $A$  and  $B$  independent events?

**Solution**

- a. Since  $P(A|B) \neq P(A)$ ,  $A$  and  $B$  are dependent events.
- b. Since  $P(A|B) = P(A) = 0.8$ ,  $A$  and  $B$  are independent events.

**PARTICIPATION ACTIVITY**

1.17.3: Determining whether two events are independent or dependent.



©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Determine whether  $A$  and  $B$  are independent events using the information provided.

1)  $P(A|B) = 0.3$

$P(B) = 0.6$

$P(A) = 0.5$

 Independent Dependent

2)  $P(A|B) = 0.3$

$P(B) = 0.8$

$P(A) = 0.3$

 Independent Dependent

3)  $P(B|A) = 0.8$

$P(B) = 0.8$

$P(A) = 0.3$

 Independent Dependent

4)  $P(A) = 0.3$

$P(B) = 0.3$



$A$  and  $B$  are mutually exclusive events.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

 Independent Dependent**The law of total probability**

The **law of total probability** states that if the sample space is partitioned into two or more mutually exclusive subevents, the probability of an event  $B$  can be expressed in terms of conditional probabilities given each of the subevents.

#### PARTICIPATION ACTIVITY

#### 1.17.4: The law of total probability.



#### Animation content:

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

undefined

#### Animation captions:

1. Let  $B$  be an event in  $S$ .
2. Subdivide  $S$  into three mutually exclusive subregions,  $S_1$ ,  $S_2$ , and  $S_3$ .
3. Thus,  $B$  is the union of the events  $B \cap S_1$ ,  $B \cap S_2$ , and  $B \cap S_3$ .
4. Next, the addition rule is applied.
5. The definition of conditional probability is used to rewrite  $P(B)$  in terms of conditional probabilities.

When a sample space is partitioned into two subevents, the subevents can be considered as one event and the complement of that event.

#### Example 1.17.4: The law of total probability with a partition of two subevents.

In a certain hospital's cancer ward, 30% of patients have a certain kidney disease that complicates the use of a certain cancer drug. The success rate of the drug is 90% in patients without the kidney disease and 65% in patients with the kidney disease. What is the probability that a randomly selected patient in the cancer ward is successfully treated with the drug?

#### Solution

Let  $S$  be the sample space of all patients in the cancer ward. Partition  $S$  into two subevents,  $K$  representing selecting a patient with the kidney disease, and  $\bar{K}$ , representing selecting a patient without the kidney disease. Let  $D$  be the event that the drug treatment is successful. Thus,  $P(D)$  is the probability that a randomly selected patient in the cancer ward is successfully treated with the drug.

The probability of selecting a patient with the kidney disease is  $P(K) = 30\% = 0.3$ . Thus, the probability of selecting a patient without the kidney disease is

$$P(\bar{K}) = 1 - P(K) = 1 - 0.3 = 0.7.$$

Given that a patient has the kidney disease, the probability that the drug treatment is successful is  $P(D|K) = 65\% = 0.65$ . Given that the patient does not have the kidney disease, the probability that the drug treatment is successful is  $P(D|\bar{K}) = 90\% = 0.9$ .

Finally, the law of total probability is used to find  $P(D)$ .

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

$$\begin{aligned} P(D) &= P(D|K)P(K) + P(D|\bar{K})P(\bar{K}) \\ &= (0.65)(0.3) + (0.9)(0.7) \\ &= 0.195 + 0.63 \\ &= 0.825 \end{aligned}$$

The probability that a randomly selected patient in the cancer ward is successfully treated with the drug regardless of kidney disease status is **0.825**, or **82.5%**.

The law of total probability can be generalized to a sample space partitioned into any number of subregions.

