

5.7 FURTHER DISCOVERY

This chapter's epigraph is a famous "Yogiism," from Hall of Fame catcher, coach, and manager Yogi Berra [6].

If you would like to learn more about Robert Brown's experiments, and the history and science behind them, visit the following web site, titled "What Brown Saw and You Can Too."

<http://physerver.hamilton.edu/Research/Brownian/index.html>

The Drunkard's Walk by Leonard Mlodinow [34] is a very accessible book about how randomness and chance affect our lives. For more information about generating random numbers, and the differences between PRNGs and true random number generators, visit

<https://www.random.org/randomness/> .

The Park-Miller random number generator is due to Keith Miller and the late Steve Park [37].

The Roper Center for Public Opinion Research, at the University of Connecticut, maintains some helpful educational resources about random sampling and errors in the context of public opinion polling at

<http://www.ropercenter.uconn.edu/education.html> .

5.8 PROJECTS

Project 5.1 The magic of polling

According to the latest poll, the president's job approval rating is at 45%, with a margin of error of $\pm 3\%$, based on interviews with approximately 1,500 adults over the weekend.

We see news headlines like this all the time. But how can a poll of 1,500 randomly chosen people claim to represent the opinions of millions in the general population? How can the pollsters be so certain of the margin of error? In this project, we will investigate how well random sampling can really estimate the characteristics of a larger population. We will assume that we know the true percentage of the overall population with some characteristic or belief, and then investigate how accurate a much smaller poll is likely to get.

Suppose we know that 30% of the national population agrees with the statement, "Animals should be afforded the same rights as human beings." Intuitively, this means that, if we randomly sample ten individuals from this population, we should, on average, find that three of them agree with the statement and seven do not. But does it follow that **every** poll of ten randomly chosen people will mirror the percentage of the larger **population**? Unlike a Monte Carlo simulation, a poll is taken just once (or maybe twice) at any particular point in time. To have confidence

in the poll results, we need some assurance that the results would not be drastically different if the poll had queried a different group of randomly chosen individuals. For example, suppose you polled ten people and found that two agreed with the statement, then polled ten more people and found that seven agreed, and then polled ten more people and found that all ten agreed. What would you conclude? There is too much variation for this poll to be credible. But what if we polled more than ten people? Does the variation, and hence the trustworthiness, improve?

In this project, you will write a program to investigate questions such as these and determine empirically how large a sample needs to be to reliably represent the sentiments of a large population.

1. Simulate a poll

In conducting this poll, the pollster asks each randomly selected individual whether he or she agrees with the statement. We know that 30% of the population does, so there is a 30% chance that each individual answers “yes.” To simulate this polling process, we can iterate over the number of individuals being polled and count them as a “yes” with probability 0.3. The final count at the end of the loop, divided by the number of polled individuals, gives us the poll result. Implement this simulation by writing a function

```
poll(percentage, pollSize)
```

that simulates the polling of `pollSize` individuals from a large population in which the given `percentage` (between 0 and 100) will respond “yes.” The function should return the percentage (between 0 and 100) of the poll that actually responded “yes.” Remember that the result will be different every time the function is called. Test your function with a variety of poll sizes.

2. Find the polling extremes

To investigate how much variation there can be in a poll of a particular size, write a function

```
pollExtremes(percentage, pollSize, trials)
```

that builds a list of `trials` poll results by calling `poll(percentage, pollSize)` `trials` times. The function should return the minimum and maximum percentages in this list. For example, if five trials give the percentages [28, 35, 31, 24, 31], the function should return the minimum 24 and maximum 35. If the list of poll results is named `pollResults`, you can return these two values with

```
return min(pollResults), max(pollResults)
```

Test your function with a variety of poll sizes and numbers of trials.

3. What is a sufficient poll size?

Next, we want to use your previous functions to investigate how increasing poll sizes affect the variation of the poll results. Intuitively, the more people you poll, the more accurate the results should be. Write a function

```
plotResults(percentage, minPollSize, maxPollSize, step, trials)
```

that plots the minimum and maximum percentages returned by calling the function `pollExtremes(percentage, pollSize, trials)` for values of `pollSize` ranging from `minPollSize` to `maxPollSize`, in increments of `step`. For each poll size, call your `pollExtremes` function with

```
low, high = pollExtremes(percentage, pollSize, trials)
```

and then append the values of `low` and `high` each to its own list for the plot. Your function should return the margin of error for the largest poll, defined to be $(\text{high} - \text{low}) / 2$. The poll size should be on the x -axis of your plot and the percentage should be on the y -axis. Be sure to label both axes.

Question 5.1.1 *Assuming that you want to balance a low margin of error with the labor involved in polling more people, what is a reasonable poll size? What margin of error does this poll size give?*

Write a main function (if you have not already) that calls your `plotResults` function to investigate an answer to this question. You might start by calling it with `plotResults(30, 10, 1000, 10, 100)`.

4. Does the error depend on the actual percentage?

To investigate this question, write another function

```
plotErrors(pollSize, minPercentage, maxPercentage, step, trials)
```

that plots the margin of error in a poll of `pollSize` individuals, for actual percentages ranging from `minPercentage` to `maxPercentage`, in increments of `step`. To find the margin of error for each poll, call the `pollExtremes` function as above, and compute $(\text{high} - \text{low}) / 2$. In your plot, the percentage should be on the x -axis and the margin of error should be on the y -axis. Be sure to label both axes.

Question 5.1.2 *Does your answer to the previous part change if the actual percentage of the population is very low or very high?*

You might start to investigate this question by calling the function with `plotErrors(1500, 10, 80, 1, 100)`.