

ECON 480-3
LECTURE 17: THE BOOTSTRAP

Ivan A. Canay
Northwestern University

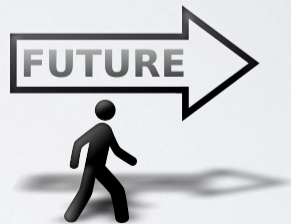
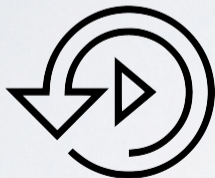


PART III SO FAR

- ▶ HC Standard Errors
- ▶ HAC Standard Errors
- ▶ CCE Standard Errors

TODAY

- ▶ Confidence Sets and Pivots
- ▶ Bootstrap: Algorithm
- ▶ Bootstrap: Sample Mean
- ▶ Discussion



CONFIDENCE SETS

- ▶ **Data:** let $X_i, i = 1, \dots, n$ be an i.i.d. sample of observations with distribution $P \in \mathbf{P}$.
- ▶ The **family** \mathbf{P} may be a parametric, nonparametric, or semiparametric family of distributions.
- ▶ We are interested in making inferences about some **parameter** $\theta(P) \in \Theta = \{\theta(P) : P \in \mathbf{P}\}$.
- ▶ Examples of $\theta(P)$ are the **mean** of P , the **median** of P , of a **regression coefficient**, among others.
- ▶ **Confidence set for $\theta(P)$:** a random set $C_n = C_n(X_1, \dots, X_n)$ such that
$$P\{\theta(P) \in C_n\} \approx 1 - \alpha,$$
at least for n sufficiently large.

- ▶ The typical way of constructing such sets is based off of approximating the distribution of a **root**,

$$R_n = R_n(X_1, \dots, X_n, \theta(P)) .$$

- ▶ **Root:** any real-valued function depending on both the data and the parameter of interest, $\theta(P)$.

SAMPLING DISTRIBUTION

- ▶ **Idea:** if the distribution of the root were **known**, then one could straightforwardly construct a confidence set for $\theta(P)$.

- ▶ Let $J_n(P)$ denote the sampling distribution of R_n and define the corresponding cdf as,

$$J_n(x, P) = P\{R_n \leq x\}.$$

- ▶ **Note:** the distribution of R_n depends on **both** the sample size, n , and the distribution of the data, P .

- ▶ Knowing $J_n(x, P)$ we may choose a **constant** c such that

$$P\{R_n \leq c\} \approx 1 - \alpha.$$

- ▶ The set $C_n = \{\theta \in \Theta : R_n(X_1, \dots, X_n, \theta) \leq c\}$ is a confidence set in the sense described above.

- ▶ We may also choose c_1 and c_2 so that

$$P\{c_1 \leq R_n \leq c_2\} \approx 1 - \alpha,$$

and construct the desired confidence set as

$$C_n = \{\theta \in \Theta : c_1 \leq R_n(X_1, \dots, X_n, \theta) \leq c_2\}.$$

PIVOTS

- ▶ In some **rare** instances, $J_n(x, P)$ does **not depend** on P .
- ▶ In these instances, the root is said to be **pivotal** or a **pivot**.
- ▶ **Example:** if $\theta(P)$ is the mean of P and $\mathbf{P} = \{N(\theta, 1) : \theta \in \mathbf{R}\}$, then the root

$$R_n = \sqrt{n}(\bar{X}_n - \theta(P))$$

is a **pivot** because $R_n \sim N(0, 1)$.

- ▶ In this case, we may construct confidence sets C_n with **finite-sample validity**; that is,

$$P\{\theta(P) \in C_n\} = 1 - \alpha$$

for all n and $P \in \mathbf{P}$.

ASYMPTOTIC PIVOTS

- ▶ Sometimes, the root may not be pivotal in the sense described above, but it may be **asymptotically pivotal** or an **asymptotic pivot**.
- ▶ **Asymptotic pivot:** $J_n(x, P)$ converges in distribution to a limit distribution $J(x, P)$ that **does not depend** on P .
- ▶ **Example:** $\theta(P)$ is the mean of P and \mathbf{P} is the set of all distributions on \mathbf{R} with a finite, nonzero variance, then

$$R_n = \frac{\sqrt{n}(\bar{X}_n - \theta(P))}{\hat{\sigma}_n}$$

is asymptotically pivotal because it converges in distribution to $J(x, P) = \Phi(x)$.