### Theorem 1.17.1: The generalized law of total probability.

Let a sample space  $S$  be partitioned into  $k$  subregions  $S_1, S_2, S_3, \dots, S_k$ . Let  $B$  be an event in  $S$ . Then

$$P(B) = P(B|S_1)P(S_1) + P(B|S_2)P(S_2) + P(B|S_3)P(S_3) + \dots + P(B|S_k)P(S_k)$$

### Example 1.17.5: The law of total probability.

At a certain university, the student body is composed of **24%** first-year students, **21%** second-year students, **27%** third-year students, and **28%** fourth-year students. **70%** of first-year students support increasing a student body association fee for improved gym facilities. **63%** of second-year students, **41%** of third-year students, and **51%** of fourth-year students support the fee increase.

What is the probability that a randomly selected student at the university supports the fee increase?

### Solution

Let  $S$  be the sample space of the entire student body. Partition  $S$  into the subevents  $Y_1$ ,  $Y_2$ ,  $Y_3$ , and  $Y_4$ , representing selecting a first-year, second-year, third-year, and fourth-year student, respectively. Let  $F$  be the event that a randomly selected student supports the fee increase. Thus,  $P(F)$  is the probability that a randomly selected student supports the fee increase.

The probability of selecting a first-year student is  $P(Y_1) = 24\% = 0.24$ . Similarly,  $P(Y_2) = 0.21$ ,  $P(Y_3) = 0.27$ , and  $P(Y_4) = 0.28$ .

©zyBooks 01/08/23 20:15 1267703  
Traver Yates

Given that a student is a first-year student, the probability that student supports the fee increase is 70%. Thus,  $P(F|Y_1) = 70\% = 0.70$ . Similarly,  $P(F|Y_2) = 0.63$ ,  $P(F|Y_3) = 0.41$ , and  $P(F|Y_4) = 0.51$ .

Finally, the law of total probability is used to find  $P(F)$ .

$$\begin{aligned} P(F) &= P(F|Y_1)P(Y_1) + P(F|Y_2)P(Y_2) + P(F|Y_3)P(Y_3) + P(F|Y_4)P(Y_4) \\ &= (0.70)(0.24) + (0.63)(0.21) + (0.41)(0.27) + (0.51)(0.28) \\ &= 0.168 + 0.132 + 0.111 + 0.143 \\ &= 0.554 \end{aligned}$$

The probability that a randomly selected student regardless of year supports the fee increase is 0.554, or 55.4%.

#### PARTICIPATION ACTIVITY

##### 1.17.5: The law of total probability.



A camera manufacturer is transitioning to a new model and is producing 20% new-model cameras and 80% old-model cameras. The manufacturer estimates that 1 in every 5000 new-model cameras are defective and 1 in every 2000 old-model cameras are defective.

Let  $D$  be the event that a randomly selected camera is defective and  $N$  be the event that a new-model camera is selected.

- 1) The manufacturer would like to determine the overall probability that a camera is defective regardless of model. Which probability should be calculated?



- $P(N)$
- $P(D)$
- $P(D|N)$

©zyBooks 01/08/23 20:15 1267703  
Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- 2) Which partition of the sample space should be used to find the probability that a camera is defective?



- $D$  and  $\bar{D}$
- $N$  and  $\bar{N}$
- $D$  and  $\bar{N}$

3) What is  $P(D|N)$ ? □

- 0.0002
- 0.0005
- 0.0000001

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

4) What is the probability that a camera is defective regardless of the model? □

- 0.0007
- 0.00026
- 0.00044

## References

(\*1) Weigold, Michelle. "A Rumpy Manx Cat" *Wikimedia Commons*. 13 June 2019, [https://commons.wikimedia.org/wiki/File:A\\_Rumpy\\_Manx\\_Cat.ipa](https://commons.wikimedia.org/wiki/File:A_Rumpy_Manx_Cat.ipa).

## 1.18 Bayes' Theorem



This section has been set as optional by your instructor.

### Bayes' Theorem

Sometimes, the converse of a conditional probability is of interest. Ex: A spam filter measures the frequency of known spam messages and known legitimate messages containing words from a list of words more likely to appear in spam. The converse probability is the likelihood that a message containing words from the spam list is actually spam. ©zyBooks 01/08/23 20:15 1267703  
Traver Yates

**Bayes' Theorem** relates the probability of an event  $A$  given a condition  $B$  to the probability of the condition  $B$  given that the event  $A$  occurred. That is, Bayes' Theorem allows  $P(B|A)$  to be calculated from  $P(A|B)$ . MAT-243-J3996-OL-TRAD-UG.23EW3

Theorem 1.18.1: Bayes' Theorem.

