

2025 BIOSTAT 523 STATISTICAL INFERENCE FOR BIOMETRY II

SLIDE SET 14: THE BOOTSTRAP

Jon Wakefield

Departments of Statistics and Biostatistics
University of Washington
jonno@uw.edu

MOTIVATION

THE BOOTSTRAP

VARIANCE ESTIMATION VIA THE BOOTSTRAP

CONSTRUCTING A CONFIDENCE INTERVAL VIA THE
BOOTSTRAP

FAILURE OF THE BOOTSTRAP

DISCUSSION

MOTIVATION

MOTIVATION AND OVERVIEW

We have seen that the jackknife sampling method can produce estimates of bias and variance, by resampling from the observed data.

But, the jackknife does not work for all statistics, for example, it does not work for the median.

Conformal prediction is a method for a particular type of prediction problem.

We now study the bootstrap, which is more computationally expensive than the jackknife to use, but offers some advantages.

The bootstrap is extremely popular, but is sometimes used, when it shouldn't (because there are rules!).

THE BOOTSTRAP

The bootstrap is a popular **resampling technique** that is useful in situations in which:

- ▶ We wish to **relax the assumptions** of a parametric modeling approach.
- ▶ The **asymptotic sampling distribution** of the estimator is difficult to derive.

References: Efron and Tibshirani (1993); Davison and Hinkley (1997); Chernick (2011); Efron and Hastie (2016).

We may be interested in bootstrap methods for:

- ▶ Estimating the variance of an estimator.
- ▶ Constructing CIs for parameters.

For simplicity, suppose we are in a population setting where we have independent and identically distributed (iid) data Y_1, \dots, Y_n , and denote the (unknown) cumulative distribution function of Y by F .

PARAMETERS OF INTEREST

Let $\theta = T(F)$ be a parameter of interest.

Examples:

1. Suppose $Y_1, \dots, Y_n \sim F$ where $F \in (F_\theta, \theta \in \Theta)$ and $\hat{\theta}_n$ be the **MLE** of θ . We would like to estimate $\text{var}(\hat{\theta}_n)$ and a $1 - \alpha$ confidence interval (CI) for θ .
2. Suppose $Y_1, \dots, Y_n \sim F$, and $\theta = T(F)$ is the mean of F , i.e., $\theta = E[Y] = \int y \, dF(y)$ and $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n Y_i$ is the **sample mean**. Again, we would like to estimate $\text{var}(\hat{\theta}_n)$ and a $1 - \alpha$ CI for θ .
3. Suppose $Y_1, \dots, Y_n \sim F$, and $\theta = T(F)$ is the median of F , i.e., $\Pr(Y_i \leq \theta) = \Pr(Y_i > \theta) = 1/2$ and $\hat{\theta}_n$ is the **sample median**. Again, we would like to estimate $\text{var}(\hat{\theta}_n)$ and a $1 - \alpha$ CI for θ .

PARAMETERS OF INTEREST

In 1., θ is a parameter in a **parametric model** – in this case applying the delta method may be cumbersome, and we might also want to find a CI that doesn't depend on the model being correct (and evaluating the sandwich would be tricky).

In 2. and 3., we are in a **nonparametric situation**, as we have no specific model.

INFERENCE FOR THE SAMPLE MEDIAN

For the sample mean, a $1 - \alpha$ confidence interval for the population mean μ is

$$\bar{Y}_n \pm z_{1-\alpha/2} \times \frac{s_n}{\sqrt{n}},$$

where \bar{Y}_n and s_n are the sample mean and sample standard deviation.

Can we do the same thing for the sample median, which we denote by $\hat{\theta}_n$?

Note that if θ is the population median (and F^{-1} has a continuous derivative in a neighborhood of $1/2$):

$$\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow_d N\left(0, \frac{1}{4f^2(\theta)}\right),$$

where $f(\cdot)$ is the density of Y . This is not straightforward to use in practice, because of the dependence on the unknown $f(\cdot)$.

Monte Carlo Methods

If F is known, it is straightforward to mimic frequentist inference; for simplicity, suppose we have a univariate parameter θ .

For $b = 1, \dots, B$ samples:

- ▶ Generate a random sample $y_1^{*(b)}, \dots, y_n^{*(b)} \sim_{iid} F$.
- ▶ Compute $\hat{\theta}_n^{*(b)}$ using $y_1^{*(b)}, \dots, y_n^{*(b)}$.

Use $\hat{\theta}_n^{*(b)}, b = 1, \dots, B$, to estimate the sampling distribution of $\hat{\theta}_n$.

If $B \rightarrow \infty$, we approach the theoretical sampling distribution of $\hat{\theta}_n$.

Of course, in practice, F is unknown and the idea behind the bootstrap is to resample datasets from an estimate of F .

