

Lecture 23: Variance estimation, replication, jackknife, and bootstrap

Motivation

To evaluate and compare different estimators, we need consistent estimators of variances or asymptotic variances of estimators.

This is also important for hypothesis testing and confidence sets.

Let $\text{Var}(\hat{\theta})$ be the variance or asymptotic variance of an estimator $\hat{\theta}$.

Traditional approach to estimate $\text{Var}(\hat{\theta})$: Derivation and substitution

- First, we derive a theoretical formula
- Approximation (asymptotic theory) is usually needed
- The formula may depend on unknown quantities
- We then substitute unknown quantities by estimators

Example: the δ -method

Y_1, \dots, Y_n are iid (k -dimensional)

$\theta = g(\mu)$ (e.g., a ratio of two components of μ), $\hat{\theta} = g(\bar{Y})$

$\text{Var}(\hat{\theta}) \approx [\nabla g(\mu)]^T \text{Var}(\bar{Y}) \nabla g(\mu)$

An estimator of $\text{Var}(\hat{\theta})$ is $\hat{V}_n = [\nabla g(\bar{Y})]^T (S^2/n) \nabla g(\bar{Y})$

By the SLLN, $\bar{X} \rightarrow_{a.s.} \mu$ and $S^2 \rightarrow_{a.s.} \text{Var}(X_1)$.

Hence, \hat{V}_n is strongly consistent for $\text{Var}(\hat{\theta}_n)$, provided that $\nabla g(\mu) \neq 0$ and ∇g is continuous at μ .

Example 5.15

Let Y_1, \dots, Y_n be i.i.d. random variables with finite $\mu_y = EY_1$, $\sigma_y^2 = \text{Var}(Y_1)$, $\gamma_y = EY_1^3$, and $\kappa_y = EY_1^4$.

Consider the estimation of $\theta = (\mu_y, \sigma_y^2)$.

Let $\hat{\theta}_n = (\bar{X}, \hat{\sigma}_y^2)$, where $\hat{\sigma}_y^2 = n^{-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$.

If $X_i = (Y_i, Y_i^2)$, then $\hat{\theta}_n = g(\bar{X})$ with $g(x) = (x_1, x_2 - x_1^2)$

$$\text{Var}(X_1) = \begin{pmatrix} \sigma_y^2 & \gamma_y - \mu_y(\sigma_y^2 + \mu_y^2) \\ \gamma_y - \mu_y(\sigma_y^2 + \mu_y^2) & \kappa_y - (\sigma_y^2 + \mu_y^2)^2 \end{pmatrix}$$

and

$$\nabla g(x) = \begin{pmatrix} 1 & 0 \\ -2x_1 & 1 \end{pmatrix}.$$

The estimator \hat{V}_n is strongly consistent, since $\nabla g(x)$ is continuous.

Is the derivative ∇g always easy to derive?

An alternative method?

Suppose we can independently obtain B copies of the data set X

Say X^1, \dots, X^B

Then we can calculate $\hat{\theta}^b = \hat{\theta}(X^b)$, $b = 1, \dots, B$

Variance of $\hat{\theta}$ can be estimated as

$$\frac{1}{B} \sum_{b=1}^B \left(\hat{\theta}^b - \frac{1}{B} \sum_{l=1}^B \hat{\theta}^l \right)^2$$

No derivation is needed (at the expense of more computations)

These estimators are valid for large B ($B \rightarrow \infty$, law of large numbers).

But typically, we only have one dataset, X .

Pseudo-replications

Can we apply the same idea by creating B pseudo-replicate datasets?

This means X^1, \dots, X^B are “copies” of X , but they are not independent of X (in fact, they are dependent on X)

These methods are called *resampling* or *replication* methods.

Consider pseudo replicates $X^i = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$, $i = 1, \dots, n$. Let $\hat{\theta}_{-i}$ be the same estimator as $\hat{\theta}_n$ but based on X^i , $i = 1, \dots, n$. Since $\hat{\theta}_n$ and $\hat{\theta}_{-1}, \dots, \hat{\theta}_{-n}$ estimate the same quantity, the following "sample variance" can be used as a measure of the variation of $\hat{\theta}_n$:

$$\frac{1}{n-1} \sum_{i=1}^n \left(\hat{\theta}_{-i} - \bar{\theta}_n \right)^2, \quad \bar{\theta}_n = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{-i}$$

Two issues:

- $\hat{\theta}_{-i}$'s are not independent.
- $\hat{\theta}_{-i} - \hat{\theta}_{-j}$ usually converges to 0 at a fast rate (such as n^{-1}).