- ▶ We may then construct confidence sets that are asymptotically valid in the sense that

$$\lim_{n \rightarrow \infty} P\{\theta(P) \in C_n\} = 1 - \alpha$$

for all $P \in \mathbf{P}$.

ASYMPTOTIC APPROXIMATIONS

- ▶ Typically, the root will be **neither** a pivot nor an asymptotic pivot.
- ▶ The distribution $J_n(x, P)$ and the limiting distribution $J(x, P)$ will typically depend on P .
- ▶ **Example:** $\theta(P)$ is the mean of P and \mathbf{P} is the set of all distributions on \mathbf{R} with a finite, nonzero variance, then

$$R_n = \sqrt{n}(\bar{X}_n - \theta(P))$$

converges in distribution to $J(x, P) = \Phi(x/\sigma(P))$.

- ▶ We can approximate this limit distribution with $\Phi(x/\hat{\sigma}_n)$, which will lead to confidence sets that are asymptotically valid in the sense described above.
- ▶ This approach depends very heavily on the limit distribution $J(x, P)$ being both **known** and **tractable**. Even if it is known, the limit distribution may be **difficult to work with** (e.g., it could be the supremum of some complicated stochastic process with many nuisance parameters).
- ▶ Even if it is **known** and **manageable**, the method may be poor in finite-samples because it essentially relies on a **double approximation**: (1) $J_n(x, P)$ is approximated by $J(x, P)$, then (2) $J(x, P)$ is approximated in some way by estimating the unknown parameters of the limit distribution.

THE BOOTSTRAP

- ▶ The **bootstrap** is a fourth, more general approach to approximating $J_n(x, P)$.
- ▶ The **idea** is very simple: replace the unknown P with an **estimate** \hat{P}_n .
- ▶ Given \hat{P}_n it is possible to compute $J_n(x, \hat{P}_n)$ (either **analytically** or using **simulation** to any desired degree of accuracy).
- ▶ **Estimate:** ① In the case of i.i.d. data, a typical choice is the **empirical distribution**.

$$\begin{array}{cccccc} \text{Data} & X_1 & X_2 & X_3 & \cdots & X_n \\ \hat{P}_n & \frac{1}{n} & \frac{1}{n} & \frac{1}{n} & \cdots & \frac{1}{n} \end{array} \cdot$$

- ② If $P = P(\psi)$ for some **finite-dimensional** parameter ψ , then one may also use $\hat{P}_n = P(\hat{\psi}_n)$ for some estimate $\hat{\psi}_n$ of ψ .
- ▶ The hope is that whenever \hat{P}_n is **close** to P , $J_n(x, \hat{P}_n)$ is **close** to $J_n(x, P)$.
- ▶ Essentially, this requires that $J_n(x, P)$, when viewed as a function of P , is **continuous** in an appropriate neighborhood of P .
- ▶ Often, this turns out to be true, but, unfortunately, it is **not true in general**.

QUESTIONS?



THE CASE OF THE MEAN

- ▶ We will now consider the case where P is a distribution on \mathbf{R} and $\theta(P)$ is the **mean** of P .
- ▶ We will consider first the **root** $R_n = \sqrt{n}(\bar{X}_n - \theta(P))$.
- ▶ Let \hat{P}_n denote the **empirical distribution** of the $\{X_i : i = 1, \dots, n\}$.

Data	X_1	X_2	X_3	\dots	X_n	
\hat{P}_n	$\frac{1}{n}$	$\frac{1}{n}$	$\frac{1}{n}$	\dots	$\frac{1}{n}$.

- ▶ **Question:** under what conditions is $J_n(x, \hat{P}_n)$ “close” to $J_n(x, P)$?
- ▶ The sequence of distributions \hat{P}_n is a **random sequence**, so it is more convenient to answer the question first for a **nonrandom sequence** P_n and then extend to random sequences.
- ▶ Before presenting formal results, we discuss how to implement the bootstrap in practice for this case.

IMPLEMENTATION OF THE BOOTSTRAP

- ▶ Most often, the bootstrap approximation $J_n(x, \hat{P}_n)$ cannot be calculated **exactly**.
- ▶ We can **approximate** this distribution to an **arbitrary degree of accuracy** by taking samples from \hat{P}_n
- ▶ Let's consider two cases for the case of the mean:
$$\mathbf{P}_{\text{np}} = \{\text{distributions with finite second moments}\} \quad \text{and} \quad \mathbf{P}_{\text{p}} = \{P : \exp(1/\theta)\}.$$
- ▶ Implementation of the bootstrap requires **4 steps**.

