

Lecture 16: The Bootstrap

Jacob Steinhardt

March 11, 2021

- Overdispersion
 - Parametric confidence intervals \implies overly narrow uncertainty
 - Last time: can fix with negative binomial model
 - Are there more model-agnostic ways to fix this?

- Yes! (Sort of)
 - The bootstrap: a nonparametric method for generating confidence intervals
 - Can work even if CLT doesn't hold
 - But can sometimes fail, and need β to at least be meaningful

Recap of frequentist inference

Data $X_1, \dots, X_n \sim p$, parameter $\theta(p)$

Confidence interval at level α : $I(X_1, \dots, X_n)$ (interval on real line) such that

$$\mathbb{P}[\theta(p) \in I(X_{1:n})] \geq 1 - \alpha$$

More generally: **confidence region** satisfies $\theta(p) \in R(X_{1:n})$ w.p. $1 - \alpha$.

Note probability is over random draw of X_1, \dots, X_n (for fixed p).

Wald confidence ellipsoids for GLMs

Last time looked at `statsmodels` package, which uses the Wald ellipsoid:

$$R_\alpha(X_{1:n}) = \{z \mid (z - \hat{\beta}_n)^T I_n (z - \hat{\beta}_n) \leq F^{-1}(\alpha)\},$$

where $\hat{\beta}_n = \operatorname{argmin}_\beta L_n(\beta)$ is the maximum likelihood estimate, and $I_n = \nabla^2 L_n(\hat{\beta}_n)$ is the Fisher information.

Asymptotic normality implies that F is the cdf of the χ^2 distribution.

The above form is specific to maximum likelihood estimators, but similar confidence ellipsoids exist for any M-estimator.

Escaping model mis-specification

Saw last time that Wald confidence interval can be wrong if model is wrong

We'll escape this with a **non-parametric** tool for producing frequentist CIs

Non-parametric \implies doesn't rely on model \implies more robust

Key tool: the **bootstrap**

The Bootstrap

Idea for computing confidence intervals by resampling the data

Without bootstrap:

- Often rely on model assumptions
- Wald test, chi-square test, student-t test, . . .
- Lots of algebra, need different formula for each setting

With bootstrap:

- Fewer assumptions
- Single unified approach
- Computer simulation

Bootstrap: formal setting

Data: $X^{(1)}, \dots, X^{(n)} \sim p$

Estimator: $\hat{\theta} = \hat{\theta}(X^{(1)}, \dots, X^{(n)})$

- θ^* : population parameter (that $\hat{\theta}$ converges to as $n \rightarrow \infty$)

Question: How close is θ^* to $\hat{\theta}$?

- Typically framed as computing distribution of $\frac{1}{\sqrt{n}}(\hat{\theta} - \theta^*)$

The ideal hypothetical: re-sampling

Population distribution p^*

- $X^{(1)}, \dots, X^{(n)} \sim p^*$

The ideal hypothetical: re-sampling

Population distribution p^*

- $X^{(1)}, \dots, X^{(n)} \sim p^*$

Noise in $\hat{\theta}$ due to randomness in $X^{(1)}, \dots, X^{(n)}$

The ideal hypothetical: re-sampling

Population distribution p^*

- $X^{(1)}, \dots, X^{(n)} \sim p^*$

Noise in $\hat{\theta}$ due to randomness in $X^{(1)}, \dots, X^{(n)}$

Imagine hypothetically sampling fresh data:

$$X^{(1)}, \dots, X^{(n)} \rightarrow \hat{\theta} \text{ (Original sample)}$$

$$X^{(1)'}, \dots, X^{(n)'} \rightarrow \hat{\theta}' \text{ (Re-sample)}$$

$$X^{(1)''}, \dots, X^{(n)''} \rightarrow \hat{\theta}''$$

$$X^{(1)'''}, \dots, X^{(n)'''} \rightarrow \hat{\theta}'''$$

⋮

The ideal hypothetical: re-sampling

Population distribution p^*

- $X^{(1)}, \dots, X^{(n)} \sim p^*$

Noise in $\hat{\theta}$ due to randomness in $X^{(1)}, \dots, X^{(n)}$

Imagine hypothetically sampling fresh data:

$$X^{(1)}, \dots, X^{(n)} \rightarrow \hat{\theta} \text{ (Original sample)}$$

$$X^{(1)'}, \dots, X^{(n)'} \rightarrow \hat{\theta}' \text{ (Re-sample)}$$

$$X^{(1)''}, \dots, X^{(n)''} \rightarrow \hat{\theta}''$$

$$X^{(1)'''}, \dots, X^{(n)'''} \rightarrow \hat{\theta}'''$$

⋮

Implicit commitment: distribution of $\hat{\theta}$ roughly centered on θ^* (low bias)

