

# Introduction to confidence intervals

Millie Harris

## Summary

In statistics, to estimate an unknown parameter you can construct a confidence interval. This is the range of values you expect the true estimate to fall between if you were to repeat the study several times, with a certain level of confidence. This study guide introduces confidence intervals, confidence levels, and  $Z$ -values using the normal distribution.

*Before reading this guide, it is recommended that you read [Guide: Introduction to probability](#), [Guide: Introduction to hypothesis testing](#), and [Guide: Expected value, variance, standard deviation](#).*

## What is a confidence interval?

If you were conducting a study about a population of people and took several different samples of data, the sample means could be different for each sample. This makes estimating a population mean very difficult, as this variability affects your confidence that the sample is a true reflection of the population. A population mean is a fixed, unknown constant, and providing a confident estimate of this quantity that cannot be measured is one of the fundamental goals of statistics.

### Initial example and set up

For instance, suppose that you were investigating average weights of chocolate bars. To do this, you had taken five samples of 30 chocolate bars from the massive conglomerate that is Cantor's Confectionery, and worked out the average weights of each of these five samples:

Sample number	$S_m$	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$
Sample mean $\bar{x}_m$		28.2	29.6	27.9	29.1	28.8

Which of these is most reflective of the average weight of all of the chocolate bars produced by Cantor's Confectionery? It's impossible to say! So what you can do is consider a **range of values**. This range of values is called a **confidence interval (CI)**, which is specified up to

a **confidence level** (usually 95%, but this can vary). This range of values is centred at the sample mean  $\bar{x}$ , and the bounds are given by the **margin of error** of your sample.

So each of these sample means will have a confidence interval associated to it, which is worked out using the data from the sample. You can then expand the above table to include the confidence intervals:

Sample number $S_m$	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$
Sample mean $\bar{x}_m$	28.2	29.6	27.9	29.1	28.8
95% CI for sample $m$	[27.6, 28.8]	[29.1, 30.1]	[27.6, 28.2]	[28.3, 29.9]	[28.1, 29.5]

What does a confidence interval do? It should give you an idea of where the population mean should be. But given the various confidence intervals for the samples above, it's not even clear where the population mean could be in any one of these. For instance, the confidence intervals for samples 1 and 2 do not meet; the population mean cannot be in both! This is because it is a fixed, unknown constant which does not vary.

It's important to note that there is one confidence interval **per sample** of a population, and so if you take many samples (as has been done above), there are many confidence intervals that each estimate the population mean  $\mu$ .

### What a confidence interval tells you

Suppose that CL is the confidence level of your confidence interval. The general rule is that for CL% confidence intervals, then **if you sample the population 100 times, then you can expect the population mean to lie within approximately CL many of the calculated CL % CIs.** You can replace CL with your given confidence level.

So you can say

It is likely (with CL% confidence) that the population mean is within this confidence interval.

because the population mean will be in approximately CL/100 of all CL% CI's, but you **cannot say**

For *this particular* confidence interval, the probability that the population mean is within this confidence interval is CL%.

This is because the population mean  $\mu$  is a **fixed, unknown constant**; and so the probability that it is in a given CL% CI is either 0 (it's not in there) or 1 (it's in there). Note the important difference between 'probability' and 'confidence'!

The concept that does this is known as a **credible interval**; see [Guide: Credible intervals] for more.

### What's the point?

Confidence intervals are a vital tool used to measure uncertainties in everyday life. For example:

- In politics, confidence intervals can be used to show the uncertainty in polling estimates.
- In economics, confidence intervals can be used to show uncertainties in market trends and inflation.
- In medicine and biology, confidence intervals can be used to show uncertainties around effects like mean weight loss, drug effectiveness, or survival rates.
- In sports, confidence intervals can be used by coaches to measure the true performance levels of athletes.

### In this guide

This guide will focus on how to construct and interpret confidence intervals using the **normal distribution only**. This will include looking at  $Z$ -values, two-tailed alpha values, and confidence levels. Then, the guide will discuss how to interpret what your confidence interval means.

For information on confidence intervals using other distributions see [Guide: More on confidence intervals].

## Constructing confidence intervals

### The normal distribution

The **normal distribution**, or sometimes called the "bell-curve" because of its shape, is a function used to model the probability of various naturally occurring measurements, such as the average height of trees, average ages of cats, average IQs of humans, and so on. You would expect these measurements to have lots of measurements close to the middle (the average) and fewer measurements towards the extremes (the tails).

