

Subsampling vs Bootstrap

Dimitris N. Politis, Joseph P. Romano, Michael Wolf

$$R_n(x_n, \theta(P)) = \tau_n(\hat{\theta}_n - \theta(P))$$

Example:

$$\begin{aligned} \hat{\theta}_n &= \bar{X}_n, \tau_n = \sqrt{n}, \theta = EX = \mu(P) \\ \text{or } \hat{\theta} &= \min X_n, \tau_n = n, \theta(P) = \sup\{x : F(x) \leq 0\} \end{aligned}$$

Define: $J_n(P)$, the distribution of $\tau_n(\hat{\theta}_n - \theta(P))$ under P . For real $\hat{\theta}_n$,

$$J_n(x, P) \equiv \text{Prob}_P(\tau_n(\hat{\theta}_n - \theta(P)) \leq x)$$

Since P is unknown, $\theta(P)$ is unknown, and $J_n(x, P)$ is also unknown. The bootstrap estimate $J_n(x, P)$ by $J_n(x, \hat{P}_n)$, where \hat{P}_n is a consistent estimate of P in some sense. For example, take $\hat{P}_n(x) = \frac{1}{n} \sum_{i=1}^n 1(X_i \leq x)$ the empirical distribution:

$$\sup_x \left| \hat{P}_n(x) - P(x) \right| \xrightarrow{a.s.} 0$$

Similarly estimate $(1 - \alpha)$ th quantile of $J_n(x, P)$ by $J_n(x, \hat{P}_n)$: i.e. Estimate $J_n^{-1}(x, P)$ by $J_n^{-1}(x, \hat{P}_n)$.

Usually $J_n(x, \hat{P}_n)$ can't be explicitly calculated(although in some simple case it can be), use Monte Carlo approximation:

$$J_n(x, \hat{P}_n) \approx \frac{1}{B} \sum_{i=1}^B 1(\tau_n(\hat{\theta}_{n,i} - \hat{\theta}_n) \leq x)$$

for $\hat{\theta}_{n,i} = \hat{\theta}(X_{1,i}^*, \dots, X_{n,i}^*)$.

When bootstrap works(the meaning of "works"), for each x ,

$$\begin{aligned} J_n(x, \hat{P}_n) - J_n(x, P) &\xrightarrow{p} 0 \\ \implies J_n^{-1}(1 - \alpha, \hat{P}_n) - J_n^{-1}(1 - \alpha, P) &\xrightarrow{p} 0 \end{aligned}$$

When should Bootstrap "work"? Need local uniformity in weak convergence:

1. Usually $J_n(x, P) \rightarrow J(x, P)$.
2. Also usually $\hat{P}_n \rightarrow P$ a.s. in some sense, say $\sup_x \left| \hat{P}_n(x) - P(x) \right| \xrightarrow{a.s.} 0$.
3. Suppose for each sequence P_n s.t. $P_n \rightarrow P$, say $\sup_x \left| P_n - P \right| \rightarrow 0$, it is also true that $J_n(x, P_n) \rightarrow J(x, P)$, then it must be true that a.s. $J_n(x, \hat{P}_n) \rightarrow J(x, P)$
4. So it ends up having to show for $P_n \rightarrow P$, $J_n(x, P_n) \rightarrow J(x, P)$, use triangular array formulation.

Case when it works: sample mean with finite variance. It is known that:

1. $\sup_x \left| \hat{F}_n(x) - F(x) \right| \xrightarrow{a.s.} 0$.
2. $\theta(\hat{F}_n) = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{a.s.} \theta(F) = EX$.
3. $\sigma^2(\hat{F}_n) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \xrightarrow{a.s.} \sigma^2(F) = VarX$.
4. Use Linderberg-Feller for the triangular array, applied to the deterministic sequence of P_n such that: 1) $\sup_x \left| P_n(x) - P(x) \right| \rightarrow 0$; 2) $\theta(P_n) \rightarrow \theta(P)$; 3) $\sigma^2(P_n) \rightarrow \sigma^2(P)$, it can be shown that $\sqrt{n}(\bar{X}_n - \theta(P_n)) \xrightarrow{d} N(0, \sigma^2)$ under P_n .
5. Since \hat{P}_n satisfies 1,2,3 a.s., therefore a.s. $J_n(x, \hat{P}_n) \rightarrow J(x, P)$.