BOOTSTRAP ALGORITHM

- 1 Conditional on the data (X_1, \dots, X_n) , draw **B samples** of **size n** from \hat{P}_n . Denote the j th sample by
$$(X_{1,j}^*, \dots, X_{n,j}^*) \quad \text{for } j = 1, \dots, \mathbf{B}.$$

Non-parametric: When $\mathbf{P} = \mathbf{P}_{\text{np}}$, \hat{P}_n is typically the empirical distribution and this involves resampling the original observations **with replacement** (each observation has probability $1/n$).

Parametric: When $\mathbf{P} = \mathbf{P}_{\text{p}}$, $\hat{P}_n = \exp(1/\hat{\theta}_n)$ where $\hat{\theta}_n$ is, for example, the MLE of θ : $\hat{\theta}_n = \bar{X}_n$.

NOTE! Step 1, how you take your samples, determines the **type** of bootstrap (parametric, non-parametric, Wild-bootstrap, clustered samples, and many others).

IMPLEMENTATION OF THE BOOTSTRAP

BOOTSTRAP ALGORITHM

- ② For each bootstrap sample j , compute **the root**, i.e.,

$$R_{j,n}^* = R_n(X_{1,j}^*, \dots, X_{n,j}^*, \hat{\theta}_n) \quad \text{for } j = 1, \dots, B.$$

Note: $\theta(\hat{P}_n) = \hat{\theta}_n$, so in the bootstrap distribution $\theta(P)$ becomes $\hat{\theta}_n$. In the case of the mean,

$$R_{j,n}^* = \sqrt{n}(\bar{X}_{j,n}^* - \bar{X}_n).$$

- ③ Compute the **empirical cdf** of $(R_{1,n}^*, \dots, R_{B,n}^*)$ as

$$L_n(x) = \frac{1}{B} \sum_{j=1}^B I\{R_{j,n}^* \leq x\}.$$

By the Glivenko-Cantelli Theorem, $\sup_{x \in \mathbf{R}} |L_n(x) - J(x, \hat{P}_n)| \rightarrow 0$ as $B \rightarrow \infty$. Thus, the

Glivenko-Cantelli Theorem implies we can achieve an arbitrary degree of accuracy by choosing B sufficiently big. In practice, researcher use $B = 1,000$ or above, given enough computational power.

IMPLEMENTATION OF THE BOOTSTRAP

BOOTSTRAP ALGORITHM

- 4 Compute the desired function of $L_n(x)$, for example, a quantile,

$$L_n^{-1}(1 - \alpha) = \inf\{x \in \mathbf{R} : L_n(x) \geq 1 - \alpha\},$$

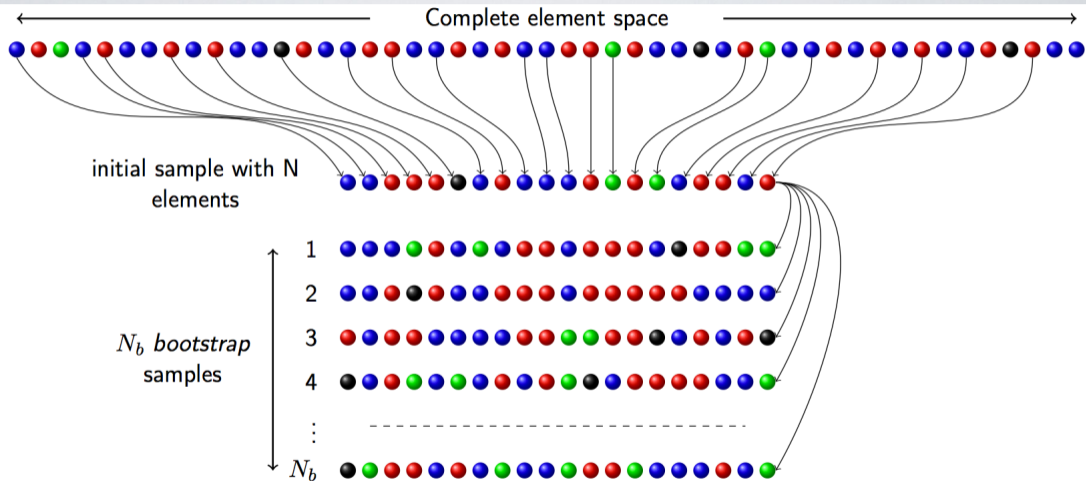
for a given significance level α .