If $\hat{\theta}_n = \bar{X}$ is the sample mean, then $\hat{\theta}_{-i} - \bar{\theta}_n = (n-1)^{-1}(\bar{X} - X_i)$ and

$$\frac{1}{n-1} \sum_{i=1}^n \left(\hat{\theta}_{-i} - \bar{\theta}_n \right)^2 = \frac{1}{(n-1)^3} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{S^2}{(n-1)^2}$$

Thus, the correction factor $(n-1)^2/n$ should be multiplied, which leads to the *jackknife variance estimator* of $\text{Var}(\hat{\theta}_n)$

$$\hat{V}_J = \frac{n-1}{n} \sum_{i=1}^n \left(\hat{\theta}_{-i} - \bar{\theta}_n \right)^2.$$

Theorem 5.17.

Let X_1, \dots, X_n be i.i.d. random d -vectors from F with finite $\mu = E(X_1)$ and $\text{Var}(X_1)$, and let $\hat{\theta}_n = g(\bar{X})$. Suppose that ∇g is continuous at μ and $\nabla g(\mu) \neq 0$. Then the jackknife variance estimator \hat{V}_J is strongly consistent for $\text{Var}(\hat{\theta}_n)$.

Proof.

Let \bar{X}_{-i} be the sample mean based on $X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n$. From the mean-value theorem, we have

$$\begin{aligned}\hat{\theta}_{-i} - \hat{\theta}_n &= g(\bar{X}_{-i}) - g(\bar{X}) \\ &= [\nabla g(\xi_{n,i})]^\tau (\bar{X}_{-i} - \bar{X}) \\ &= [\nabla g(\bar{X})]^\tau (\bar{X}_{-i} - \bar{X}) + R_{n,i},\end{aligned}$$

where $R_{n,i} = [\nabla g(\xi_{n,i}) - \nabla g(\bar{X})]^\tau (\bar{X}_{-i} - \bar{X})$ and $\xi_{n,i}$ is a point on the line segment between \bar{X}_{-i} and \bar{X} .

From $\bar{X}_{-i} - \bar{X} = (n-1)^{-1}(\bar{X} - X_i)$, it follows that $\sum_{i=1}^n (\bar{X}_{-i} - \bar{X}) = 0$ and

$$\frac{1}{n} \sum_{i=1}^n (\hat{\theta}_{-i} - \hat{\theta}_n) = \frac{1}{n} \sum_{i=1}^n R_{n,i} = \bar{R}_n.$$

From the definition of the jackknife estimator,

$$\widehat{V}_J = A_n + B_n + 2C_n,$$

where

$$A_n = \frac{n-1}{n} [\nabla g(\bar{X})]^\tau \sum_{i=1}^n (\bar{X}_{-i} - \bar{X})(\bar{X}_{-i} - \bar{X})^\tau \nabla g(\bar{X}),$$

$$B_n = \frac{n-1}{n} \sum_{i=1}^n (R_{n,i} - \bar{R}_n)^2,$$

$$C_n = \frac{n-1}{n} \sum_{i=1}^n (R_{n,i} - \bar{R}_n) [\nabla g(\bar{X})]^\tau (\bar{X}_{-i} - \bar{X}).$$

By $\bar{X}_{-i} - \bar{X} = (n-1)^{-1}(\bar{X} - X_i)$, the SLLN, and the continuity of ∇g at μ ,

$$A_n / \text{Var}(\widehat{\theta}_n) \rightarrow_{a.s.} 1.$$

Also,

$$(n-1) \sum_{i=1}^n \|\bar{X}_{-i} - \bar{X}\|^2 = \frac{1}{n-1} \sum_{i=1}^n \|X_i - \bar{X}\|^2 = O(1) \text{ a.s.} \quad (1)$$

Hence $\max_{i \leq n} \|\bar{X}_{-i} - \bar{X}\|^2 \rightarrow_{a.s.} 0$, which, together with the continuity of ∇g at μ and $\|\xi_{n,i} - \bar{X}\| \leq \|\bar{X}_{-i} - \bar{X}\|$, implies that

$$u_n = \max_{i \leq n} \|\nabla g(\xi_{n,i}) - \nabla g(\bar{X})\| \rightarrow_{a.s.} 0.$$