Let  $A$  and  $B$  be independent events in the same sample space and  $P(A) \neq 0$  and  $P(B) \neq 0$ . Then,

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\bar{A})P(\bar{A})}$$

### Proof.

By definition of conditional probability,

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$P(B|\bar{A}) = \frac{P(B \cap \bar{A})}{P(\bar{A})}$$

Rearranging gives

$$P(B \cap A) = P(B|A)P(A)$$

$$P(B \cap \bar{A}) = P(B|\bar{A})P(\bar{A})$$

By the law of total probability,

$$P(B) = P(B \cap A) + P(B \cap \bar{A})$$

Substituting the expressions for  $P(B \cap A)$  and  $P(B \cap \bar{A})$  into the expression for  $P(B)$  gives

$$P(B) = P(B|A)P(A) + P(B|\bar{A})P(\bar{A})$$

Substituting the expressions for  $P(B \cap A)$  and  $P(B)$  into the expression for  $P(A|B)$  gives

$$\begin{aligned} P(A|B) &= \frac{P(A \cap B)}{P(B)} \\ &= \frac{P(B \cap A)}{P(B)} \\ &= \frac{P(B|A)P(A)}{P(B)} \\ &= \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\bar{A})P(\bar{A})} \end{aligned}$$

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



## Animation content:

undefined

## Animation captions:

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

1. Let  $D$  be the event of having a disease.
2. Let  $T$  be the event of testing positive for the disease. Most people with the disease will test positive for the disease.
3. However, some people without the disease will also test positive.
4. Bayes' Theorem is used to find the probability of  $D$  given  $T$ .
5. The probability that a positive result actually indicates that a person has the disease is about **66.7%**.

## Example 1.18.1: Bayes' Theorem.

Suppose that 1 out of every 2000 computer chips produced by a chip manufacturer has a defect. The manufacturer has developed a method to test the chips, but the test is not perfect. If the chip has a defect, the test will correctly discover the defect with probability 0.98. If the chip does not have a defect, the test will incorrectly report that the chip has a defect with probability 0.01.

Let  $D$  be the event that a chip is defective. Let  $T$  be the event that a test indicates a defect. If a particular test indicates a defect, what is the probability that the chip is actually defective?

## Solution

The probability of interest is  $P(D|T)$ , the probability that a chip is actually defective given that the test indicated a defect.

Because 1 of every 2000 chips has a defect,  $P(D) = \frac{1}{2000} = 0.0005$ . Thus,

$P(\bar{D}) = 1 - P(D) = 0.9995$ . Given that a chip is defective, the test will indicate a defect with probability 0.98, so  $P(T|D) = 0.98$ . Given that a chip is not defective, the test will incorrectly indicate a defect with probability 0.01, so  $P(T|\bar{D}) = 0.01$ .

Bayes' Theorem is used to calculate  $P(D|T)$  from  $P(D)$ ,  $P(T|D)$ , and  $P(T|\bar{D})$ .

$$\begin{aligned}
 P(D|T) &= \frac{P(T|D)P(D)}{P(T|D)P(D) + P(T|\bar{D})P(\bar{D})} \\
 &= \frac{(0.98)(0.0005)}{(0.98)(0.0005) + (0.01)(0.9995)} \\
 &= \frac{0.00049}{0.00049 + 0.009995} \\
 &= \frac{0.00049}{0.010485} \\
 &\approx 0.047
 \end{aligned}$$

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Therefore, if a test indicates a defect, the probability that the chip is actually defective is **0.047**.

**PARTICIPATION ACTIVITY**
**1.18.2: Interpreting the probabilities in Bayes' Theorem.**


In a certain population, **1** in every **5000** people has a particular genetic disease. A genetic test for the disease has been developed but is prone to some error. The probability that the genetic test is incorrect when a person has the disease, that is, the genetic test has a negative result when the person has the disease, is **0.005**. The probability that the genetic test is incorrect when a person does not have the disease, that is, the genetic test has a positive result when the person does not have the disease, is **0.01**.

