

Confidence intervals with Rcmdr

Blanca Lizarbe, blanca.lizarbe@uam.es

Contents

1	Definition of confidence intervals	1
2	Example	1

The following text and explanation is partly inspired from the book “*Discovering statistics using R*” by Z. Field, A. Field and J. Miles”.

1 Definition of confidence intervals

We’re interested in using the sample mean as an estimate of the value of the mean of the population. We have seen that different samples will give rise to different values of the mean, and we can use the standard error to get some idea of the extent to which sample means differ. An approach to assessing the accuracy of the sample mean as an estimate of the mean in the population is to calculate **boundaries that we think will likely capture the true mean**. Such boundaries are called **confidence intervals**. The basic idea behind confidence intervals is to construct a range of values that we think that, most of the time, capture the population mean.

2 Example

Example using Rcmdr: We will now create an example using the *TeachingDemo* plugin. To facilitate the learning process, we will start with an example in which we know the *true mean* of an event, which is normally not known. Imagine we know for sure that, on average, men living in the Community of Madrid (CM) that have had a heart attack in their lives, had their first heart problems at the age of 75, with a SD of 8 years. So the true mean is in this example, 75. We now imagine that we take surveys on men living in Madrid, using samples of 25 individuals. For example, we take 25 individuals living in one specific district of the city. We ask them when they had their first heart problem, and the mean value, on this specific sample, is 77. It is not 75, but it is close. We can repeat this on, for example, 100 districts/villages in the CM. Go to **Rcmdr -> Demos -> Confidence interval for the mean** and generate a dataset with the mentioned characteristics. If you set the “confidence level” to 0.95, you should obtain something *similar* to the figure below.

For each sample surveyed, we have a mean value, and a corresponding boundary. Observe that there is a black vertical line on the “true mean” (75), and that most of the times the boundaries of the CI around each sample mean cover the true mean, but there are a few cases in which this is not true. In fact, **the boundaries surrounding each sample mean are constructed in a way that 95% of the time they will cover the true mean**. So if you repeatedly run the code, you will see that, on average, about 5 of the computed CIs miss the true mean (in the plot, those that are too low have a purple color and those that are too high have a cyan color).

If we change the “confidence level” to 0.99, then this means that the boundaries are constructed in a way that they will cover the true mean 99% of the time. Go to **Rcmdr -> Demos -> Confidence interval for the mean** and generate different conditions to understand the different effects of changing the parameters.

Confidence intervals based on z distribution

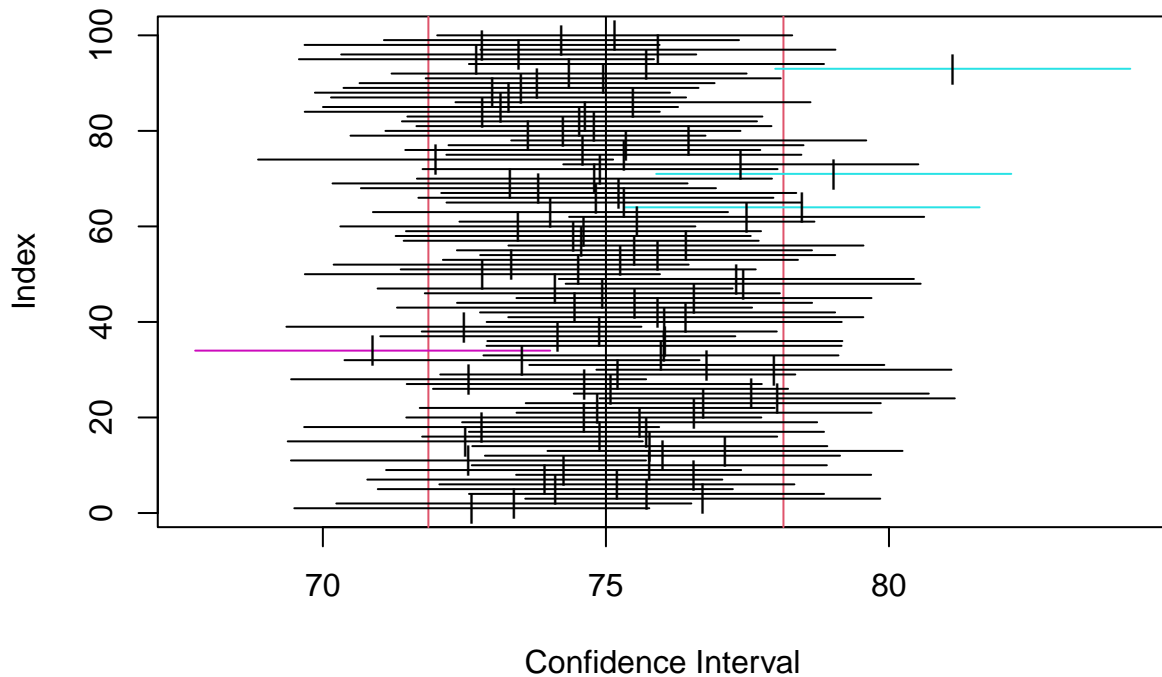


Figure 1: Mean and CI of 100 samples of $n=25$, with a true mean of 75, st dev of 8, and confidence level of 0.95

Notice that, if we leave the settings fixed (e.g., true mean is 75) the true mean never moves. Each time we run the code, we generate 100 sample means with their CIs. Each time, the 100 sample means and their CIs are different. Over those 100 repetitions, for 95% CIs, about 95% of them capture the true mean, and 5% don't. Thus, it is the case that, when we construct 95% CIs, about 95% of those CIs tend to capture the true mean (but it is NOT the case that for any single one of those CIs the true mean has a probability of 95% of falling inside the interval).

To phrase it differently: for every sample, what we compute are the limits of the confidence interval; the true mean is ... wherever it is. Again, it is incorrect to think of the CI as "the probability that the true mean is somewhere"; instead, the CI is constructed in such a way that the CI tends to capture the true mean, because the limits of the CI we computed bracket it. The limits of the CI are random variables; the true mean is whatever the true mean is (some fixed value, unknown to us).