Note! In practice you usually skip step 3 and compute step 4 directly using a `quantile` function:

$$L_n^{-1}(1 - \alpha) = \text{quantile}_{1-\alpha}\left(R_{1,n}^*, R_{2,n}^*, \dots, R_{B,n}^*\right).$$

- ▶ In Step 1, it matters the family of distributions (parameter vs nonparametric).
- ▶ In Step 2, it matters what the parameter of interest is and what the root is.
- ▶ Step 3 may be skipped if step 4 can be calculated directly as above.
- ▶ We are now ready to start with the **formal results**

THE BOOTSTRAP: ILLUSTRATION WITH MARBELS



QUESTIONS?



INTERMEDIATE STEP: NONRANDOM SEQUENCE

THEOREM

Let $\theta(P)$ be the **mean** of P and let \mathbf{P} denote the set of all distributions on \mathbf{R} with a finite, nonzero variance. Consider the root

$$R_n = \sqrt{n}(\bar{X}_n - \theta(P)) .$$

Let $P_n, n \geq 1$ be a **nonrandom sequence** of distributions such that:

- ① P_n converges in distribution to P
- ② $\theta(P_n) \rightarrow \theta(P)$
- ③ $\sigma^2(P_n) \rightarrow \sigma^2(P)$.

Then,

- (i) $J_n(x, P_n)$ converges in distribution to $J(x, P) = \Phi(x/\sigma(P))$.
- (ii) The $1 - \alpha$ quantile $J_n^{-1}(1 - \alpha, P_n) = \inf\{x \in \mathbf{R} : J_n(x, P_n) \geq 1 - \alpha\}$ converges to

$$J^{-1}(1 - \alpha, P) = z_{1-\alpha}\sigma(P) .$$

PROOF OF (I)

For each n , let $X_{i,n}, i = 1, \dots, n$ be i.i.d. with distribution P_n .

WTS: $\sqrt{n}(\bar{X}_{n,n} - \theta(P_n)) \xrightarrow{d} N(0, \sigma^2(P))$.

PROOF OF (I)

1. Suppose that Y_n and Y are real valued random variables and that $Y_n \xrightarrow{d} Y$.
If the Y_n are **uniformly bounded**, then $E[Y_n] \rightarrow E[Y]$.
2. Suppose that $Y_n \xrightarrow{d} Y$. Let g be a measurable map from \mathbf{R} to \mathbf{R} . Let C be the set of points in \mathbf{R} for which g is continuous. If $P\{Y \in C\} = 1$, then $g(Y_n) \xrightarrow{d} g(Y)$.

PROOF OF (I)

Claim that $\lim_{n \rightarrow \infty} E[Z_{n,i}^2 I\{|Z_{n,i}| > \epsilon\sqrt{n}\}] = 0$ and apply the Lindeberg-Feller CLT.

PROOF OF PART (II)

Part (II) follows from part (i) and Lemma 2 below applied to $F_n(x) = J_n(x, P)$ and $F(x) = J(x, P)$.

LEMMA

Let $F_n, n \geq 1$ and F be nonrandom distribution functions on \mathbf{R} such that F_n converges in distribution to F . Suppose F is **continuous** and **strictly increasing** at

$$F^{-1}(1 - \alpha) = \inf\{x \in \mathbf{R} : F(x) \geq 1 - \alpha\}.$$

Then,

$$F_n^{-1}(1 - \alpha) = \inf\{x \in \mathbf{R} : F_n(x) \geq 1 - \alpha\} \rightarrow F^{-1}(1 - \alpha).$$

Proof: see Lecture Notes.

MAIN THEOREM

We are now ready to pass from the **nonrandom** sequence P_n to the **random** sequence \hat{P}_n .

THEOREM

Let $\theta(P)$ be the mean of P and let \mathbf{P} denote the set of all distributions on \mathbf{R} with a finite, nonzero variance. Consider the root

$$R_n = \sqrt{n}(\bar{X}_n - \theta(P)) .$$

Then,

- (i) $J_n(x, \hat{P}_n)$ converges in distribution to $J(x, P) = \Phi(x/\sigma(P))$ a.s.
- (ii) $J_n^{-1}(1 - \alpha, \hat{P}_n)$ converges to $J^{-1}(1 - \alpha, P) = z_{1-\alpha}\sigma(P)$ a.s.