More mathematically, the function is called a **probability density function** and it is written as  $N(\mu, \sigma^2)$ , which depends on two parameters: the **mean** ( $\mu$ ) and the **standard deviation** ( $\sigma$ ). In situations where  $\mu = 0$  and  $\sigma = 1$ , you have a **standard normal distribution**. For more on  $\mu$  and  $\sigma$  see [Guide: Expected value, variance, standard deviation](#); for more on probability density functions, please see [Guide: PMFs, PDFs, CDFs](#).

**i** Definition of a normally distributed random variable

For a population where  $\mu$  is the population mean and  $\sigma$  is the population standard deviation, a random variable  $X$  which is modelled by the normal distribution with these parameters is written as

$$X \sim N(\mu, \sigma^2)$$

which you can read as 'X is normally distributed with parameters  $\mu$  and  $\sigma$ '.

So to use the mathematical tools associated with the normal distribution on a population, you need to have a fairly good idea of the two parameters, the mean  $\mu$  (mu) and standard deviation  $\sigma$  (sigma).

For more information on the normal distribution (such as what the function  $N(\mu, \sigma^2)$  actually is) see [Factsheet: Normal distribution](#).

## Tails

Because of the 'bell-shape' of the normal distribution, there are less values at each end. The extremes of the normal distribution are called the **tails**.

By drawing two vertical lines on the  $x$ -axis at the two points  $x = x_0$  and  $x = -x_0$  (where  $x_0$  is some number), and drawing upwards, you can define the area under the curve bounded to the right of  $x = x_0$  and to the left of  $x = -x_0$ . These areas at each end of the curve are often also called the **tails** of the distribution.

The total area of both tails combined is called  $\alpha$  (alpha), and this is always a number between 0 and 1. Since the graph is symmetric and the lines are drawn at  $x = x_0$  and  $x = -x_0$ , the area of each tail is exactly  $\alpha/2$ .

To construct a confidence interval you need both tails, because you are looking at values both above and below the population mean, which is unknown. So you will construct what is called a **two-tailed test**.

## Confidence level

### **i** Definition of a confidence level

A **confidence level** (CL) suggests that if you were to repeat the sample and construction of a confidence interval 100 times, you would expect the true value of the population mean to fall within CL of the constructed confidence intervals. A CL is typically represented using a percentage.

For example, a 95% CL suggests that the true value of the population mean would fall within 95 out of 100 computed confidence intervals.

### **!** $\alpha$ is 1 minus the CL

$1 - \alpha$  is the confidence level. So you only need **one** of either  $\alpha$  or the confidence level in order to generate a confidence interval.

## Z value

The value of  $\alpha$  (and/or the CL) is decided before constructing the confidence interval. Every value of  $\alpha$  gives scores on the  $x$  axis which leaves that much  $\alpha/2$  in each tail. Because of the symmetry of the normal distribution, these scores are plus and minus each other. This is called the  $Z$ -value.

### **i** Definition of the $Z$ value using the normal distribution

Using the normal distribution, a  **$Z$ -value** (sometimes called  **$Z$ -score** or **standard score**) is a known test statistic. It shows how many standard deviations above or below the mean an observed data point is.

For the purposes of constructing confidence intervals, the  $Z$ -value is written as  $Z_{\alpha/2}$ . To work out  $Z_{\alpha/2}$ , you need to specify the  $\alpha$  (and/ or CL) value. You can then use the calculator below to find out  $Z_{\alpha/2}$ .

### **i** Example 1

Use the  $Z$ -value calculator below to find the  $Z$  values for  $\alpha = 0.1$ ,  $\alpha = 0.05$ , and  $\alpha = 0.01$ . These values have been chosen as these are the most commonly used alpha values in statistics, with  $\alpha = 0.05$  in particular giving 95% confidence intervals.

Using the normal distribution, alpha value (and/or the confidence level), and the corresponding  $Z$  value, you can then construct a confidence interval.

## How do you construct a confidence interval?

**i** Definition of confidence interval using the normal distribution

The **sample margin of error** is defined to be the following:

$$ME = Z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$$

where

- $(Z_{\alpha/2})$  is the  $Z$ -value corresponding to the chosen confidence level  $(1 - \alpha)$ ,
- $s$  is the sample standard deviation
- $n$  is the sample size.

The fraction  $s/\sqrt{n}$  is often known as the **standard error** of the sample.