BOOTSTRAP METHODS: UNKNOWN F

Two obvious choices for estimating F :

- ▶ Use the **empirical distribution function** of the data, which we denote \hat{F}_n – this is by far the most common approach used.
- ▶ If one has some faith in the **assumed model** then we may use this model, call this $F_{\hat{\theta}_n}$, where the notation emphasizes that the distribution function is fully specified by the parameter θ , which we estimate by $\hat{\theta}_n$.

Then we may:

- ▶ Sample data with replacement from \hat{F}_n to give the **nonparametric** bootstrap.
- ▶ Sample data from $F_{\hat{\theta}_n}$ to give the **parametric** bootstrap.

THE EMPIRICAL DISTRIBUTION FUNCTION

The empirical distribution function (EDF) is an estimator of the cumulative distribution function (CDF).

The CDF at a fixed value y_0 is,

$$F(y_0) = \Pr(Y_i \leq y_0),$$

for $i = 1, \dots, n$, so that $F(y_0)$ is the probability of the event $\{Y_i \leq y_0\}$.

The natural estimator of this probability is the empirical proportion:

$$\hat{F}_n(y_0) = \frac{\text{Number of } Y_i \leq y_0}{\text{Total number of observations}} = \frac{\sum_{i=1}^n I(Y_i \leq y_0)}{n}.$$

We can do this for all y to give the EDF (this is also the nonparametric MLE of F).

In Figure 1 we see that as n increases we approach the true CDF.

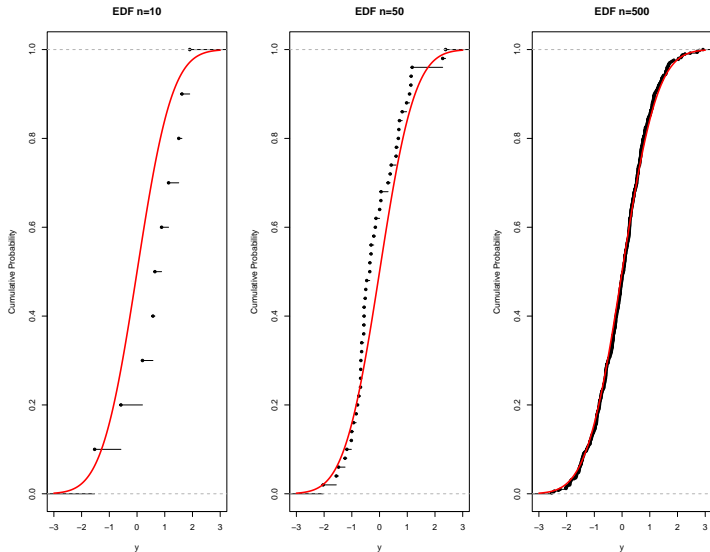


FIGURE 1: Empirical distribution function estimate along with true CDF (in red). There are jumps of $1/n$ at every data point, and between data points, the EDF $\hat{F}_n(y)$ is flat.

THE EMPIRICAL DF

The EDF at y is an average of $Z_i = I(Y_i \leq y)$ so that

$$Z_i = \begin{cases} 1 & \text{if } Y_i \leq y \\ 0 & \text{if } Y_i > y \end{cases}$$

is such that $Z_i \sim \text{Bernoulli}(F(y))$ so that

$$\begin{aligned} E[I(Y_i \leq y)] &= E[Z_i] = F(y) \\ \text{var}(I(Y_i \leq y)) &= \text{var}(Z_i) = F(y)(1 - F(y)) \end{aligned}$$

for a given y .

Recall that $\hat{F}_n(y) = \sum_{i=1}^n I(Y_i \leq y) = \sum_{i=1}^n Z_i$, so

$$\begin{aligned} E[\hat{F}_n(y)] &= E[I(Y_1 \leq y)] = F(y) \\ \text{var}(\hat{F}_n(y)) &= \frac{\sum_{i=1}^n \text{var}(Z_i)}{n^2} = \frac{F(y)(1 - F(y))}{n} \end{aligned}$$

THE EMPIRICAL DF AS AN ESTIMATOR OF $F(y)$

Properties as an estimator,

$$\text{Bias}(\hat{F}_n(y)) = E[\hat{F}_n(y)] - F(y) = 0.$$

The variance converges to 0 as $n \rightarrow \infty$ so that

$$\hat{F}_n(y) \rightarrow_p F(y),$$

so that $\hat{F}_n(y)$ is a **consistent estimator** of $F(y)$.

In addition,

$$\sqrt{n}(\hat{F}_n(y) - F(y)) \rightarrow_d N(0, F(y)(1 - F(y))).$$

These are **pointwise** properties of the estimator, i.e., at y .

It can also be shown that the complete empirical DF converges to F , i.e.,

$$\sqrt{n}(\hat{F}_n - F) \rightarrow_d B,$$

where B is a **Brownian bridge**.