From (1), $\sum_{i=1}^n \|\bar{X}_{-i} - \bar{X}\|^2 / \text{Var}(\hat{\theta}_n) = O(1)$ a.s. and

$$\frac{B_n}{\text{Var}(\hat{\theta}_n)} \leq \frac{n-1}{\text{Var}(\hat{\theta}_n)n} \sum_{i=1}^n R_{n,i}^2 \leq \frac{u_n}{\text{Var}(\hat{\theta}_n)} \sum_{i=1}^n \|\bar{X}_{-i} - \bar{X}\|^2 \rightarrow_{a.s.} 0.$$

By the Cauchy-Schwarz inequality,

$(C_n/V_n)^2 \leq (A_n/V_n)(B_n/V_n) \rightarrow_{a.s.} 0$, which completes the proof.

Discussion

- A key step in the proof is that $\hat{\theta}_{-i} - \hat{\theta}_n$ can be approximated by $[\nabla g(\bar{X})]^\tau (\bar{X}_{-i} - \bar{X})$, which indicates that \hat{V}_J is consistent for $\hat{\theta}_n$ that can be well approximated by some linear statistic.
- More results can be found in Shao and Tu (1995, Chapter 2).
- The jackknife method can be applied to non-i.i.d. problems.

Bootstrap

Create bootstrap pseudo-replicate datasets X^{*1}, \dots, X^{*B} randomly generated from X .

Let $\hat{\theta}^{*b}$ be the same as an estimator $\hat{\theta}$ but based on X^{*b} , $b = 1, \dots, B$.

Is

$$\frac{1}{B} \sum_{b=1}^B \left(\hat{\theta}^{*b} - \frac{1}{B} \sum_{l=1}^B \hat{\theta}^{*l} \right)^2 \quad (2)$$

still a valid estimator of $\text{Var}(\hat{\theta})$?

In fact, the cdf $G(t) = P(\hat{\theta} - \theta \leq t)$ can be estimated as

$$\frac{1}{B} \sum_{b=1}^B I(\hat{\theta}^{*b} - \hat{\theta} \leq t) = \frac{\# \text{ of } b\text{'s such that } \hat{\theta}^{*b} - \hat{\theta} \leq t}{B} \quad (3)$$

where $I(A)$ is the indicator function of A .

The answer to this question depends on

- how the sample X is taken
- how X^1, \dots, X^B are constructed
- the type of the estimator, $\hat{\theta}$

A heuristic description for the bootstrap

\mathcal{P} : the population producing data X

$\widehat{\mathcal{P}}$: an estimated of the population based on data X

X^* : the bootstrap data produced by $\widehat{\mathcal{P}}$

$$\text{real world} \quad \mathcal{P} \Rightarrow X \Rightarrow \hat{\theta} = \hat{\theta}(X)$$

$$\text{bootstrap} \quad \widehat{\mathcal{P}} \Rightarrow X^* \Rightarrow \hat{\theta}^* = \hat{\theta}(X^*)$$

$\text{Var}(\hat{\theta})$ and $G(t) = P(\hat{\theta} - \theta \leq t)$ can be approximated by $\text{Var}_*(\hat{\theta}^*)$ and $\widehat{G}(t) = P_*(\hat{\theta}^* - \hat{\theta} \leq t)$, respectively, where the variance and probability are taken under the bootstrap sampling conditioned on X .

If $\widehat{\mathcal{P}}$ is close to \mathcal{P} , then

- $\widehat{G}(t)$ is close to $G(t)$;
- $\text{Var}_*(\hat{\theta}^*)$ is close to $\text{Var}(\hat{\theta})$.

Note that $\text{Var}_*(\hat{\theta}^*)$ and $\widehat{G}(t)$ are functions of X and are estimators.

If they have explicit forms, then they can be directly used.

If not, then we approximate them by the Monte Carlo approximations (2) and (3), respectively, based on bootstrap data sets X^{*1}, \dots, X^{*B} (copies of X^*).

How do we generate X^* based on X ?

