

# An Essay on Linear Regression Bootstrapping

Pham Viet Hung



M Û E G Y E T E M 1 7 8 2

# Content

- ▶ Resampling methods
  - ▶ Introduction
  - ▶ Jackknife vs Bootstrap
  - ▶ Bootstrap Confidence Interval
- ▶ Linear Regression
  - ▶ Gauss-Markov
  - ▶ Heteroscedasticity
- ▶ Bootstrap Regression with Simulations
  - ▶ Methods
  - ▶ Simulations

# Resampling

# Introduction

- ▶ Extract all the info from the data
- ▶ Earlier: a statistic is a RV having a probability distribution, that is, a sampling distribution of the statistic. Sometimes, no need to approximate the properties ("relative accuracy") of an estimator with some strong prior hypothesis. E.x.: OLS is a BLUE (Gauss-Markov Theorem).
- ▶ However, most of the time, those hypothesis do not apply. An estimation of a statistic is needed  $\implies$  Jackknife, Bootstrap

# Jackknife

---

**Algorithm 1:** The (Leave-One-Out) Jackknife

---

**Input** :  $S_n = (X_1, X_2, \dots, X_n)$  – observed samples of size  $n$

**Output:**  $J_{\theta_{jack}} = [\hat{\theta}_{(-1)}^*, \dots, \hat{\theta}_{(-n)}^*]$  – jackknifed parameters

**Init:**  $N :=$  sample size

$J_{\theta_{jack}} := []$  – an empty array of parameter estimates

$T$  – A function of a statistic

**for**  $i := 1$  to  $N$  **do**

$S_{(-i)}^* :=$  remove  $i^{\text{th}}$  element from  $S_n$

$\hat{\theta}_{(-i)}^* := T(S_{(-i)}^*)$

$J_{\theta_{jack}}[i] := \hat{\theta}_{(-i)}^*$

**end**

---

# Bootstrap

---

**Algorithm 2:** The Bootstrap

---

**Input** :  $\hat{F}_n = (X_1, X_2, \dots, X_n)$  – observed sample of size  $n$

**Output:**  $L_{\theta_{boot}} = [\hat{\theta}_1^*, \dots, \hat{\theta}_n^*]$  – bootstrapped parameters

**Init:**  $B$  – # of bootstrap repetitions

$b := 1$

$L_{\theta_{boot}} := []$  – an empty array of parameter estimates

$T$  – A function of a statistic

**while**  $b < B + 1$  **do**

$\hat{F}_n^{*(i)} :=$  a sample size  $n$  from  $S_n$  by **sampling with replacement**

$\hat{\theta}_b^* := T\left(\hat{F}_n^{*(b)}\right)$

$L_{\theta_{boot}}[b] := \hat{\theta}_b^*$

$b++$

**end**

---

## Bias-Corrected and accelerated (BCa)

- ▶ bootstrap CDF:  $\hat{G}(x) = \frac{\#\{\hat{\theta}_b^* \leq x\}}{B}$
- ▶  $\alpha^{\text{th}}$  percentile point  $\hat{\theta}^{*(\alpha)}$  of the bootstrap distribution:  $\hat{\theta}^{*(\alpha)} = \hat{G}^{-1}(\alpha)$
- ▶ bias-correction factor:  $\hat{z}_0 = \Phi^{-1}\left(\frac{\#\{\hat{\theta}_b^* < \hat{\theta}\}}{B}\right) \rightarrow$  to measure the difference between the median of  $\hat{\theta}^*$  and  $\hat{\theta}$  in normal units.
- ▶ acceleration factor:  $\hat{a} = \frac{1}{6} \frac{\sum_{i=1}^n (\hat{\theta}_{(\cdot)}^* - \hat{\theta}_{(-i)}^*)^3}{\left[ \sum_{i=1}^n (\hat{\theta}_{(\cdot)}^* - \hat{\theta}_{(-i)}^*)^2 \right]^{3/2}} \rightarrow$  corrects for skewness in the bootstrap sampling distribution using **jackknife**.

## Bias-Corrected and accelerated (BCa)

The  $100(1 - \alpha)\%$  **BCa** confidence interval is given by

$$\left[ \hat{\theta}^{*(\alpha_1)}, \hat{\theta}^{*(\alpha_2)} \right] = \left[ \hat{G}^{-1}(\alpha_1), \hat{G}^{-1}(\alpha_2) \right] \quad (1)$$

where

$$\begin{aligned} \alpha_1 &= \Phi \left( \hat{z}_0 + \frac{\hat{z}_0 + z_{\alpha/2}}{1 - \hat{\alpha}(\hat{z}_0 + z_{\alpha/2})} \right) \\ \alpha_2 &= \Phi \left( \hat{z}_0 + \frac{\hat{z}_0 + z_{1-\alpha/2}}{1 - \hat{\alpha}(\hat{z}_0 + z_{1-\alpha/2})} \right), \end{aligned} \quad (2)$$

and  $z_\alpha$  is the  $100\alpha$ th percentile point of the standard normal distribution.