THE NON-PARAMETRIC BOOTSTRAP

Again, consider a univariate parameter, $\theta = T(F)$.

For $b = 1, \dots, B$ **samples** the **non-parametric bootstrap** samples:

- ▶ Generate a random sample (known as a **bootstrap sample**) $y_1^{*(b)}, \dots, y_n^{*(b)} \sim_{iid} \hat{F}_n$ – this is equivalent to drawing n observations, **with replacement** from the original data $\{Y_1, \dots, Y_n\}$.
- ▶ Compute $\hat{\theta}_n^{*(b)}$ using $y_1^{*(b)}, \dots, y_n^{*(b)}$.

Use $\hat{\theta}_n^{*(b)}, b = 1, \dots, B$, to estimate the sampling distribution of $\hat{\theta}_n$.

Note that we can't enumerate all possible bootstrap samples, as there are $\binom{2n-1}{n}$ of them!

There are two approximations/sources of error in the bootstrap:

1. **Statistical:** $\hat{F}_n \neq F$.
2. **Simulation:** $B \neq \infty$, but we can take B large.

For 1., if n is small, the approximation will be poorer when we use \hat{F}_n .

For 2., for some targets such as the bias and variance, we can get away with smaller B (e.g., $B \geq 200$), but for others such as confidence intervals we need larger B (e.g., $B \geq 1000$).

Lehmann (1999, p.426), refers to $T(\hat{F}_n^*)$ as an **approximator** of $T(\hat{F}_n)$ rather than an **estimator**.

BOOTSTRAP METHODS

Theoretical: The population is to the sample

Bootstrap: The sample is to the bootstrap sample

As an example, the theoretical bias of as estimator is

$$E[\hat{\theta}(\mathbf{Y})|F] - \theta.$$

The bootstrap attempts to estimate this by

$$\underbrace{E[\hat{\theta}(\mathbf{Y}^*)|\hat{F}_n]}_{\text{Average over Samples}} - \underbrace{\theta}_{\text{Population}}.$$

And in practice this is estimated by

$$\underbrace{\frac{1}{B} \sum_{b=1}^B \hat{\theta}^{*(b)}}_{\text{Average over Bootstrap Samples}} - \underbrace{\hat{\theta}(\mathbf{y})}_{\text{Sample Estimate}}.$$

HEURISTIC OF WHY THE BOOTSTRAP WORKS FOR THE MEDIAN

- ▶ We have $Y_1, \dots, Y_n \sim F$.
- ▶ The distribution of the median, based on a sample of size n , will be denoted $F_{M_n}(y)$ which we write as,

$$F_{M_n}(y) = \Psi(y; F, n),$$

where Ψ is some complicated function.

- ▶ We know that $\hat{F}_n \approx F$ when n is large so as long as Ψ is **smooth** with respect to F then $\Psi(y; \hat{F}_n, n)$ will be similar to $\Psi(y; F, n)$, i.e.,

$$\begin{aligned} \hat{F}_n \approx F \quad \implies \quad F_{M_n^{*(b)}}(y) &= \Psi(y; \hat{F}_n, n) \\ &\approx \Psi(y; F, n) = F_{M_n}(y), \end{aligned}$$

where $M_n^{*(b)}$ is the function evaluated for the b -th bootstrap sample.

HEURISTIC OF WHY THE BOOTSTRAP WORKS FOR THE MEDIAN

- ▶ Hence, since the **bootstrap samples** are all from \hat{F}_n and we have

$$\Psi(y; F, n) = F_{M_n^{\star(1)}}(y) = F_{M_n^{\star(2)}}(y) = \cdots = F_{M_n^{\star(B)}}(y).$$

- ▶ What this means is:

The CDF of the median based on the bootstrap samples, $F_{M_n^{\star}}(y)$ is approximating the CDF of the true CDF of the median in a sample of size n , $F_{M_n}(y)$.

OVERVIEW OF THE THEORY

There is a great deal of theory on when the bootstrap does and does not work.

A key observation is to note that, while the bootstrap is applicable in many situations, it is not valid in **all** situations, and so care should be taken in when it is applied.

A starting rule is that if we are in a situation where the delta method is valid, then the bootstrap will also work – this needs asymptotic normality and smoothness.

From Section 4 of Chapter 18 at:

<https://sites.stat.washington.edu/jaw/COURSES/580s/581/lectnotes.18.html>

OVERVIEW OF THE THEORY

The theory needs to show that the asymptotic behavior of the distribution of the nonparametric bootstrap “mimics” the behavior of the original estimator in probability or almost surely (a.s.).