A **CL% confidence interval (CI)** is defined to be the interval

$$[\bar{x} - ME, \bar{x} + ME]$$

where  $\bar{x}$  is the sample mean and  $ME$  is the sample margin of error.

Here,  $\bar{x} - ME$  is called your **lower bound** and  $\bar{x} + ME$  is called your **upper bound**.

The

Here's a step-by-step guide to working out a confidence interval.

**💡** General steps for constructing a confidence interval using the normal distribution

**Step 1:** Write down everything you need to work out a confidence interval, which is

- your sample size  $n$
- a sample mean  $\bar{x}$  and sample standard deviation  $s$
- your alpha value ( $\alpha$ ) (or the confidence level (CL)) together with its corresponding  $Z$  value

**Step 2:** Use your  $\alpha$  (or CL) and the  $Z$ -value calculator to find  $Z_{\alpha/2}$ .

**Step 3:** Find the margin of error

$$ME = Z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$$

and then construct the confidence interval for this sample mean:

$$[\bar{x} - ME, \bar{x} + ME]$$

**Step 4:** Check your work! The average of your confidence interval should equal your sample mean.

### **i** Example 2

Cantor's Confectionery have purchased a new computer to help monitor the quality of their products. The computer allows them to input the mean and standard deviation of a sample of their best-selling products, and then put them into the following program.

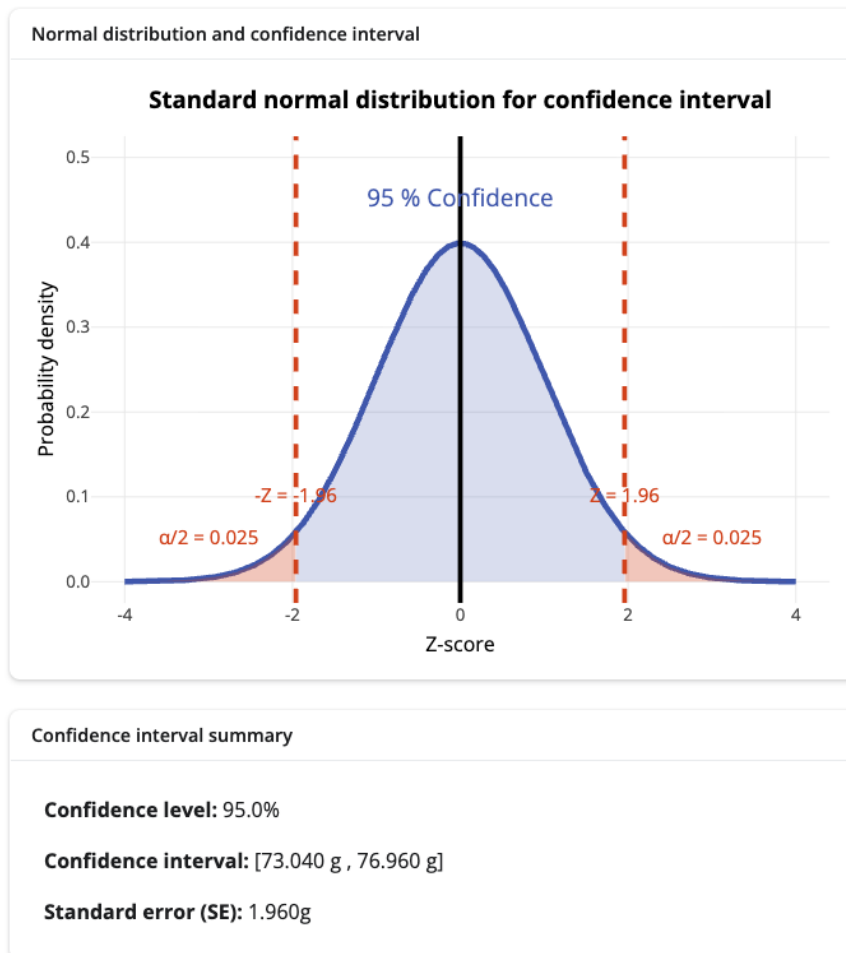


Figure 1: A page of output from the confidence interval program at Cantor's Confectionery.

**i Example 2 (continued)**

What can you tell from this?

