

An Elaboration of the Motivation and Theory Underlying the Bootstrap- t Interval

Imagine we have a sample \mathcal{S} of size n that we assume was drawn from the normal distribution $N(\mu, \sigma^2)$, and the attribute that we care about is the population average (μ). Then we know that $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ and also that

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

and also that

$$Z \equiv \frac{\bar{X} - \mu}{\tilde{\sigma}/\sqrt{n}} \sim t_{(n-1)}.$$

Thus we can write

$$Pr(c_{lower} \leq Z \leq c_{upper}) = 1 - p$$

where c_{lower} is the $p/2$ quantile of the $t_{(n-1)}$ distribution and c_{upper} is the $1 - p/2$ quantile of the $t_{(n-1)}$ distribution.

The probability statement above can be arranged to yield a *random interval* for μ as follows:

$$Pr(c_{lower} \leq Z \leq c_{upper}) = 1 - p$$

$$Pr\left(c_{lower} \leq \frac{\bar{X} - \mu}{\tilde{\sigma}/\sqrt{n}} \leq c_{upper}\right) = 1 - p$$

$$Pr\left(\bar{X} - c_{upper} \times \frac{\tilde{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{X} - c_{lower} \times \frac{\tilde{\sigma}}{\sqrt{n}}\right) = 1 - p$$

and so a $(1 - p) \times 100\%$ random interval for μ is

$$\left[\bar{X} - c_{upper} \times \frac{\tilde{\sigma}}{\sqrt{n}}, \bar{X} - c_{lower} \times \frac{\tilde{\sigma}}{\sqrt{n}}\right].$$

The corresponding $(1 - p) \times 100\%$ *observed interval* is found by substituting *random* quantities above by their corresponding *observed* sample estimates:

$$\left[\bar{x} - c_{upper} \times \frac{\hat{\sigma}}{\sqrt{n}}, \bar{x} - c_{lower} \times \frac{\hat{\sigma}}{\sqrt{n}}\right].$$

Now, the development above assumed that the attribute we care about is the population average, and with that assumption came two convenient consequences:

- We know how to calculate $SD[\bar{X}]$
- The distribution of Z is known

So, how would we proceed if we wanted to calculate a confidence interval for an attribute $a(\mathcal{P})$ other than the average?

Well, analogously the estimator $\tilde{a}(\mathcal{S})$ has some distribution, just like \bar{X} did above, but now we don't know what it is. Similarly, the standardized ratio

$$Z = \frac{\tilde{a}(\mathcal{S}) - a(\mathcal{P})}{\widetilde{SD}[\tilde{a}(\mathcal{S})]}$$

has some distribution – no longer the $t_{(n-1)}$ distribution – but some distribution nonetheless. From this distribution (whatever it is) we can in theory determine c_{lower} and c_{upper} which are respectively its $p/2$ and $1 - p/2$ quantiles and write:

$$Pr(c_{lower} \leq Z \leq c_{upper}) = 1 - p$$

As before this statement can be pivoted, (socially) isolating for $a(\mathcal{P})$ in the middle:

$$Pr(c_{lower} \leq Z \leq c_{upper}) = 1 - p$$

$$Pr\left(c_{lower} \leq \frac{\tilde{a}(\mathcal{S}) - a(\mathcal{P})}{\widetilde{SD}[\tilde{a}(\mathcal{S})]} \leq c_{upper}\right) = 1 - p$$

$$Pr\left(\tilde{a}(\mathcal{S}) - c_{upper} \times \widetilde{SD}[\tilde{a}(\mathcal{S})] \leq a(\mathcal{P}) \leq \tilde{a}(\mathcal{S}) - c_{lower} \times \widetilde{SD}[\tilde{a}(\mathcal{S})]\right) = 1 - p$$

This yields the $(1 - p) \times 100\%$ *random interval*

$$\left[\tilde{a}(\mathcal{S}) - c_{upper} \times \widetilde{SD}[\tilde{a}(\mathcal{S})], \tilde{a}(\mathcal{S}) - c_{lower} \times \widetilde{SD}[\tilde{a}(\mathcal{S})]\right].$$

As above, we want to convert this into an *observed interval*. To do so, we could replace $\tilde{a}(\mathcal{S})$ with the sample estimate $a(\mathcal{S})$. However, we also need

- (1) to be able to estimate $SD[\tilde{a}(\mathcal{S})]$, and
- (2) to know the distribution of Z in order to determine c_{lower} and c_{upper}

To resolve each of these issues, we'll employ the bootstrap. In particular, to address (1), we will approximate the distribution of $\tilde{a}(\mathcal{S})$ using the bootstrap distribution determined by $a(\mathcal{S}_1^*), a(\mathcal{S}_2^*), \dots, a(\mathcal{S}_B^*)$, and thus estimate $SD[\tilde{a}(\mathcal{S})]$ by the bootstrap standard deviation

$$\widehat{SD}_\star[\tilde{a}(\mathcal{S})] = \sqrt{\frac{\sum_{b=1}^B (a(\mathcal{S}_b^*) - \bar{a}^*)^2}{B - 1}}$$

where $\bar{a}^* = \frac{1}{B} \sum_{b=1}^B a(\mathcal{S}_b^*)$ is the average attribute value over all of the bootstrap samples.

To address (2), we will approximate the distribution of Z by the bootstrap distribution of

$$Z^* = \frac{\tilde{a}(\mathcal{S}^*) - a(\mathcal{P}^*)}{\widetilde{SD}[\tilde{a}(\mathcal{S}^*)]} = \frac{\tilde{a}(\mathcal{S}^*) - a(\mathcal{S})}{\widetilde{SD}[\tilde{a}(\mathcal{S}^*)]}.$$

Thus c_{lower} and c_{upper} will be *estimated* as the $p/2$ and $1 - p/2$ quantiles of $z_1^*, z_2^*, \dots, z_B^*$ where

$$z_b^* = \frac{a(\mathcal{S}_b^*) - a(\mathcal{S})}{\widetilde{SD}[\tilde{a}(\mathcal{S}_b^*)]}$$

(Note that a double bootstrap will likely be necessary to determine $\widehat{SD}[a(\mathcal{S}_b^*)]$).

Thus the $(1 - p) \times 100\%$ *observed interval* is given by

$$\left[a(\mathcal{S}) - c_{upper} \times \widehat{SD}_*[\tilde{a}(\mathcal{S})], a(\mathcal{S}) - c_{lower} \times \widehat{SD}_*[\tilde{a}(\mathcal{S})] \right].$$

This is what is referred to as the bootstrap- t confidence interval for $a(\mathcal{P})$.