Therefore “local uniformity” of weak convergence is satisfied here.

Cases when bootstrap fails:

1. Order Statistics: $F \sim U(0, \theta)$, and $X_{(1)}, \dots, X_{(n)}$ is the order statistics of the sample, so $X_{(n)}$ is the maximum:

$$\begin{aligned} P\left(n \frac{\theta - X_{(n)}}{\theta} > x\right) &= P\left(X_{(n)} < \theta - \frac{\theta x}{n}\right) = P\left(X_i < \theta - \frac{\theta x}{n}\right)^n = \left(\frac{1}{\theta} \left(\theta - \frac{\theta x}{n}\right)\right)^n \\ &= \left(1 - \frac{x}{n}\right)^n \xrightarrow{n \rightarrow \infty} e^{-x} \end{aligned}$$

The bootstrap version:

$$P\left(n \left(X_{(n)} - X_{(n)}^*\right) / X_{(n)} = 0\right) = \left(1 - \left(1 - \frac{1}{n}\right)^n\right) \xrightarrow{n \rightarrow \infty} (1 - e^{-1}) \approx 0.63$$

2. Degenerate U-statistics: Take $w(x, y) = xy$, $\theta(F) = \int \int w(x, y) dF(x) dF(y) = \mu(F)^2$.

$$\hat{\theta}_n = \theta(\hat{F}_n) = \frac{1}{n(n-1)} \sum \sum_{i \neq j} X_i X_j$$

$$S(x) = \int xy dF(y) = x\mu(F)$$

If $\mu(F) \neq 0$ it is known that

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, 4\text{Var}(S(X))) = N(0, 4(\mu^2 EX^2 - \mu^4))$$

The bootstrap works.

But if $\mu(F) = 0 \implies \theta(F) = 0$:

$$\theta(\hat{F}_n) = \frac{1}{n(n-1)} \sum \sum_{i \neq j} X_i X_j = \bar{X}_n^2 - \frac{1}{n} \frac{1}{n-1} \sum_i (X_i - \bar{X}_n)^2 = \bar{X}_n^2 - \frac{S_n^2}{n}$$

$$n(\theta(\hat{F}_n) - \theta(F)) = n\bar{X}_n^2 - S_n^2 \xrightarrow{d} N(0, \sigma^2) - \sigma^2$$

However the bootstrap version of $n[\theta(\hat{F}_n^*) - \theta(\hat{F}_n)]$:

$$n\left(\left[\bar{X}_n^{*2} - \frac{1}{n}S_n^{*2}\right] - \left[\bar{X}_n^2 - \frac{1}{n}S_n^2\right]\right) = n\bar{X}_n^{*2} - S_n^{*2} - n\bar{X}_n^2 + S_n^2 \approx n(\bar{X}_n^{*2} - \bar{X}_n^2)$$

$$= [\sqrt{n}(\bar{X}_n^* - \bar{X}_n)]^2 + 2\sqrt{n}(\bar{X}_n^* - \bar{X}_n)\sqrt{n}\bar{X}_n \xrightarrow{d} N(0, \sigma^2)^2 + 2N(0, \sigma^2)\sqrt{n}\bar{X}_n$$

Subsampling: iid case: Y_i block of size b from (X_1, \dots, X_n) , $i = 1, \dots, q$, for $q = \binom{n}{b}$. Let $\hat{\theta}_{n,b,i} = \hat{\theta}(Y_i)$ calculated with the i th block of data. Use the empirical distribution of $\tau_b(\hat{\theta}_{n,b,i} - \hat{\theta})$ over the q pseudo-estimates to approximate the distribution of $\tau_n(\hat{\theta} - \theta)$: Approximate

$$J_n(x, P) = P(\tau_n(\hat{\theta}_n - \theta) \leq x)$$

by

$$L_{n,b}(x) = q^{-1} \sum_{i=1}^q 1(\tau_b(\hat{\theta}_{n,b,i} - \hat{\theta}_n) \leq x)$$