If we are estimating $T(F)$ by $T(\hat{F}_n)$ and we know (perhaps from a delta method argument) that

$$\sqrt{n} \underbrace{(T(\hat{F}_n) - T(F))}_{\hat{\theta}_n - \theta} \rightarrow_d N(0, \text{var}_F(T)),$$

then we need to show that the bootstrap estimator, $T(\hat{F}_n^*)$, satisfies:

$$\sqrt{n}(T(\hat{F}_n^*) - T(\hat{F}_n)) \rightarrow_d N(0, \text{var}_F(Y)), \quad \text{in probability or a.s.}$$

OVERVIEW OF THE THEORY

For the sample mean of a distribution F on \mathbb{R} , if $Y \sim F$ and $E[Y^2] < \infty$ then for

$$T(F) = \int y dF(y) = \mu(F)$$

we know that

$$\sqrt{n}(T(\hat{F}_n) - T(F)) = \sqrt{n}(\bar{Y}_n - \mu(F)) \rightarrow_d N(0, \text{var}(Y)).$$

For the bootstrap, the corresponding statement is: If $E[Y^2] < \infty$ then for Y_1, Y_2, \dots ,

$$\sqrt{n}(T(\hat{F}_n^*) - T(\hat{F}_n)) = \sqrt{n}(\bar{Y}_n^* - \bar{Y}_n) \rightarrow_d N(0, \text{var}(Y)).$$

This can be proved using a central limit theorem, see Bickel and Freedman (1981).

VARIANCE ESTIMATION VIA THE BOOTSTRAP

BOOTSTRAP VARIANCE ESTIMATOR

The **bootstrap variance estimator** is

$$\widehat{\text{var}}(\hat{\theta}_n^*) = \frac{1}{B} \sum_{b=1}^B \left(\hat{\theta}_n^{*(b)} - \bar{\theta}^* \right)^2,$$

where $\bar{\theta}^* = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_n^{*(b)}$.

When B is large the sample variance of the bootstrap estimators

$$\widehat{\text{var}}(\hat{\theta}_n^*) \approx \text{var}(\hat{\theta}_n^* | \hat{F}_n), \quad (1)$$

where we have conditioned on \hat{F}_n being fixed.

BOOTSTRAP VARIANCE ESTIMATOR

To argue that the bootstrap variance is a good estimate of the target variance $v(\hat{\theta}_n)$, we need to have

$$\widehat{\text{var}}(\hat{\theta}_n^*) \approx \text{var}(\hat{\theta}_n^* | \hat{F}_n) \approx \text{var}(\hat{\theta}_n),$$

but because in (1) can be controlled with B large what really matters is

$$\text{var}(\hat{\theta}_n^* | \hat{F}_n) \approx \text{var}(\hat{\theta}_n),$$

or, more formally,

$$\frac{\text{var}(\hat{\theta}_n^* | \hat{F}_n)}{\text{var}(\hat{\theta}_n)} \rightarrow_p 1 \quad \text{as } n \rightarrow \infty.$$

The ratio is often used when both quantities converge to 0 as $n \rightarrow \infty$.

Not too tricky to show this for many statistics $\theta = T(F)$.

SANDWICH ESTIMATION AND THE BOOTSTRAP

We heuristically show why we would often expect sandwich and bootstrap variance estimates to be in close correspondence.

For simplicity, we consider a univariate parameter θ , and let $\hat{\theta}_n$ denote the MLE arising from a sample of size n .

In a change of notation we denote the score by

$$\mathbf{S}(\theta) = [S_1(\theta), \dots, S_n(\theta)]^\top,$$

where $S_i(\theta) = d\ell_i/d\theta$ is the contribution to the score from observation Y_i , $i = 1, \dots, n$.

Hence,

$$\mathbf{S}(\theta) = \sum_{i=1}^n S_i(\theta) = \mathbf{S}(\theta)^\top \mathbf{1}$$

where $\mathbf{1}$ is an $n \times 1$ vector of 1's.

SANDWICH ESTIMATION AND THE BOOTSTRAP

The sandwich form of the asymptotic variance of $\hat{\theta}_n$ is

$$\text{var}(\hat{\theta}_n) = \frac{1}{n} \frac{B_1}{A_1^2}$$

where

$$A_1(\theta) = E \left[\frac{\partial S_1}{\partial \theta} \right], \quad B_1(\theta) = E [S_1(\theta)^2]$$

which may be empirically estimated via

$$\begin{aligned} \hat{A}_n &= \left. \frac{1}{n} \frac{dS}{d\theta} \right|_{\hat{\theta}_n} = \frac{1}{n} \sum_{i=1}^n \left. \frac{dS_i}{d\theta} \right|_{\hat{\theta}_n} \\ \hat{B}_n &= \left. \frac{1}{n} \mathbf{S}(\theta)^\top \mathbf{S}(\theta) \right|_{\hat{\theta}_n} = \frac{1}{n} \sum_{i=1}^n \left. S_i(\theta)^2 \right|_{\hat{\theta}_n}. \end{aligned}$$

