

Martian Alphabet

1.1 Learning outcomes

- Describe the statistical investigation process
- Identify observational units, variables, and variable types in a statistical study

1.2 Terminology review

Statistics is the study of how best to collect, analyze, and draw conclusions from data. Today in class you will be introduced to the following terms:

- Observational units or cases
- Variables: categorical or quantitative
- Proportions
- Graphs: frequency bar plot and relative frequency bar plot
- Distribution

For more on these concepts, read Sections 1.2 and 2.1 in the textbook.

1.3 Can you read “Martian”?

How well can humans distinguish one “Martian” letter from another? In today’s activity, we’ll find out. When shown the two Martian letters, Kiki and Bumba, write down whether you think Bumba is on the left or on the right.

Steps of the statistical investigation process

Step 1: The first step of any statistical investigation is to *ask a research question*. In this study the research question is: can we as a class read Martian? (We will refine this later on!).

Step 2: To answer any research question, we must *design a study and collect data*. For our question, the study consists of each student being presented with two Martian letters and asking which was Bumba. Your responses will become our observed data that we will explore.

Observational units or **cases** are the subjects data are collected on. In a spreadsheet of the data set, each row will represent a single observational unit.

1. What are the observational units in this study?
2. How many students are in class today? This is the *sample size*.

A **variable** is information collected or measured on each observational unit or case. Each column in a data set will represent a different variable. Today we are only measuring one variable on each observational unit.

3. Identify the variable we are collecting on each observational unit in this study, i.e., what are we measuring on each student?

We will look at two types of variables: **quantitative** and **categorical** (see Figure 1.1).

Quantitative variables are numerical measurements that can be discrete (whole, non-negative numbers) or continuous (any value within an interval). The number of students in a class would be a discrete variable as you can not have a partial student. GPA would be a continuous variable ranging from 0 to 4.0.

Categorical variables are data that are in groups or categories such as eye color, state of residency, or whether or not a student lives on campus. Categorical variables with a natural ordering are considered ordinal variables while those without a natural ordering are considered a nominal variable. All categorical variables will be treated as nominal for analysis in this course.

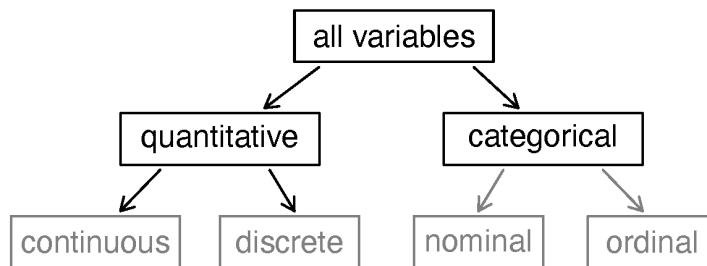


Figure 1.1: Types of variables.

4. Is the variable identified in question 3 categorical or quantitative?
5. Were you correct or incorrect in identifying Bumba?

Step 3: Once we have collected data, the next step is to *summarize and visualize the data*.

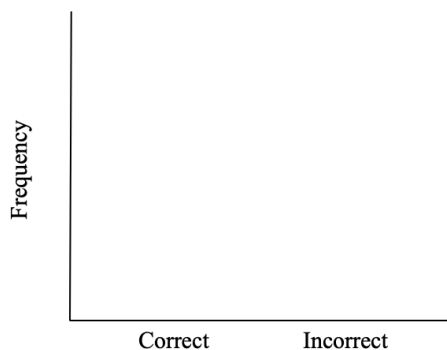
6. How many people in your class were correct in identifying Bumba? Using the class size from question 2, calculate the proportion of students who correctly identified Bumba.

$$\text{proportion} = \frac{\text{number of students who correctly identified Bumba}}{\text{total number of students}}$$

The proportion in question 6 is called a **summary statistic**—a single value that summarizes the data set. It is important to note that a variable is different than a summary statistic. A *variable* is measured on a *single observational unit* while a summary statistic is calculated from a group of observational units. For example, the variable “whether or not a student lives on campus” can be measured on each individual student. In a class of 50 students we can calculate the proportion of students who live on campus, the summary statistic. Look back and make sure you wrote the variable in question 3 as a variable, NOT a summary statistic.