Proof: it follows from the previous theorem if we show that \hat{P}_n satisfies the requirements imposed on P_n a.s.:

- ① \hat{P}_n converges in distribution to P a.s.
- ② $\theta(\hat{P}_n) \rightarrow \theta(P)$ a.s.
- ③ $\sigma^2(\hat{P}_n) \rightarrow \sigma^2(P)$ a.s.

REMARKS

- ▶ Similar results hold for the **studentized root** where $\hat{\sigma}_n$ is consistent for $\sigma(P)$.
- ▶ Using this root leads to the so-called **Bootstrap- t** , as the root is just the t -statistic.
- ▶ A **key step** in the proof of this result is to show that $\hat{\sigma}_n$ converges in probability to $\sigma(P)$ under an appropriate sequence of distributions. We do this in 481.
- ▶ The **advantage** of working with a studentized root is that the limit distribution of R_n is **pivotal**, which affects the properties of the bootstrap approximation as discussed in the next section.
- ▶ **By Slutsky's Theorem** confidence sets of the form

$$C_n = \{ \theta \in \mathbf{R} : R_n(X_1, \dots, X_n, \theta) \leq J_n^{-1}(1 - \alpha, \hat{P}_n) \},$$

which are known as **symmetric** confidence sets, or

$$C_n = \left\{ \theta \in \mathbf{R} : J_n^{-1} \left(\frac{\alpha}{2}, \hat{P}_n \right) \leq R_n(X_1, \dots, X_n, \theta) \leq J_n^{-1} \left(1 - \frac{\alpha}{2}, \hat{P}_n \right) \right\},$$

which are known as **equi-tailed** confidence sets, satisfy

$$P\{\theta(P) \in C_n\} \rightarrow 1 - \alpha \quad \text{for all } P \in \mathbf{P}.$$

ASYMPTOTIC REFINEMENTS

- ▶ **Q:** is the bootstrap better than an asymptotically normal approximation?
- ▶ **Yes:** under certain conditions (ensuring existence of so-called Edgeworth expansions of $J_n(x, P)$) it follows that one-sided confidence sets C_n based off such an **asymptotic approximation** satisfy

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

- ▶ One-sided confidence sets based off of **the bootstrap** and the root $R_n = \sqrt{n}(\bar{X}_n - \theta(P))$ also satisfy (1), though there is some evidence to suggest that it does a bit better in the size of $O(n^{-1/2})$ term.
- ▶ On the other hand, one-sided confidence sets based off the bootstrap- t , i.e., using the root

$$R_n = \frac{\sqrt{n}(\bar{X}_n - \theta(P))}{\hat{\sigma}_n}$$

satisfy

$$|P\{\theta(P) \in C_n\} - (1 - \alpha)| = O\left(\frac{1}{n}\right). \quad (2)$$

- ▶ **Refinement:** the coverage error of the bootstrap- t interval is $O(n^{-1})$ and is of **smaller order** than that provided by the normal approximation or the bootstrap based on a nonstudentized root.

ASYMPTOTIC REFINEMENTS

- ▶ A heuristic reason why the bootstrap based on the studentized root outperforms the bootstrap based on the nonstudentized root is as follows.
 - ▶ **Nonstudentized case:** the bootstrap is estimating a distribution that has mean 0 and unknown variance $\sigma^2(P)$. The main contribution to the estimation error is the implicit estimation of $\sigma^2(P)$ by $\sigma^2(\hat{P}_n)$.
 - ▶ **Studentized case:** the studentized root has a distribution that is nearly independent of P since it is an asymptotic pivot.
- ▶ The bootstrap may also provide a refinement in **two-sided tests**. For example, symmetric intervals based on the absolute value of the studentized root are $O(n^{-2})$, versus the asymptotic approximation that is of order $O(n^{-1})$.
- ▶ Note that, by construction, such intervals are symmetric about $\hat{\theta}_n$.
- ▶ **Final comment:** The bootstrap **may fail** and work **poorly** in finite samples when $J(x, P)$ is **not** continuous in P . This happens in some settings relevant to economics: Auction models, Entry Game Models, Missing data settings, etc.

THE END!