SANDWICH ESTIMATION AND THE BOOTSTRAP

A convenient representation of a bootstrap sample is $\mathbf{Y}^* = \mathbf{Y} \times \mathbf{D}$ where $\mathbf{D} = \text{diag}(D_1, \dots, D_n)$ is a diagonal matrix consisting of multinomial random variables

$$\begin{bmatrix} D_1 \\ \vdots \\ D_n \end{bmatrix} \sim \text{Multinomial} \left[n, \left(\frac{1}{n}, \dots, \frac{1}{n} \right) \right]$$

with

$$\begin{aligned} E([D_1, \dots, D_n]^T) &= \mathbf{1} \\ \text{var}([D_1, \dots, D_n]^T) &= \mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T \rightarrow \mathbf{I}_n \end{aligned}$$

as $n \rightarrow \infty$.

The MLE of θ in the bootstrap sample is denoted $\hat{\theta}_n^*$ and satisfies $S^*(\hat{\theta}_n^*) = 0$, where $S^*(\theta)$ is the score corresponding to \mathbf{Y}^* . Note that,

$$S^*(\theta) = \sum_{i=1}^n S_i^*(\theta) = \sum_{i=1}^n S_i(\theta) D_i.$$

SANDWICH ESTIMATION AND THE BOOTSTRAP

We consider a Taylor series expansion

$$0 = S^*(\hat{\theta}_n^*) \approx S^*(\hat{\theta}_n) + (\hat{\theta}_n^* - \hat{\theta}_n) \left. \frac{dS^*}{d\theta} \right|_{\hat{\theta}_n}$$

which leads to the one-step approximation

$$\hat{\theta}_n^* = \hat{\theta}_n - \frac{S^*(\hat{\theta})}{\left. \frac{d}{d\theta} S^*(\theta) \right|_{\hat{\theta}_n}}.$$

The bootstrap score evaluated at $\hat{\theta}_n$ is

$$\sum_{i=1}^n S_i^*(\hat{\theta}_n) = \sum_{i=1}^n S_i(\hat{\theta}_n) D_i \neq 0,$$

unless the bootstrap sample coincides with the original sample, i.e.,
unless $\mathbf{D} = \mathbf{I}_n$.

SANDWICH ESTIMATION AND THE BOOTSTRAP

We replace $\mathbf{S}^*(\hat{\theta}) \left[\frac{d}{d\theta} \mathbf{S}^*(\theta) \Big|_{\hat{\theta}_n} \right]^{-1}$ by its limit

$$\mathbb{E} \left[\frac{d}{d\theta} \mathbf{S}^*(\theta) \Big|_{\hat{\theta}_n} \right] = \mathbb{E} \left[\sum_{i=1}^n \frac{d}{d\theta} \mathbf{S}_i(\theta) \mathbf{D}_i \Big|_{\hat{\theta}_n} \right] = \frac{d}{d\theta} \mathbf{S}(\theta) \Big|_{\hat{\theta}_n} \mathbb{E}[\mathbf{D}] = n \times \hat{\mathbf{A}}_n$$

where $\hat{\mathbf{A}}_n = \frac{d}{d\theta} \mathbf{S}(\theta) \Big|_{\hat{\theta}_n}$.

Therefore, the one-step bootstrap estimator is approximated by

$$\hat{\theta}_n^* \approx \hat{\theta}_n - \frac{\mathbf{S}(\hat{\theta}_n) \mathbf{D}}{n \hat{\mathbf{A}}_n}$$

and is approximately unbiased as an estimator since

$$\mathbb{E}[\hat{\theta}_n^* - \hat{\theta}_n] \approx - \frac{\mathbf{S}(\hat{\theta}_n) \mathbb{E}[\mathbf{D}]}{n \hat{\mathbf{A}}_n} = - \frac{\mathbf{S}(\hat{\theta}_n) \mathbf{1}}{n \hat{\mathbf{A}}_n} = 0$$

and, recall, $\hat{\theta}_n$ is being held constant.

SANDWICH ESTIMATION AND THE BOOTSTRAP

The variance is

$$\begin{aligned}\text{var}(\hat{\theta}_n^* - \hat{\theta}_n) &\approx \frac{\mathbf{S}(\hat{\theta}_n) \text{var}([D_1, \dots, D_n]) \mathbf{S}(\hat{\theta}_n)}{(n\hat{A}_n)^2} = \frac{\mathbf{S}(\hat{\theta}_n) \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) \mathbf{S}(\hat{\theta}_n)}{(n\hat{A}_n)^2} \\ &\approx \frac{\mathbf{S}(\hat{\theta}_n) \mathbf{S}(\hat{\theta}_n)}{(n\hat{A}_n)^2} = \frac{n\hat{B}_n}{(n\hat{A}_n)^2} = \frac{\hat{B}_n}{n\hat{A}_n^2},\end{aligned}$$

which is the sandwich estimator.