Claim: If $b \rightarrow \infty$, $b/n \rightarrow 0$, $\tau_b/\tau_n \rightarrow 0$, as long as $\tau_n (\hat{\theta} - \theta) \xrightarrow{d}$ something,

$$J_n(x, P) - L_{n,b}(x) \xrightarrow{p} 0$$

DIFFERENT MOTIVATION FOR SUBSAMPLING VS. BOOTSTRAP:

Subsampling: each subset of size b comes from the TRUE model. Since $\tau_n (\hat{\theta}_n - \theta) \xrightarrow{d} J(x, P)$, so as long as $b \rightarrow \infty$:

$$\tau_b (\hat{\theta}_b - \theta) \xrightarrow{d} J(x, P)$$

For n large, the distributions of $\tau_n (\hat{\theta}_n - \theta)$ and $\tau_b (\hat{\theta}_b - \theta)$ should be close.
But

$$\tau_b (\hat{\theta}_b - \theta) = \tau_b (\hat{\theta}_b - \hat{\theta}_n) + \tau_b (\hat{\theta}_n - \theta)$$

Since

$$\tau_b (\hat{\theta}_n - \theta) = O_p \left(\frac{\tau_b}{\tau_n} \right) = o_p(1)$$

The distributions of $\tau_b (\hat{\theta}_b - \theta)$ and $\tau_b (\hat{\theta}_b - \hat{\theta}_n)$ should be close. The distribution of $\tau_b (\hat{\theta}_b - \hat{\theta}_n)$ is estimated by the empirical distribution over $q = \binom{n}{b}$ pseudo-estimates.

Bootstrap: Recalculate the statistics from the ESTIMATED model \hat{P}_n . Given that \hat{P}_n is close to P , hopefully $J_n(x, \hat{P}_n)$ is close to $J_n(x, P)$ (Or to $J(x, P)$, the limit distribution).
But when bootstrap fails

$$\hat{P}_n \rightarrow P \not\Rightarrow J_n(x, \hat{P}_n) \rightarrow J(x, P)$$

Formal Proof of consistency of subsampling:

Assumptions: $\tau_n (\hat{\theta}_n - \theta) \xrightarrow{d} J(x, P)$, $b \rightarrow \infty$, $\frac{b}{n} \rightarrow 0$, $\frac{\tau_b}{\tau_n} \rightarrow 0$.

Need to show: $L_{n,b}(x) - J(x, P) \xrightarrow{p} 0$.

Since $\tau (\theta_n - \theta) \xrightarrow{p} 0$, it is enough to show

$$U_{n,b}(x) = q^{-1} \sum_{i=1}^q 1 \left(\tau_b (\hat{\theta}_{n,b,i} - \theta) \leq x \right) \xrightarrow{p} J(x, P)$$

$$U_{n,b}(x) - J(x, P) = U_{n,b}(x) - EU_{n,b}(x) + EU_{n,b}(x) - J(x, P)$$

Enough to show

$$U_{n,b}(x) - EU_{n,b}(x) \xrightarrow{p} 0$$

and

$$EU_{n,b}(x) - J(x, P) \rightarrow 0$$

But

$$EU_{n,b}(x) - J(x, P) = J_b(x, P) \rightarrow 0$$

$U_{n,b}(x)$ is a b th order U-statistics with kernel function bounded by $(-1, 1)$. Use Hoeffding exponential-type inequality(Serfling(1980), Thm A. p201):

$$P(U_{n,b}(x) - J_b(x, P) \geq \epsilon) \leq \exp\left(-2\frac{n}{b}\epsilon^2/[1 - (-1)]\right) = \exp\left(-\frac{n}{b}t^2\right) \rightarrow 0$$

as $\frac{n}{b} \rightarrow \infty$.