- First of all, you can see that this is a normal distribution for the mean weight of one of Cantor's Confectionery's best selling products.
- This takes an input of a sample mean and a sample standard deviation.
- When different values for  $\alpha$  are selected, the tails of the normal distribution change. This means that when constructing a confidence interval, the  $Z$ -values will be different depending on your confidence level.
- By changing any of the parameters, the values of the confidence interval changes. This is because the value of the confidence intervals depend on all of these parameters.

### **i** Example 3

Cantor's Confectionery uses the normal distribution from Example 1 and the computer from Example 2 to construct a 95% confidence interval for the mean weight of their bags of sweets. They take a sample of 100 bags which has an average weight of 75 grams. The calculated standard deviation is 10 grams.

They ask you, an eminent and capable statistician, to check the results of their computation.

**Step 1:** What do you need?

- The sample size is 100 bags of sweets, so  $n = 100$ .
- The average weight of the bags of sweets is 75 grams so  $\bar{x} = 75$  g. The sample standard deviation of the sample is 10 grams so  $s = 10$  g.
- The confidence level is 95%; so  $\alpha = 0.05$  and  $\alpha/2 = 0.025$ .

**Step 2:** Use the Z value calculator to identify  $Z_{\frac{0.05}{2}} = Z_{0.025} = 1.960$ .

**Step 3:** Construct the confidence interval. First, you'll need to work out the margin of error; which since the sample standard deviation is in grams, and the Z-value and square root of sample size are unitless, should be expressed in grams. Here,

$$ME = Z_{0.025} \cdot \frac{s}{\sqrt{n}} = 1.960 \cdot \frac{10 \text{ g}}{\sqrt{100}} = 1.960 \cdot \frac{10 \text{ g}}{10} = 1.960 \text{ g}.$$

Then a 95% CI can be calculated:

$$95\% \text{ CI} = [(75 - 1.960) \text{ g}, (75 + 1.960) \text{ g}] = [73.04 \text{ g}, 76.96 \text{ g}]$$

**Step 4:** Check your work. You know that the sample mean should be the exact centre of your confidence interval. Here

$$73.04 \text{ g} + 76.96 \text{ g} = 150 \text{ g}$$

and since

$$\frac{150 \text{ g}}{2} = 75 \text{ g}$$

you can see that the centre of the confidence interval is the sample mean  $\bar{x} = 75$  g.

#### **i** Example 4

Using the sample mean and sample standard deviation from Example 3, you can ask the computer from Cantor's Confectionery in Example 2 to work out a 90% confidence interval and 99% confidence interval. It outputs the following results:

$$\text{A 90\% CI} = [73.36 \text{ g}, 76.64 \text{ g}]$$

$$\text{A 99\% CI} = [72.424 \text{ g}, 77.576 \text{ g}]$$

What does this suggest about the confidence levels and the corresponding confidence intervals?

For all three examples, the sample mean falls in the centre of the confidence interval. But, as the confidence level **increases**, so does the **width** of the confidence interval. This makes sense; if you want to be more confident as to where the population mean  $\mu$  is, then you need to give a wider set of values.

This is also linked to the tails of the normal distribution. If you were to think about the tails of these normal graphs, as the CL increases the amount of area underneath each extreme decreases. This means you have more values in the middle of the graph, which is why you have more values in the confidence interval. For more on this, see [Guide: More on confidence intervals].

Here's another example which shows that you do not need to follow the precise steps.

### **i** Example 5

There is a new shop in town! Lovelace's Lollies are claiming their products are better than Cantor's Confectionery. They have employed you to construct some 95% confidence intervals for them.

They take a sample of 77 of their best selling boxes of lollies. The average weight of these boxes is 84 grams with a standard deviation of 9.5 grams. So

- The sample size is  $n = 77$ .
- $\bar{x} = 84$  g is the sample mean, and  $s = 9.5$  g is the sample standard deviation.
- It's a 95% CI, and so  $Z_{\frac{0.05}{2}} = Z_{0.025} = 1.960$ .

So to construct the confidence interval you need to first compute the standard error:

$$ME = Z_{\alpha/2} \cdot \frac{s}{\sqrt{n}} = 1.960 \cdot \frac{9.5 \text{ g}}{\sqrt{77}} = 2.122 \text{ g} \quad \text{to 3dp.}$$

Which means that the 95% CI for this sample is

$$95\% \text{ CI} = [(84 - 2.122) \text{ g}, (84 + 2.122) \text{ g}] = [81.878 \text{ g}, 86.122 \text{ g}]$$

This completes your work for Lovelace's Lollies.