Let  **$D$**  be the event that the person has the disease and  **$T$**  be the event that the person tests positive for the disease. Match each probability to the correct description.

Select the definition that matches each term

1)  $P(D)$

- The probability that a randomly selected person in a population has the genetic disease
- The probability that the genetic test is incorrect when a person does not have the disease
- The probability that the genetic test is incorrect when a person has the disease
- The probability that the genetic test is correct when a person has the disease

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- The probability that a person who tests positive for the disease actually has the disease

2)  $P(D|T)$

- The probability that a randomly selected person in a population has the genetic disease
- The probability that the genetic test is incorrect when a person does not have the disease
- The probability that the genetic test is incorrect when a person has the disease
- The probability that the genetic test is correct when a person has the disease
- The probability that a person who tests positive for the disease actually has the disease

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

3)  $P(\bar{T}|D)$

- The probability that a randomly selected person in a population has the genetic disease
- The probability that the genetic test is incorrect when a person does not have the disease
- The probability that the genetic test is incorrect when a person has the disease
- The probability that the genetic test is correct when a person has the disease
- The probability that a person who tests positive for the disease actually has the disease

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

4)  $P(T|D)$

- The probability that a randomly selected person in a population has the genetic disease
- The probability that the genetic test is incorrect when a person does not have the disease
- The probability that the genetic test is incorrect when a person has the disease
- The probability that the genetic test is correct when a person has the disease
- The probability that a person who tests positive for the disease actually has the disease

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

5)  $P(T|\bar{D})$

- The probability that a randomly selected person in a population has the genetic disease
- The probability that the genetic test is incorrect when a person does not have the disease
- The probability that the genetic test is incorrect when a person has the disease
- The probability that the genetic test is correct when a person has the disease
- The probability that a person who tests positive for the disease actually has the disease

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

**PARTICIPATION ACTIVITY**

1.18.3: Applying Bayes' Theorem.



Use the information in the previous participation activity to answer the following.

1) What is  $P(D)$ ? Type as: #.####

**Check****Show answer**

2) What is  $P(\bar{D})$ ? Type as: #.####

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

3) What is  $P(\bar{T}|D)$ ? Type as: #.###

**Check****Show answer**

4) What is  $P(T|D)$ ? Type as: #.###

**Check****Show answer**

5) What is  $P(T|\bar{D})$ ? Type as: #.##

**Check****Show answer**

6) What is the probability that a person who tests positive for the disease actually has the disease?  
Type as: #.###

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## 1.19 Combinations and permutations



This section has been set as optional by your instructor.

## Multiplication rule

When the outcomes of an experiment all have the same probability of occurring, the probability of an event is computed as the number of outcomes in the event divided by the number of outcomes in the sample space. Determining the number of possible outcomes in both the event and sample space is not always simple. **Counting methods** or **combinatorics** are mathematical approaches to determining the number of possible outcomes.

The **multiplication rule of counting** states that if one event can occur in  $A$  ways and another event in  $B$  ways, the number of possible ways both events can occur is  $A \times B$ . Ex: If a car buyer can choose one of 3 possible colors and one of 4 possible entertainment systems, then  $3 \times 4 = 12$  different car configurations are possible. If a third event occurs  $C$  ways, the number of possibilities is computed  $A \times B \times C$ . Ex: If the buyer may also pick one of 2 seating arrangements,  $3 \times 4 \times 2 = 24$  different car configurations are possible. Additional events are similarly computed.

### PARTICIPATION ACTIVITY

1.19.1: Multiplication rule.



## Animation captions:

1. For each of 3 sandwich choices, a customer can select one of two soups.
2. 6 possible sandwich and soup choices exist.
3. The multiplication rule for the possible sandwich and soup choices yields:  $3 \times 2 = 6$ .
4. The multiplication rule for possible meals includes 4 dessert choices, so:  $3 \times 2 \times 4 = 24$  meal choices exist in total.