So

$$L_{n,b}(x) - J(x, P) = L_{n,b}(x) - U_{n,b}(x) + U_{n,b}(x) - J_b(x, P) + J_b(x, P) - J(x, P) \xrightarrow{p} 0.$$

Q.E.D.

Time Series(!): Respect the ordering of the data to preserve correlation.

$$\hat{\theta}_{n,b,t} = \hat{\theta}_b(X_t, \dots, X_{t+b-1}), \quad q = T - b + 1.$$

$$L_{n,b}(x) = \frac{1}{q} \sum_{i=1}^q 1\left(\tau_b\left(\hat{\theta}_{n,b,t} - \hat{\theta}_n\right) \leq x\right)$$

Assumption: $\tau_n\left(\hat{\theta}_n - \theta\right) \xrightarrow{d} J(x, P)$, $b \rightarrow \infty$, $\frac{b}{n} \rightarrow 0$, $\frac{\tau_b}{\tau_n} \rightarrow 0$, $\alpha(m) \rightarrow 0$.

Result: $L_{n,b}(x) - J(x, P) \xrightarrow{p} 0$.

Most difficult part: To show $\tau_n\left(\hat{\theta}_n - \theta\right) \xrightarrow{d} J(x, P)$.

Can treat iid data as time series, or even using non-overlapping blocks $k = \lfloor \frac{n}{b} \rfloor$, but using $\binom{n}{b}$ more efficient. For example, if $\bar{U}_n(x) = k^{-1} \sum_{j=1}^k 1\left(\tau_b[R_{n,b,j} - \theta(P)] \leq x\right)$ then

$$U_{n,b}(x) = E\left[\bar{U}_n(x) | \mathcal{X}_n\right] = E\left[1\left(\tau_b[R_{n,b,j} - \theta(P)] \leq x\right) | \mathcal{X}_n\right]$$

for $\mathcal{X}_n = (X_{(1)}, \dots, X_{(n)})$. $U_{n,b}(x)$ is better than $\bar{U}_n(x)$ since \mathcal{X}_n is sufficient statistics for iid data.

Hypothesis Testing: $T_n = \tau_n t_n(X_1, \dots, X_n)$,

$$G_n(x, P) = \text{Prob}_p(\tau_n \leq x) \xrightarrow{P \in P_0} J(x, P)$$

$$\hat{G}_{n,b}(x) = q^{-1} \sum_{i=1}^q 1(T_{n,b,i} \leq x) = q^{-1} \sum_{i=1}^q 1(\tau_b t_{n,b,i} \leq x)$$

As long as $b \rightarrow \infty$, $\frac{b}{n} \rightarrow 0$, then under $P \in P_0$:

$$\hat{G}_{n,b}(x) \rightarrow G(x, P)$$

If under $P \in P_1$, $T_n \rightarrow \infty$, then $\forall x$, $\hat{G}_{n,b}(x) \rightarrow 0$.

Key difference with confidence interval: don't need $\frac{\tau_b}{\tau_n} \rightarrow 0$, because don't need to estimate θ_0 but assumed known under the null hypothesis.

Estimating the unknown rate of convergence: Assume that $\tau_n = n^\beta$, for some $\beta > 0$, but β is unknown. Estimate β using different size of subsampling distribution. Key idea: Compare the shape of the empirical distributions of $\hat{\theta}_b - \hat{\theta}_n$ for different values of b to infer the value of β .

Let $q = \binom{n}{b}$ for iid data, or $q = (T - b + 1)$ for time series data:

$$L_{n,b}(x|\tau_b) \equiv q^{-1} \sum_{a=1}^q 1\left(\tau_b \left(\hat{\theta}_{n,b,a} - \hat{\theta}_n\right) \leq x\right)$$

$$L_{n,b}(x|1) \equiv q^{-1} \sum_{a=1}^q 1\left(\hat{\theta}_{n,b,a} - \hat{\theta}_n \leq x\right)$$

This implies

$$L_{n,b}(x|\tau_b) = L_{n,b}(\tau_b^{-1}x|1) \equiv t$$

$$x = L_{n,b}^{-1}(t|\tau_b) = \tau_b (\tau_b^{-1}x) = \tau_b L_{n,b}^{-1}(t|1)$$

