

Focus on integrity

Data integrity and analytics objectives

Overcoming the challenges of insufficient data

▶ **Video:** Dealing with insufficient data  
3 min

📖 **Reading:** What to do when you find an issue with your data  
10 min

▶ **Video:** The importance of sample size  
3 min

📖 **Reading:** Calculating sample size  
20 min

📖 **Practice Quiz:** Self-Reflection: Why pre-cleaning activities are important  
1 question

📖 **Practice Quiz:** Test your knowledge on insufficient data  
3 questions

Testing your data

Consider the margin of error

Weekly challenge 1

# Calculating sample size

Before you dig deeper into sample size, familiarize yourself with these terms and definitions:

| Terminology              | Definitions  |
|--------------------------|--|
| Population               | The entire group that you are interested in for your study. For example, if you are surveying people in your company, the population would be all the employees in your company.   |
| Sample                   | A subset of your population. Just like a food sample, it is called a sample because it is only a taste. So if your company is too large to survey every individual, you can survey a representative sample of your population.   |
| Margin of error          | Since a sample is used to represent a population, the sample's results are expected to differ from what the result would have been if you had surveyed the entire population. This difference is called the margin of error. The smaller the margin of error, the closer the results of the sample are to what the result would have been if you had surveyed the entire population. |
| Confidence level         | How confident you are in the survey results. For example, a 95% confidence level means that if you were to run the same survey 100 times, you would get similar results 95 of those 100 times. Confidence level is targeted before you start your study because it will affect how big your margin of error is at the end of your study.   |
| Confidence interval      | The range of possible values that the population's result would be at the confidence level of the study. This range is the sample result +/- the margin of error.  |
| Statistical significance | The determination of whether your result could be due to random chance or not. The greater the significance, the less due to chance.   |

## Things to remember when determining the size of your sample

When figuring out a sample size, here are things to keep in mind:

- Don't use a sample size less than 30. It has been statistically proven that 30 is the smallest sample size where an average result of a sample starts to represent the average result of a population.
- The confidence level most commonly used is 95%, but 90% can work in some cases.

Increase the sample size to meet specific needs of your project:

- For a **higher** confidence level, use a larger sample size
- To **decrease** the margin of error, use a larger sample size
- For **greater** statistical significance, use a larger sample size

**Note:** Sample size calculators use statistical formulas to determine a sample size. More about these are coming up in the course! Stay tuned.

### Why a minimum sample of 30?

This recommendation is based on the **Central Limit Theorem (CLT)** in the field of probability and statistics. As sample size increases, the results more closely resemble the normal (bell-shaped) distribution from a large number of samples. A sample of 30 is the smallest sample size for which the CLT is still valid. Researchers who rely on **regression analysis** – statistical methods to determine the relationships between controlled and dependent variables – also prefer a minimum sample of 30.

Still curious? Without getting too much into the math, check out these articles:

- [Central Limit Theorem \(CLT\)](#) [🔗](#): This article by Investopedia explains the Central Limit Theorem and briefly describes how it can apply to an analysis of a stock index.
- [Sample Size Formula](#) [🔗](#): This article by Statistics Solutions provides a little more detail about why some researchers use 30 as a minimum sample size.

## Sample sizes vary by business problem

Sample size will vary based on the type of business problem you are trying to solve.

For example, if you live in a city with a population of 200,000 and get 180,000 people to respond to a survey, that is a large sample size. But without actually doing that, what would an acceptable, smaller sample size look like?

Would 200 be alright if the people surveyed represented every district in the city?

**Answer:** It depends on the stakes.

- A sample size of 200 might be large enough if your business problem is to find out how residents felt about the new library
- A sample size of 200 might not be large enough if your business problem is to determine how residents would vote to fund the library

You could probably accept a larger margin of error surveying how residents feel about the new library versus surveying residents about how they would vote to fund it. For that reason, you would most likely use a larger sample size for the voter survey.



## Larger sample sizes have a higher cost

You also have to weigh the cost against the benefits of more accurate results with a larger sample size. Someone who is trying to understand consumer preferences for a new line of products wouldn't need as large a sample size as someone who is trying to understand the effects of a new drug. For drug safety, the benefits outweigh the cost of using a larger sample size. But for consumer preferences, a smaller sample size at a lower cost could provide good enough results.



## Knowing the basics is helpful

Knowing the basics will help you make the right choices when it comes to sample size. You can always raise concerns if you come across a sample size that is too small. A sample size calculator is also a great tool for this. Sample size calculators let you enter a desired confidence level and margin of error for a given population size. They then calculate the sample size needed to statistically achieve those results.

Refer to the [Determine the Best Sample Size](#) [🔗](#) video for a demonstration of a sample size calculator, or refer to the [Sample Size Calculator](#) [🔗](#) reading for additional information.



Mark as completed