### PARTICIPATION ACTIVITY

1.19.2: Multiplication rule of counting.



A new interdisciplinary data analytics major requires students to choose 3 courses, by selecting one course from each of 3 departments. The 3 selected courses are called a "curriculum plan". The Psychology department offers 2 possible courses, Computer Science offers 4, and Math offers 5.

- 1) How many different curriculum plans for the three courses are possible?

- 3
- 11
- 40

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EV

3



2) The Computer Science department

proposes a 5th course. How many  
more curriculum plans are now  
possible?

- 10
- 48
- 50

3) All 3 departments propose 1 new  
course each. The program director  
wants to add only one new course.  
Which department should the director  
increase to produce the most possible  
curriculum plans?

- Computer Science
- Math
- Psychology

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



## Permutations and factorials

A **permutation** is a possible ordering of a set of objects. Determining the number of possible permutations is possible using the multiplication rule. Ex: Joshua is planning a violin concert and has 5 pieces of music to play in any order. The number of possible concert programs is a permutation. The first piece in the concert is selected from 5 choices. Once the first piece is chosen, 4 choices remain for the second. Continuing with the same logic, the multiplication rule results in  $5 \times 4 \times 3 \times 2 \times 1$  possible concert programs. The mathematical expression for the number of permutations of  $n$  objects is known as a factorial. A **factorial**, denoted  $n!$ , is the multiplication of a non-negative integer  $n$  and all non-negative integers less than  $n$ ;  $n! = n \times (n - 1) \times \cdots \times 2 \times 1$ .  $0!$  is defined to be 1.

Python-Function 1.19.1: `math.factorial()`.

The `math.factorial()` function returns the factorial of a non-negative integer. An error results if the argument is not an integer or is negative. The built-in `math` library must be imported to use the function.

```
import math

# Find 5!
print(math.factorial(5))

# Find 0!
print(math.factorial(0))
```

120  
1

[Run example](#)

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

PARTICIPATION ACTIVITY

1.19.3: Permutations of a concert program.



### Animation captions:

1. The multiplication rule can compute the number of possible concert programs (permutations) for **5** songs. **5** choices exist for the first song.
2. **4** choices are available for the second piece once the opening number is selected.  
 $5 \times 4 = 20$  sets of the first two pieces are possible.
3. **3** choices are available for the third piece once the first two numbers are selected.
4.  $5! = 5 \times 4 \times 3 \times 2 \times 1 = 120$  different concert orders, or permutations, are possible with five musical pieces.

PARTICIPATION ACTIVITY

1.19.4: Permutations and factorials.



- 1) A magician asks an audience member to shuffle **6** objects. What is the number of possible orderings, or permutations?

- 120
- 720
- 46,656

- 2) The magician of the previous question is not able to see the order of the **6** objects. What is the probability the magician correctly guesses the order by chance?

- $\frac{1}{2}$
- $\frac{1}{6}$

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

$\frac{1}{720}$

- 3) In the concert program animation above, the number of possible concert programs with **5** pieces of music is computed. Joshua hopes to perform well and prepares **3** additional pieces as a possible encore. He will choose just one encore. Including the encore, how many possible programs are possible, for the **5** main songs plus the **1** encore?

- 360
- 720
- 40,320

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

## Computing the number of permutations

A form of permutation that often occurs, sometimes referred to as a sub-permutation, is when not all the possible objects are included in the ordering. Ex: Joshua will choose only **3** of the **5** pieces for the concert. The formula for the number of ways to produce ordered sets of size  $k$  from a collection of  $n$  objects is

$$P_{n,k} = \frac{n!}{(n-k)!}$$

Ex: The number of arrangements of size  $k = 5$  for  $n = 5$  pieces of music is

$$P_{5,5} = \frac{5!}{(5-5)!} = \frac{5!}{0!} = 5!. \text{ The number of arrangements of } k = 3 \text{ pieces selected from } n = 5 \text{ is}$$

$$P_{3,5} = \frac{5!}{(5-3)!} = \frac{5!}{2!} = \frac{5 \times 4 \times 3 \times 2 \times 1}{2 \times 1} = 5 \times 4 \times 3 = 60.$$

PARTICIPATION ACTIVITY

1.19.5: Permutation formula.