Since $L_{n,b}(x|\tau_b) \xrightarrow{P} J(x, P)$, if $J(x, P)$ is continuous and increasing, it can be inferred that

$$L_{n,b}^{-1}(t|\tau_b) = J^{-1}(t, P) + o_p(1)$$

Same as

$$\tau_b L_{n,b}^{-1}(t|1) = J^{-1}(t, P) + o_p(1)$$

So

$$b^\beta L_{n,b}^{-1}(t|1) = J^{-1}(t, P) + o_p(1)$$

take $\log(\text{Assuming } J^{-1}(t, P) > 0, \text{ or } t > J(0, P))$, for different b_1 and b_2 , then this becomes

$$\begin{aligned}\beta \log b_1 + \log(L_{n,b_1}^{-1}(t|1)) &= \log J^{-1}(t, P) + o_p(1) \\ \beta \log b_2 + \log(L_{n,b_2}^{-1}(t|1)) &= \log J^{-1}(t, P) + o_p(1)\end{aligned}$$

Different out the “fixed effect”

$$\beta (\log b_1 - \log b_2) = \log(L_{n,b_2}^{-1}(t|1)) - \log(L_{n,b_1}^{-1}(t|1)) + o_p(1)$$

So estimate β by

$$\hat{\beta} = (\log b_1 - \log b_2)^{-1} (\log(L_{n,b_2}^{-1}(t|1)) - \log(L_{n,b_1}^{-1}(t|1))) = \beta + (\log b_1 - \log b_2)^{-1} \times o_p(1)$$

Take $b_1 = n^{\gamma_1}$, $b_2 = n^{\gamma_2}$, ($1 \geq \gamma_1 > \gamma_2 > 0$)

$$\hat{\beta} - \beta = ((\gamma_1 - \gamma_2) \log n)^{-1} o_p(1) = o_p((\log n)^{-1})$$

How to know $t > J(0, P)$

$$L_{n,b}(0|\tau_b) = L_{n,b}(0|1) = J(0, P) + o_p(1)$$

So estimating $J(0, P)$ not a problem.

Alternatively, take $t_2 \in (0.5, 1)$, take $t_1 \in (0, 0.5)$

$$b^\beta (L_{n,b}^{-1}(t_2|1) - L_{n,b}^{-1}(t_1|1)) = J^{-1}(t_2|P) - J^{-1}(t_1|P) + o_p(1)$$

$$\beta \log b + \log(L_{n,b}^{-1}(t_2|1) - L_{n,b}^{-1}(t_1|1)) = \log(J^{-1}(t_2|P) - J^{-1}(t_1|P)) + o_p(1)$$

$$\hat{\beta} = (\log b_1 - \log b_2)^{-1} [\log(L_{n,b_2}^{-1}(t_2|1) - L_{n,b_2}^{-1}(t_1|1)) - \log(L_{n,b_1}^{-1}(t_2|1) - L_{n,b_1}^{-1}(t_1|1))]$$

Take $b_1 = n^{\gamma_1}$, $b_2 = n^{\gamma_2}$, ($1 > \gamma_1 > \gamma_2 > 0$), As before

$$\hat{\beta} - \beta = o_p((\log n)^{-1})$$

Two step subsampling: $\hat{\tau}_n = n^{\hat{\beta}}$

$$L_{n,b}(x|\hat{\tau}_b) = q^{-1} \sum_{a=1}^q 1 \left(\hat{\tau}_b \left(\hat{\theta}_{n,b,a} - \hat{\theta}_n \right) \leq x \right)$$

Can show that

$$\sup_x \left| L_{n,b}(x|\hat{\tau}_b) - J(x, P) \right| \xrightarrow{p} 0.$$

Problem: imprecise in small samples. E.g. in variation estimation, best choice of b gives error rate of $O(n^{-1/3})$ but parameter estimates, if model is true, gives $O(n^{-1/2})$ error rate. Bootstrap pivotal statistics, when applicable, gives even better than $O(n^{-1/2})$ error rate.