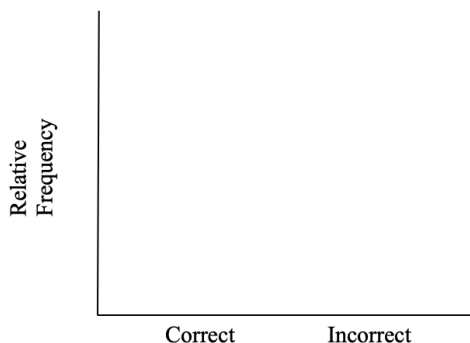
Looking at the data set and the summary statistics is only one way to display the data. We will also want to create a visualization or picture of the data. A **frequency bar plot** is used to display categorical data as a count or frequency. Since our variable has two levels, correct or incorrect, we will create two bars—one for each level.

7. Plot the observed class data using a frequency bar plot. Be sure to add a scale to the y-axis.



We can also visualize the data as a proportion in a **relative frequency bar plot**. Relative frequency is the proportion calculated for each level of the categorical variable.

8. Plot the observed class data using a relative frequency bar plot. Be sure to add a scale to the y-axis.



Step 4: The next step is to *use statistical analysis methods to draw inferences from the data*. To answer the research question, we will simulate what *could* have happened in our class given random chance, repeat that many times to understand the expected *variability* between different “randomly guessing” classes, then compare our class’s observed data to the simulation. This gives us an estimate of how often (or the probability of) our class’s result would occur if we were all merely guessing, allowing us to determine if the data provides evidence that we as a class can in fact read Martian.

9. If humans really don’t know Martian and are just guessing which is Bumba, what are the chances of getting it right?

How could we use a coin to simulate each student “just guessing” which Martian letter is Bumba?

How could we use coins to simulate the entire class “just guessing” which Martian letter is Bumba?

How many people in your class would you expect to choose Bumba correctly just by chance? Explain your reasoning.

10. Each of you will flip a coin one time to simulate your “guess”. Let Heads = correct, Tails = incorrect. What was the result of your simulation?

What was the result from your class’s simulation? What proportion of students “guessed” correctly in the simulation?

11. If students really don’t know Martian and are just guessing which is Bumba, which seems more unusual: the result from your class’s **simulation** or the observed proportion of students in your class that were correct (this is your summary statistic from question 6)? Explain your reasoning.

12. While your observed class data is likely far different from the simulated “just-guessing” class, comparing our class data to a single simulation does not seem to give enough information. The differences seen could just be due to that set of coin flips! Let’s simulate another class. Each student should flip your coin again. What was the result from your class’s second simulation? What proportion of students “guessed” correctly in the second simulation? Create a plot to compare the two simulated results with the observed class result.
13. We still unfortunately only have a couple of simulations to compare our class data to. It would be much better to be able to see how our class compared to hundreds or thousands of “just-guessing” classes. Since we don’t want to flip coins all class period, your instructor will use a computer simulation to get 1000 trials. Fill in the following blanks to describe how we would create a simulation of random guessing with 1000 trials.
- Probability of correct guesses: _____
- Sample size: _____
- Number of repetitions: _____
14. Sketch the distribution displayed by your instructor here, being sure to label each axis appropriately.
15. Is your class particularly good or bad at Martian? How can you use the plot in question 14 to tell?
16. Is it *possible* that we could see our class results just by chance if everyone was just guessing? Explain your reasoning.
17. Is it *likely* that we could see our class results just by chance if everyone was just guessing? Explain your reasoning.

Step 5: The next step in the statistical investigation process is to *communicate the results and answer the research question*.

18. Does this activity provide strong evidence that students were not just guessing at random? If so, what do you think is going on here? Can we as a class read Martian?¹

Step 6: The final step of any statistical investigation is to *revisit and look ahead*.

19. Can you think of any limitations of our study? Can you think of a new topic that might be of interest based on the results of our study?

1.4 Take home messages

1. In this course we will learn how to evaluate a claim by comparing observed results (classes' "guesses" when asked to identify Bumba) to a distribution of many simulated results under an assumption like "blind guessing."
2. Blind guessing between two outcomes will be correct only about half the time. We can create data (via computer simulation) to fit the assumption of blind guessing.
3. Unusual observed results will make us doubt the assumptions used to create the simulated distribution. A large number of correct "guesses" is evidence that a person was not just blindly guessing.

1.5 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity, and to write down the names and contact information of your team mates.

¹Reference for "Martian alphabet" is a TED talk given by Vilayanur Ramachandran in 2007. The synesthesia part begins at roughly 17:30 minutes: http://www.ted.com/talks/vilayanur_ramachandran_on_your_mind.