## Animation captions:

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

1. The formula for the number of permutations computes the number of program orders of **3** pieces of music selected from **5** choices.
2. With only **3** available positions, the number of permutations are based on the first **3** numbers of the factorial.
3. The permutation formula removes the last two terms of the factorial.

**PARTICIPATION  
ACTIVITY**

1.19.6: Computing the number of permutations.



8 runners are competing in the Olympic 100-meter dash final.

- 1) What is the number of possible winners of the race?

  
/ /**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



- 2) What is the number of possible orders of finish for all 8 runners?

  
/ /**Check****Show answer**

- 3) What is the number of possible sets of gold, silver, and bronze medalists, that is, the number of possible orderings of the first 3 finishers in the race?

  
/ /**Check****Show answer**

- 4) What is the number of possible orderings of the bottom 4 finishers in the race?

  
/ /**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



## Combinations

A permutation places selected objects in order. Sometimes just being a selected object is important, regardless of order. Ex: In the Olympics, the order of the first three finishers of the 100-meter race is important, earning **3** different medals: Gold, Silver, and Bronze. But in contrast, in the U.S. Olympic trials, the top **3** finishers make the Olympic team, and the order doesn't matter. A **combination** is a possible set of objects chosen in any order; the order of selection in a combination does not matter.

The formula for the number of ways to select unordered sets of size  $k$  from a collection of  $n$  objects is  $C_{n,k} = \frac{n!}{k!(n-k)!}$ .  $C_{n,k}$  is read "n choose k"; since order does not matter, the expression counts the number of ways to simply "choose" a set of  $k$  objects at once from a group of  $n$  objects. The difference in the formula for number of combinations and that for permutations above is the term  $k!$  in the denominator. The term divides out the permutations where the same set of objects are listed but in different order. The result is that the number of combinations is always less than or equal to the number of possible permutations.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

**PARTICIPATION ACTIVITY**

1.19.7: Combinations.

**Animation captions:**

1. 12 possible permutations of the first place and repechage boats are possible.
2. If, instead of a repechage race, both of the top two boats qualify, then the order of the top two does not matter.
3. The possible combinations of qualified boats, ignoring duplicate orders, is 6.
4. Many of the ordered pairs are removed since they are duplicated if order does not matter.

**PARTICIPATION ACTIVITY**

1.19.8: Computing the number of combinations.



40 runners compete in the Olympic 400-meter race. After a series of qualifying heats, 10 of the 40 runners compete in the final race for a gold, a silver, and a bronze medal.

- 1) What is the number of possible sets of gold, silver, and bronze medalists in the final race?

**Check****Show answer**

- 2) What is the number of possible sets of any three medalists in the final race?

**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

- 3) What is the number of possible sets of any three medalists based



on all runners in the competition?



- 4) What is the number of possible sets of runners in the final race based on all runners in the competition?



©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

## Probability using combinatorics

Permutations, combinations, factorials, and the multiplication rule of counting can be used to count the sizes of an event and a sample space that are too large to count explicitly.

Example 1.19.1: Probability using permutations and the multiplication rule of counting.

A security code for unlocking a smartphone is 6 digits long. What is the probability of selecting a code with no repeated digits?

### Solution

The sample space of selecting a 6-digit code is

$10 \cdot 10 \cdot 10 \cdot 10 \cdot 10 \cdot 10 = 10^6 = 1,000,000$ . The number of 6-digit codes with no repeated digits is  $P_{10,6}$ , as the order of the digits matter.

$$\begin{aligned} P_{10,6} &= \frac{10!}{(10-6)!} \\ &= \frac{10!}{4!} \\ &= 10 \cdot 9 \cdot 8 \cdot 7 \cdot 6 \cdot 5 \\ &= 151,200 \end{aligned}$$

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

Let  $N$  be the event of selecting a code with no repeated digits. Thus,  $P(N)$  is the number of outcomes in  $N$  divided by the total number of outcomes in the sample space.

$$\begin{aligned} P(N) &= \frac{151,200}{1,000,000} \\ &= 0.151 \end{aligned}$$

The probability of selecting a 6-digit code with no repeated digits is **0.151**.