Parametric bootstrap

Let X_1, \dots, X_n be iid with a cdf F_θ where θ is an unknown parameter vector and F_θ is known when θ is known.

Let $\hat{\theta}$ be an estimator of θ based on $X = (X_1, \dots, X_n)$.

Parametric bootstrap data set $X^* = (X_1^*, \dots, X_n^*)$ is obtained by generate iid X_1^*, \dots, X_n^* from $F_{\hat{\theta}}$.

Example: location-scale problems

Let $F_\theta(x) = F_0\left(\frac{x-\mu}{\sigma}\right)$, where $\mu = E(X_1)$, $\sigma^2 = \text{Var}(X_1)$ and F_0 is a known cdf.

Let \bar{X} be the sample mean, S^2 be the sample variance, and

$$T = \frac{\sqrt{n}(\bar{X} - \mu)}{S} = \sqrt{n} \sum_{i=1}^n \frac{X_i - \mu}{S}$$

The distribution of T does not depend on any parameter.

It is the t-distribution with degrees of freedom $n - 1$ if F_0 is normal.

Otherwise its explicit form is unknown.

Let $\hat{\theta} = (\bar{X}, S^2)$ Generate iid X_i^* , $i = 1, \dots, n$, from $F_{\hat{\theta}}$.
Then $(X_i^* - \bar{X})/S \sim F_0$

$$T^* = \sqrt{n} \sum_{i=1}^n \frac{X_i^{*b} - \bar{Y}}{S} \sim T$$

The parametric bootstrap is perfect: $\text{Var}^*(T^*) = \text{Var}(T)$.

If we calculate $\text{Var}^*(T^*)$ by Monte Carlo approximation, then the parametric bootstrap is exactly the same as the simulation approach. In general, if there is a function τ such that

$$\text{Var}_{\theta}(\hat{\theta}) = \tau(\theta), \quad X_1, \dots, X_n \text{ are iid from } F_{\theta}$$

then

$$\text{Var}_{\hat{\theta}}^*(\hat{\theta}^*) = \tau(\hat{\theta}), \quad X_1^*, \dots, X_n^* \text{ are iid from } F_{\hat{\theta}}$$

Hence, the parametric bootstrap is simply the substitution approach.

If $\hat{\theta}$ is consistent and τ is continuous, then $\text{Var}_{\hat{\theta}}^*(\hat{\theta}^*)$ is consistent.

If τ does not have a close form, we apply Monte Carlo approximation.

In the location-scale example, $\tau =$ a constant and hence the bootstrap is perfect.

Example

Let X_1, \dots, X_n be iid from F_θ .

Define $\mu = \mu(\theta) = E_\theta(X_1)$, $\mu_j = \mu_j(\theta) = E_\theta(X_1 - \mu)^j$, $j = 2, 3, 4$.

Consider the estimation of μ^2 by \bar{X}^2 .

A direct calculation shows that

$$\text{Var}_\theta(\bar{X}^2) = \frac{4[\mu(\theta)]^2\mu_2(\theta)}{n} + \frac{4\mu(\theta)\mu_3(\theta)}{n^2} + \frac{\mu_4(\theta)}{n^3}$$

Based on the previous discussion, the parametric bootstrap variance estimator is

$$\text{Var}_\theta^*(\bar{X}^{*2}) = \frac{4[\mu(\hat{\theta})]^2\mu_2(\hat{\theta})}{n} + \frac{4\mu(\hat{\theta})\mu_3(\hat{\theta})}{n^2} + \frac{\mu_4(\hat{\theta})}{n^3}$$

It is a consistent estimator if μ , μ_j , $j = 2, 3, 4$, are continuous functions.

If we apply the asymptotic approach, then we estimate $\text{Var}_\theta(\bar{X}^2)$ by

$$\frac{4[\mu(\hat{\theta})]^2\mu_2(\hat{\theta})}{n}$$

Nonparametric bootstrap

Without any model, we can apply the simple nonparametric bootstrap. If $X = (X_1, \dots, X_n)$, X_1, \dots, X_n are iid, then \mathcal{P} is the cdf of X_1 and $\widehat{\mathcal{P}}$ is the empirical cdf based on X_1, \dots, X_n .

If we generate iid bootstrap data X_1^*, \dots, X_n^* from $\widehat{\mathcal{P}}$, then it is the same as taking a simple random sample with replacement from X .

