

Bootstrap Methods

Typically we use (asymptotic) theory to derive the sampling distribution of a statistic. From the sampling distribution, we can obtain the variance, construct confidence intervals, perform hypothesis tests, and more.

Challenge:

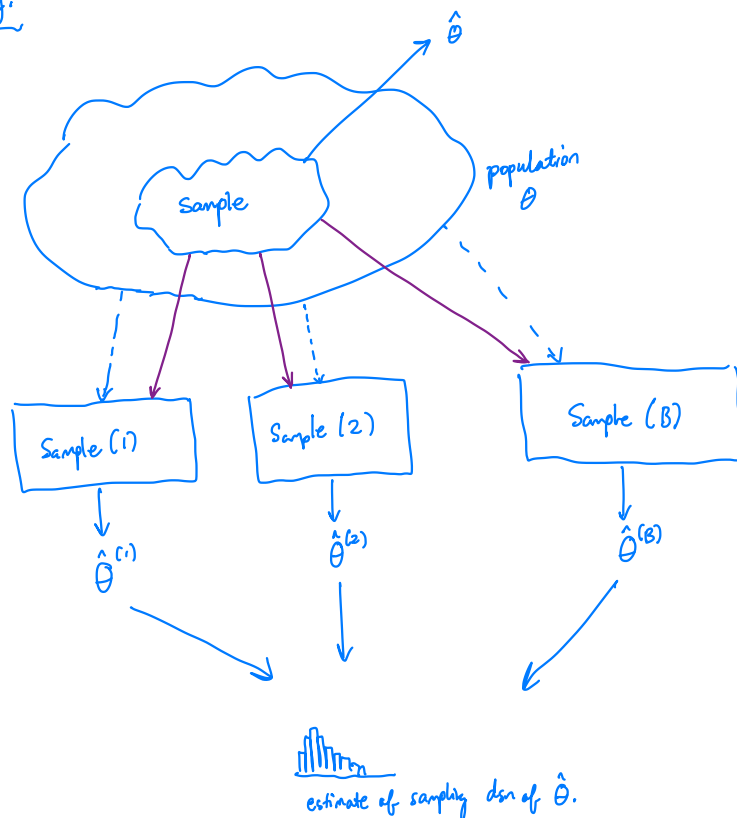
what if the sampling distribution is impossible to obtain or asymptotic theory doesn't hold?

Basic idea of bootstrapping:

Use the data to approximate the sampling distribution of the statistic.

How? Estimate the sampling distribution by creating a large # of data sets that we might have seen, and compute the statistic on each of the data sets.

E.g.



Goals of bootstrap:

estimate bias, se, CI's when

- 1) there is doubt about if distributional assumptions are met
- 2) there is doubt about whether asymptotic results hold.
- 3) The theory to derive the disn of a statistic is too hard.

In reality, we only have a sample and need to make $\text{sample}^{(1)}, \dots, \text{sample}^{(B)}$.

"Bootstrap World" where the data analyst knows everything.

Idea: treat the sample Y_1, \dots, Y_n as the population

\approx the population consists of an infinite number of Y_1 values + infinite # of Y_2 values, each occurring $\frac{1}{n}$ proportion of the time.

\Rightarrow we can resample

e.g. we are interested in the variance of the estimator.

= In "bootstrap world" we can calculate the exact variance b/c we have access to "population"

- In practice, estimate variance by repeatedly sampling from pseudo-population.

\rightarrow we have a representative sample of the population.