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

### Example 1.19.2: Probability using combinations.

In poker, a 5-card hand has two pairs if the hand has two cards of the same rank, two cards of the same rank but different from the rank of the first two cards, and one card of a different rank from the other four cards. Ex:

$3\heartsuit, 3\spadesuit, 7\clubsuit, 7\clubsuit, J\clubsuit$  is a hand with two pairs

$3\heartsuit, 3\spadesuit, 3\clubsuit, 3\clubsuit, J\clubsuit$  is not a hand with two pairs. (This is a four-of-a-kind.)

$3\heartsuit, 3\spadesuit, 7\clubsuit, 7\spadesuit, 7\heartsuit$  is not a hand with two pairs. (This is a full house.)

A dealer at a casino deals a 5-card hand from a standard deck of 52 cards. What is the probability that the hand has two pairs?

### Solution

The sample space of selecting a hand of 5 cards from a 52-card deck without regard to order is  $C_{52,5}$ .

$$\begin{aligned} C_{52,5} &= \frac{52!}{(52-5)!5!} \\ &= \frac{52!}{47!5!} \\ &= \frac{52 \cdot 51 \cdot 50 \cdot 49 \cdot 48}{5 \cdot 4 \cdot 3 \cdot 2 \cdot 1} \\ &= 2,598,960 \end{aligned}$$

From the 13 ranks, 2 ranks are chosen as the two ranks that are paired. From the 4 suits, 2 are chosen as the suits in each pair. The last card can be any of the remaining cards that is not of the same rank of the other four cards, or  $52 - 4 - 4 = 44$ .

Let  $T$  be the event of dealing a hand with two pairs. By the multiplication rule of counting, the number of outcomes in  $T$  is

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3

$$\begin{aligned}
 C_{13,2}C_{4,2}C_{4,2} \cdot 44 &= \frac{13!}{(13-2)!2!} \cdot \frac{4!}{(4-2)!2!} \cdot \frac{4!}{(4-2)!2!} \cdot 44 \\
 &= \frac{13!}{11!2!} \cdot \frac{4!}{2!2!} \cdot \frac{4!}{2!2!} \cdot 44 \\
 &= 78 \cdot 6 \cdot 6 \cdot 44 \\
 &= 123,552
 \end{aligned}$$

Thus,  $P(T)$  is the number of outcomes in  $T$  divided by the total number of outcomes in the sample space.

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3

$$\begin{aligned}
 P(T) &= \frac{123,552}{2,598,960} \\
 &\approx 0.048
 \end{aligned}$$

The probability that the dealer deals a 5-card hand with two pairs is approximately 0.048.

**PARTICIPATION ACTIVITY**

## 1.19.9: Probability using combinatorics.



In poker, a 5-card hand is a flush if all five cards are of the same suit. (For simplicity, assume that straight flushes and royal flushes are also flushes. A straight flush is a flush with cards of consecutive ranks, ex: 2 ♠, 3 ♠, 4 ♠, 5 ♠, 6 ♠. A royal flush is a straight flush with the ace as the high card, ex: 10 ♠, J ♠, Q ♠, K ♠, A ♠.)

A dealer at a casino deals a 5-card hand from a standard deck of 52 cards.

- 1) What is the size of the sample space?

**Check**
**Show answer**


- 2) Let  $F$  be the event of dealing a flush. What is the number of outcomes in  $F$ ?

**Check**
**Show answer**


- 3) What is  $P(F)$ ? Type as: #.###

©zyBooks 01/08/23 20:15 1267703  
Traver Yates  
MAT-243-J3996-OL-TRAD-UG.23EW3



**Check****Show answer**

©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3



©zyBooks 01/08/23 20:15 1267703

Traver Yates

MAT-243-J3996-OL-TRAD-UG.23EW3