Property of $\text{Var}^*(\widehat{\theta}^*)$

Consider first $\widehat{\theta} = \bar{X}$, the sample mean, $\widehat{\theta}^* = \bar{X}^*$, the sample mean of X_1^*, \dots, X_n^* .

$$E^*(\bar{X}^*) = \frac{1}{n} \sum_{i=1}^n E^*(X_i^*) = \frac{1}{n} \sum_{i=1}^n \bar{X} = \bar{X}$$

$$\begin{aligned} \text{Var}^*(\bar{X}^*) &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}^*(X_i^*) = \frac{1}{n^2} \sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X})(X_j - \bar{X})^\tau \\ &= \frac{1}{n^2} \sum_{j=1}^n (X_j - \bar{X})(X_j - \bar{X})^\tau = \frac{n-1}{n^2} S^2 \approx \frac{S^2}{n} \end{aligned}$$

When n is small, we may make an adjustment of $\frac{n}{n-1}$.

Consider next the estimation of $g(\mu)$, where $\mu = E(X_1)$ and g is a continuously differentiable function.

Our estimator is $\hat{\theta} = g(\bar{X})$.

The bootstrap analog is $\hat{\theta}^* = g(\bar{X}^*)$.

When n is large,

$$g(\bar{X}^*) = g(\bar{X}) + \nabla g(\bar{X})(\bar{X}^* - \bar{X}) + \dots \approx g(\bar{X}) + \nabla g(\bar{X})(\bar{X}^* - \bar{X})$$

and

$$\begin{aligned}\text{Var}^*(\hat{\theta}^*) &= \text{Var}^*[g(\bar{X}^*)] \\ &\approx \nabla g(\bar{X}) \text{Var}^*(\bar{X}^* - \bar{X}) \nabla g^T(\bar{X}) \\ &\approx \frac{n-1}{n^2} \nabla g(\bar{X}) S^2 \nabla g^T(\bar{X})\end{aligned}$$

Example

Let X_1, \dots, X_n be iid from F .

Define $\mu = E_\theta(X_1)$, $\mu_j = E_\theta(X_1 - \mu)^j$, $j = 2, 3, 4$.

Consider the estimation of μ^2 by \bar{X}^2 :

$$\text{Var}(\bar{X}^2) = \frac{4\mu^2\mu_2}{n} + \frac{4\mu\mu_3}{n^2} + \frac{\mu_4}{n^3}$$

and

$$\text{Var}^*(\bar{X}^{*2}) = \frac{4\bar{X}^2 m_2}{n} + \frac{4\bar{X} m_3}{n^2} + \frac{m_4}{n^3}$$

where

$$m_j = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^j, \quad j = 2, 3, 4.$$

This is because the mean of the empirical cdf \hat{F}_n is \bar{X} and the j th central moment of \hat{F}_n is m_j .

In this case, we have an explicit form for the bootstrap variance estimator $\text{Var}^*(\bar{X}^{*2})$ so no Monte Carlo is needed.

This bootstrap variance estimator is consistent, since sample moments m_j 's are consistent for μ_j 's, by the WLLN.

Since $g'(x) = 2x$ when $g(x) = x^2$, the use of the approximation derived earlier gives that

$$\text{Var}^*(\bar{X}^{*2}) \approx \frac{4\bar{X}^2 m_2}{n}$$

which is also consistent since the terms ignored are of the orders n^{-2} .

In fact, the delta-method produces the variance estimator

$$\frac{[g'(\bar{X})]^2 S^2}{n} = \frac{4\bar{X}^2 S^2}{n}$$

Discussion

- In general, the expression $\text{Var}^*(\hat{\theta}^*)$ is usually complicated and not explicit, so Monte Carlo approximation is necessary.
- In fact, the idea of using the bootstrap is not to derive its explicit form (since it involves complex derivations).
- The bootstrap is to replace theoretical derivations by repeated computations.
- In many cases the theoretical derivations are difficult or messy.
- The user does not need to do theoretical derivations.
- However, they should be told when using the bootstrap produces correct variance estimators and how to do the bootstrap.
- The research on the bootstrap methodology still requires theoretical derivations.
- The jackknife shares the same idea as the bootstrap in some aspects, e.g., replacing theoretical derivations by repeated computations, but the bootstrap is a more complete methodology.