Counterexample

$$\hat{\theta}(x_1, \dots, x_n) = \max_{i=1}^n x_i$$

n samples: always finite

∞ samples: infinite

The Bootstrap

Want to approximate hypothetical samples $\hat{\theta}', \hat{\theta}'', \dots$

But only have actual data $x^{(1)}, \dots, x^{(n)} \rightarrow \hat{\theta}$

Idea: subsample data

- With replacement
- n points in each sample

Useful framing: approximate n samples from p by n samples from \hat{p}_n

Bootstrap: Pseudocode

B : number of bootstrap samples

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

- Sample $x^{(1)'}, \dots, x^{(n)'}$ with replacement from $x^{(1)}, \dots, x^{(n)}$
- Let $\hat{\theta}^{(b)} = \hat{\theta}(x^{(1)'}, \dots, x^{(n)'})$

Output $\{\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(B)}\}$

Bootstrap in python

[Jupyter demos]

When does the bootstrap work?

Most parametric estimators are fine

- I.e. fixed number of parameters d and $d \ll n$

When does the bootstrap work?

Most parametric estimators are fine

- I.e. fixed number of parameters d and $d \ll n$

NOT parametric:

- Decision trees
- Neural nets
- Kernel regression

These “interpolate” data, sampling with replacement \approx subsampling

When does the bootstrap work?

Most parametric estimators are fine

- I.e. fixed number of parameters d and $d \ll n$

NOT parametric:

- Decision trees
- Neural nets
- Kernel regression

These “interpolate” data, sampling with replacement \approx subsampling

Other commitments:

- $\hat{\theta}$ approximately unbiased
- θ^* is a meaningful quantity

More examples

Bootstrap works for:

- Median and other quantiles
- Cumulative distribution function
- Trimmed mean
- Most U -statistics

Doesn't work for:

- De-generate U -statistics, e.g.: $U(X_{1:n}) = \frac{1}{\binom{n}{2}} \sum_{i < j} \mathbb{I}[X_i = X_j] e^{1/X_i}$
- Estimating θ for $X \sim \text{Uniform}([0, \theta])$.

Bootstrap: Underlying Theory

We seek to approximate the distribution of some quantity $R_n(X_1, \dots, X_n; \rho)$ for $X_{1:n} \sim \rho$

Let $\mathcal{L}(\rho)$ denote the limiting distribution as $n \rightarrow \infty$

For instance, $R(X_{1:n}, \rho) = \frac{\hat{\mu}_n - \mu(\rho)}{\sqrt{n\sigma(\rho)}}$, and $\mathcal{L}(\rho) = N(0, 1)$

Bootstrap replaces $R_n(X_{1:n}, \rho)$ with $R_n(X'_{1:n}, \hat{\rho}_n)$

- Intuitively this replaces $\mathcal{L}(\rho)$ with $\mathcal{L}(\hat{\rho}_n)$

Bootstrap: Underlying Theory

Bootstrap replaces $R_n(X_{1:n}, p)$ with $R_n(X'_{1:n}, \hat{p}_n)$

- Intuitively this replaces $\mathcal{L}(p)$ with $\mathcal{L}(\hat{p}_n)$

Issue is there are two limits happening at once. To make this work need:

- $R_n(q) \rightarrow \mathcal{L}(q)$ uniformly for q in a neighborhood of p
- The mapping $p \mapsto \mathcal{L}(p)$ is continuous

Proof sketch: uniform convergence means that for large n , $R_n(\hat{p}_n)$ will be very close in law to $\mathcal{L}(\hat{p}_n)$ (need uniformity since \hat{p}_n is changing). Then $\mathcal{L}(\hat{p}_n) \rightarrow \mathcal{L}(p)$ since $\hat{p}_n \rightarrow p$ and \mathcal{L} is continuous.

See Bickel and Freedman 1981, *Some Asymptotic Theory for the Bootstrap*.

Counterexamples Revisited

Nonparametric models (i.e. neural nets) fail because \mathcal{L} is not continuous

Other estimators can fail due to lack of uniformity.

- E.g. $X \sim U([0, \theta])$, take $R_n = \frac{\theta - X^{\max}}{n\theta}$. [Here X^{\max} is the max of the X_i]
- R_n converges to exponential distribution, but bootstrap samples have $X^{\max} = X'^{\max}$ with probability $1 - e^{-1}$.

Some models with growing dimension are actually fine. E.g. can have dimension $n^{1-\delta}$ in regression models and still have bootstrap work. See Mammen 1992, *Bootstrap, wild bootstrap, and asymptotic normality*.