Hence, $\text{var}(\hat{\theta}_n^* - \hat{\theta}_n)$ approximates $\text{var}(\hat{\theta}_n - \theta)$, which is a fundamental link to the bootstrap.

For a more theoretical treatment, see Arcones and Giné (1992) and Section 10.3 of Kosorok (2008).

CONSTRUCTING A CONFIDENCE INTERVAL VIA THE BOOTSTRAP

For CI construction, many improvements on the above normal-based method have been suggested:

- ▶ Wald-type CIs.
- ▶ The percentile method – pick the appropriate sample quantiles.
- ▶ Various improved and bias corrected versions have been proposed. Figure 2 is from Section 3.1 of Chernick (2011). See also Sections 11.3–11.5 of Efron and Hastie (2016).

Table 3.1 Four Methods of Setting Approximate Confidence Intervals for a Real Valued Parameter θ

Method	Abbreviation	α -Level Endpoint	Correct if
1. Standard	$\theta_S[\alpha]$	$\hat{\theta} + \hat{\sigma} z^{(\alpha)}$	$\hat{\theta} \approx N(\theta, \sigma^2)$ σ is constant
2. Percentile	$\theta_P[\alpha]$	$\hat{G}^{-1}(\alpha)$	There exists a monotone transformation such that $\hat{\phi} = g(\hat{\theta})$, where, $\phi = g(\theta)$, $\hat{\phi} \approx N(\phi, \tau^2)$ and τ is constant.
3. Bias-corrected	$\theta_{BC}[\alpha]$	$\hat{G}^{-1}([\phi[2z_\alpha + z^{(\alpha)}]])$	There exists a monotone transformation such that $\hat{\phi} \approx N(\phi - z_0 \tau, \tau^2)$ and z_0 and τ are constant.
4. BC_a	$\theta_{BC_a}[\alpha]$	$\hat{G}^{-1}\left(\phi\left[z_0 + \frac{[z_0 + z^{(\alpha)}]}{1 - a[z_0 + z^{(\alpha)}]}\right]\right)$	There exists a monotone transformation such that $\hat{\phi} \approx N(\phi - z_0 \tau_\phi, \tau_\phi^2)$, where $\tau_\phi = 1 + a\phi$ and z_0 and a are constant.

Note: Each method is correct under more general assumptions than its predecessor. Methods 2, 3, and 4 are defined in terms of the percentile of G , the bootstrap distribution.

Source: Efron and Tibshirani (1986, Table 6) with permission from The Institute of Mathematical Statistics.

FIGURE 2: From Chernick (2011).

BOOTSTRAP METHODS: VARIANCE AND CI ESTIMATION

Recall

$$\widehat{\text{var}}(\hat{\theta}_n^*) = \frac{1}{B} \sum_{b=1}^B \left(\hat{\theta}_n^{*(b)} - \bar{\theta}^* \right)^2,$$

where $\bar{\theta}^* = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_n^{*(b)}$.

If n is sufficiently large that asymptotic normality of the estimator may be appealed to and a CI estimate may be based upon

$$\hat{\theta}_n \pm z_{1-\alpha/2} \times \sqrt{\widehat{\text{var}}(\hat{\theta}_n^*)}.$$

This is a Wald-type CI, and is not invariant to reparametrization.

BOOTSTRAP CONFIDENCE INTERVAL: THE PERCENTILE METHOD

Let

$$\hat{G}(t) = \frac{1}{B} \sum_{b=1}^B I\left(\sqrt{n}(\hat{\theta}_n^{*(b)} - \hat{\theta}_n) \leq t\right).$$

Then a $1 - \alpha$ bootstrap percentile method confidence interval is

$$C_n = \left[\hat{\theta}_n - \frac{t_{1-\alpha/2}}{\sqrt{n}}, \hat{\theta}_n - \frac{t_{\alpha/2}}{\sqrt{n}} \right],$$

where $t_{\alpha/2} = \hat{G}^{-1}(\alpha/2)$ and $t_{1-\alpha/2} = \hat{G}^{-1}(1 - \alpha/2)$.

Under appropriate regularity conditions,

$$\Pr(\theta \in C_n) = 1 - \alpha - O\left(\frac{1}{\sqrt{n}}\right).$$

Note that this method, unlike the Wald-type interval given earlier, is invariant to the parameterization adopted.

FAILURE OF THE BOOTSTRAP

Let $Y_1, \dots, Y_n \sim_{iid} \text{Uniform}(0, \theta)$, and $\hat{\theta}_n = \min\{Y_1, \dots, Y_n\}$, be the minimum of the sample and corresponds to the MLE.