## Applications of confidence intervals

It's really important to scrutinize any of your answers and check your understanding of the statistics you are presenting.

So in Examples 3 and 5, Cantor's Confectionery and Lovelace's Lollies have both constructed a 95% confidence interval from a sample of their best selling products. What does this tell you?

- If Cantor's Confectionery were to repeat the study several times, they would expect the average weight of a bag of sweets to **lie between** 73.04 grams and 76.96 grams, **with 95% confidence**.
- If Lovelace's Lollies were to repeat the study several times, they would expect the average weight of a box of lollies to **lie between** 81.88 and 86.12 grams, **with 95% confidence**.

The next example shows what these **do not** tell you.

**i Example 6**

The imminent rivalry between Cantor's Confectionery and Lovelace's Lollies has reached STARMAS News with the following headline:

"Here to compete with Cantor's Confectionery: Lovelace's Lollies **guarantee** that 95% of all boxes of lollies will weigh 86.12 grams".

There are issues with this statement.

From the definition of a confidence level, you know that a confidence level suggests that if you were to repeat the study many times, you would expect the true estimate to fall within CL% of the results.

This is **not the same** as saying CL% of the products weigh a certain amount.

Instead, the STARMAS News headline should read

"Here to compete with Cantor's Confectionery: A study on a sample of Lovelace's Lollies suggest that if more boxes were to be sampled, they expect 95% of boxes would weigh between 81.88 and 86.12 grams".

Finally, you can use a published confidence interval and some rearranging of equations to find out both the sample mean and sample standard deviation from a confidence interval.

**i Example 7**

Cantor's Confectionery release a new "family-sized" bag of sweets to compete with Lovelace's Lollies. They publish a news report on STARMAS<sup>T</sup> News to counter some of the claims made by Lovelace's Lollies.

"We've made a family-sized bag of sweets. Out of a sample of 81 bags, the 99% confidence interval was  $[100.00g, 104.00g]$  to two decimal places!

(Probably not the best PR people at Cantor's Confectionery.)

Anyway, Lovelace's Lollies have got hold of this news report and want you to investigate the sample mean and sample standard deviation of their competitor's study. This is to compare their manufacturing processes with their sample standard deviation from Example 5.

You know that the sample mean is always at the centre of a confidence interval, and so here

$$\bar{x} = \frac{100.00 + 104.00}{2} \text{ g} = 102 \text{ g.}$$

To work out the sample standard deviation  $s$ , you will need to take the equation for the standard error and rearrange for  $s$ .

- Here, the margin of error is the upper bound of the CI minus the sample mean, which is  $ME = (104 - 102) \text{ g} = 2 \text{ g}$ .
- The confidence level is 99%. This means that  $\alpha = 0.01$  and so, from the  $Z$ -value calculator:

$$Z_{\alpha/2} = Z_{0.005} = 2.576.$$

- The sample size is  $n = 81$ , which means that  $\sqrt{n} = \sqrt{81} = 9$ .

So you will have to rearrange the equation

$$2 = 2.576 \cdot \frac{s}{\sqrt{81}}$$

and doing so gives

$$s = \frac{2 \cdot 9}{2.576} = 6.99 \text{ g} \quad \text{to 2dp.}$$

This is the lowest sample standard deviation yet; it seems that Cantor's Confectionery have better manufacturing facilities!

## Quick check problems

1. Cantor's Confectionery has constructed a 99% CI for the population mean length of their pulled taffy, which is [29.1 mm, 32.9 mm]. Below are three statements about this working; two are false, and one is true. Select the true statement:
  - There is a 99% probability that the population mean is in this interval.
  - 99% of the sample data is within this interval.
  - It's likely (to a 99% confidence level) that the population mean is in this interval.
2. What would your CL be for  $\alpha = 0.05$ ?
  - (a) 85%
  - (b) 90%
  - (c) 95%
3. What is the  $Z$ -value for  $\alpha = 0.1$ , to 3 decimal places?
4. If a random variable  $X$  is normally distributed, which parameters are used?
5. What are the extremes of the normal distribution called?

## Further reading

For more questions on the subject, please go to [Questions: Introduction to confidence intervals](#).

If you would like to use a calculator to check your answers, please go to [Calculator: Confidence intervals with normal distribution](#).

## Version history

v1.0: initial version created 12/25 by Millie Harris as part of a University of St Andrews VIP project.

This work is licensed under [CC BY-NC-SA 4.0](#).