It can be shown that

$$n(\theta - \hat{\theta}_n) \rightarrow_d \text{Exponential}(\theta),$$

i.e., converges to a exponential distribution with mean θ .

However,

$$\begin{aligned} \Pr(\hat{\theta}_n^* = Y_{(n)} | \hat{F}_n) &= 1 - \Pr(Y_{(n)}^* < Y_{(n)} | \hat{F}_n) \\ &= 1 - \Pr(\text{all } Y_i^* < Y_{(n)} | \hat{F}_n) = 1 - \left(\frac{n-1}{n}\right)^n \\ &= 1 - \left(1 - \frac{1}{n}\right)^n \rightarrow 1 - e^{-1} \approx 0.632. \end{aligned}$$

Thus, the bootstrap will select the maximum in the observed data a big chunk of the time, and is not close to the exponential distribution.

Figure 3 compares the non-parametric and parametric bootstraps.

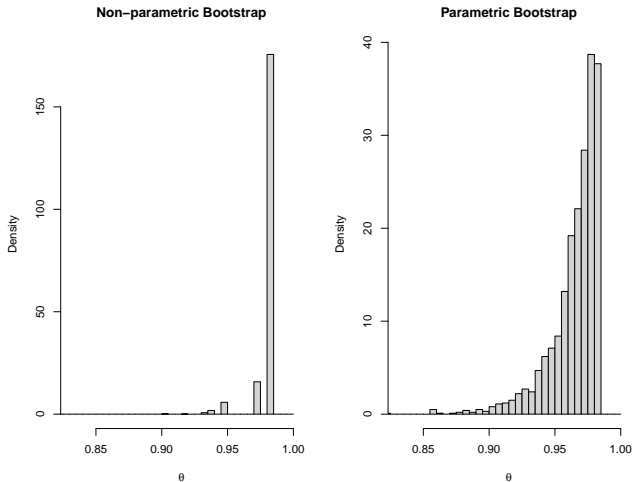


FIGURE 3: Simulated data from $\text{Uniform}(0, \theta)$ with $\theta = 1$ and $n = 50$. Bootstrap samples of size $B = 2000$ were obtained for a non-parametric bootstrap (left) and a parametric bootstrap (right), i.e., from $\text{Uniform}(0, \hat{\theta})$. The problem here is that the empirical distribution is not a good approximation to the true distribution in the tail.

EXAMPLE: LUNG CANCER AND RADON

For the lung cancer and radon example we implement the nonparametric bootstrap resampling, with replacement, $B = 1000$ sets of n case triples $[Y_{bi}^*, E_{bi}^*, x_{bi}^*]$, $b = 1, \dots, B$, $i = 1, \dots, n$.

Figure 4 displays the histogram of estimators arising from the bootstrap samples, along with the asymptotic normal approximations to the sampling distribution of the estimator under the Poisson and quasi-Poisson models.

We see that the distribution under the quasi-likelihood model is much wider than that under the Poisson model.

This is not surprising since we have already seen that the lung cancer data are overdispersed relative to a Poisson distribution.

The bootstrap histogram and quasi-Poisson sampling distribution are very similar, however.

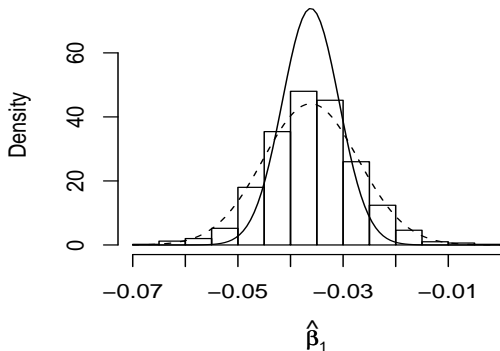


FIGURE 4: Sampling distribution of $\hat{\beta}_1$ arising from the nonparametric bootstrap samples. The solid curve is the asymptotic distribution of the MLE under the Poisson model, and the dashed line is the asymptotic distribution under the quasi-Poisson model.

EXAMPLE: LUNG CANCER AND RADON

Table 1 summarizes inference for β_1 for a number of different methods, and again confirms the similarity of asymptotic inference and the parametric bootstrap under the Poisson model.

The parametric bootstrap cannot be used with a quasi-likelihood model since there is no probability distribution for the data.

Point estimates from the Poisson, quasi-likelihood and sandwich approaches are identical.

EXAMPLE: LUNG CANCER AND RADON

Inferential Method	$\hat{\beta}_1$ ($\times 10^3$)	s.e.($\hat{\beta}_1$) ($\times 10^4$)	95% confidence interval for $e^{10\beta_1}$
Poisson	-0.036	0.0054	0.954, 0.975
Quasi-Likelihood	-0.036	0.0090	0.947, 0.982
Quadratic Variance	-0.030	0.0085	0.955, 0.987
Sandwich Estimation	-0.036	0.0080	0.949, 0.980
Bootstrap Normal	-0.036	0.0087	0.948, 0.981
Bootstrap Percentile	-0.036	0.0087	0.949, 0.981

TABLE 1: Comparison of inferential approaches for the lung cancer example.

FAILURE OF THE BOOTSTRAP

FAILURE OF THE BOOTSTRAP

- ▶ **Dependence:**

- ▶ If there are “natural” structure to the data that lead to correlated outcomes, such as over time, space, over networks or within families, or “induced” structure due to the experimental design, or complex survey structure. In these cases, the bootstrap distribution will not correspond to the true asymptotic distribution, unless the bootstrap sampling respects the data structure.

- ▶ **Lack of Smoothness (cube-root asymptotics):**

- ▶ A number of estimators converge at a rate of $n^{-1/3}$, rather than the more usual $n^{-1/2}$ rate. For example, the least median of squares estimator in a linear regression. For other examples of cube root asymptotics, see Kim and Pollard (1990). The limit distributions are non-normal, and the bootstrap does not work.

► Machine Learning Methods:

- Sparse estimators, such as the lasso (Tibshirani, 1996), are not amenable to being bootstrapped, because zero is a special case for the regression coefficients. Dezeure *et al.* (2015) with respect to the bootstrap and high-dimensional inference say, “...the asymptotic distribution of the Lasso has point mass at zero. This implies, because of noncontinuity of the distribution, that standard bootstrapping and subsampling schemes are delicate to apply and uniform convergence to the limit seems hard to achieve”.
- Chatterjee and Lahiri (2013) discuss the **residual bootstrap** in which bootstrap samples are generated from,

$$y_i^* = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}} + e_i^*, \quad i = 1, \dots, n,$$

and where $\{e_i^*\}_{i=1}^n$, are sampled with replacement from the centered residuals obtained from the initial lasso fit.

► Machine Learning Methods:

- Dezeure *et al.* (2015) critique this approach since there is non-uniform convergence to the limiting distribution and problems with the CIs both when $\beta_j = 0$ (zero-length intervals at 0) and when $\beta_j \neq 0$ (poor coverage and wide CIs).
- To obtain improved asymptotic behavior there has been much research on the **debiased (or de-sparsified)** lasso estimator that makes an adjustment to the original estimator (Van de Geer *et al.*, 2014).

FAILURE OF THE BOOTSTRAP

When does the bootstrap fail?

- ▶ **Extrema:** We have already discussed the example in which $Y \sim_{iid} \text{Uniform}(0, \theta)$, and the bootstrap points $\theta^{*(b)}$ put mass 0.632 on the largest point.
- ▶ **Small n :** The justification for the bootstrap follows an asymptotic argument, but will be accurate when we have “small” n .

This was based in part on

[https://notstatschat.rbind.io/2017/02/01/
when-the-bootstrap-doesnt-work/](https://notstatschat.rbind.io/2017/02/01/when-the-bootstrap-doesnt-work/)

Other examples can be found on this page.

DISCUSSION

- ▶ Check the small print when you want to use the bootstrap!

References

- Arcones, M. and Giné, E. (1992). On the bootstrap of M-estimators and other statistical functionals. In R. LePage and L. Billard, editors, *Exploring the Limits of Bootstrap*, New York. Wiley.
- Bickel, P. and Freedman, D. (1981). Some asymptotic theory for the bootstrap. *Annals of Statistics*, **9**, 1196–1217.
- Chatterjee, A. and Lahiri, S. N. (2013). Rates of convergence of the adaptive lasso estimators to the oracle distribution and higher order refinements by the bootstrap. *The Annals of Statistics*, **41**, 1232–1259.
- Chernick, M. R. (2011). *Bootstrap Methods: A Guide for Practitioners and Researchers*. John Wiley & Sons.
- Davison, A. and Hinkley, D. (1997). *Bootstrap Methods and their Application*. Cambridge University Press, Cambridge.
- Dezeure, R., Bühlmann, P., Meier, L., and Meinshausen, N. (2015). High-dimensional inference: confidence intervals, p -values and R-software `hdi`. *Statistical Science*, **30**, 533–558.
- Efron, B. and Hastie, T. (2016). *Computer Age Statistical Inference*. Cambridge University Press.

- Efron, B. and Tibshirani, R. (1993). *An Introduction to the Bootstrap*. Chapman and Hall/CRC, Boca Raton.
- Kim, J. and Pollard, D. (1990). Cube root asymptotics. *The Annals of Statistics*, pages 191–219.
- Kosorok, M. (2008). *Introduction to Empirical Processes and Semiparametric Inference*. Springer, New York.
- Lehmann, E. L. (1999). *Elements of large-sample theory*. Springer.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B*, **58**, 267–288.
- Van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, **42**, 1166–1202.