**Speeding up BCa:** Inner Group Jackknife, Inner Random Subset Jackknife



# Linear Regression

## Model Specification

The multivariable regression model with  $n$  observations and  $d$  independent variables:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_d X_{id} + \epsilon_i \quad (3)$$

Model assumptions:

- ▶ OLSA.1: Observations are Mutually Independent
- ▶ OLSA.2: Linear Model DGP and Strict Exogeneity
- ▶ OLSA.3: Conditional Homoskedasticity
- ▶ OLSA.4: Conditionally Uncorrelated Error Terms
- ▶ OLSA.5: No Exact Collinearity Between Regressors
- ▶ OLSA.6: Conditionally Normally Distributed Errors

# Gauss-Markov Theorems

OLS Estimator:  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

**Gauss-Markov Theorem** states that If the conditions *OLSA.1 – OLSA.5* hold true, the OLS estimator  $\hat{\beta}$  is the **Best Linear Unbiased Estimator** (BLUE). It is also consistent with the true parameter values of the multivariate linear regression model.

**Modern Gauss-Markov Theorem** says that if  $\hat{\beta}$  is an unbiased estimator of  $\beta$  then  $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1} \leq \mathbb{V}_{\hat{\beta}}$

However, what about relaxing the OLSA.3 (conditional homoskedasticity) assumption?

## White Test for Heteroscedasticity

1. Estimate the model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \epsilon$  using OLS.
2. Obtain the predicted  $\hat{\mathbf{Y}}$  and the squared residual  $\hat{\epsilon}^2$ .
3. Run the OLS for  $\hat{\epsilon} \odot \hat{\epsilon} = \delta_0 + \delta_1 \hat{\mathbf{Y}} + \delta_2 \hat{\mathbf{Y}} \odot \hat{\mathbf{Y}} + u_i$ . Retrieve the R-squared value  $\mathbf{R}_{\hat{\epsilon}^2}^2$  from this regression.
4. For the following hypothesis test,

$$\begin{cases} H_0 : \delta_1 = \delta_2 = 0 \\ H_1 : \exists i \in \{1, 2\} : \delta_i \neq 0, \end{cases}$$

apply either the  $F = \frac{\mathbf{R}_{\hat{\epsilon}^2}^2 / 2}{(1 - \mathbf{R}_{\hat{\epsilon}^2}^2) / (n - 3)} \sim F_{2, n-3}$  or  $\text{LM} = n\mathbf{R}_{\hat{\epsilon}^2}^2 \sim \chi_2^2$ .

5. Calculate the  $p$ -value.

## Heteroskedasticity-Consistent Covariance Matrix Estimator

The error covariance matrix is  $\mathbb{V}_{\hat{\beta}} := \mathbb{V}(\hat{\beta} | \mathbf{X}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \Sigma \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$ .

The estimated HCCME is  $\hat{\mathbb{V}}_{\hat{\beta}}^{\text{HCCME}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$ .

Revolutionary idea by White:  $\hat{\mathbb{V}}_n \equiv \frac{1}{n} \sum_{i=1}^n \mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T \hat{\epsilon}_i^2 \xrightarrow{a.s.} \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T \epsilon_i^2)$ .

From which the asymptotic covariance matrix is  $(\mathbf{X}^T \mathbf{X} / n)^{-1} \hat{\mathbb{V}}_n (\mathbf{X}^T \mathbf{X} / n)^{-1}$ .

The finite version becomes the HC0:  $\hat{\mathbb{V}}_{\hat{\beta}}^{\text{HC0}} = (\mathbf{X}^T \mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T \hat{\epsilon}_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1}$ .

HC3:  $\hat{\mathbb{V}}_{\hat{\beta}}^{\text{HC3}} = (\mathbf{X}^T \mathbf{X})^{-1} \left( \sum_{i=1}^n (1 - h_{ii})^{-2} \mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T \hat{\epsilon}_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1}$ , where  $h_{ii}$  from  $\mathbf{H}$ .

# Bootstrap Regression

# Methods

Applying bootstrap principles on OLS. These techniques are nonparametric bootstrapping.

3 ways to do that:

- ▶ Pairs Bootstrap
- ▶ Residual Bootstrap
- ▶ Wild Bootstrap

However, from the modelling perspective, Wild Bootstrap is the desirable one.

## Wild Bootstrap

The wild bootstrap DGP:  $\mathbf{Y}_i^* = \mathbf{X}_i^T \hat{\boldsymbol{\beta}} + \xi_i^* f(\hat{\epsilon}_i)$  where  $\xi_i^*$  is a random variable.

Mammen's two-point distribution:

$$\xi_i^* = \begin{cases} -\frac{\sqrt{5}-1}{2} & \text{with probability } \frac{\sqrt{5}+1}{2\sqrt{5}} \approx 0.7236 \\ \frac{\sqrt{5}+1}{2} & \text{with probability } \frac{\sqrt{5}-1}{2\sqrt{5}} \approx 0.2764. \end{cases} \quad (4)$$

Rademacher distribution:

$$\xi_i^* = \begin{cases} -1 & \text{with probability } 1/2 \\ 1 & \text{with probability } 1/2 \end{cases} \quad (5)$$



# Wild Bootstrap

- ▶ When conditional uncorrelatedness is assumed, for large sample size it's not worth using wild bootstrap to estimate the standard errors of the coefficients. A simple OLS with HCCME is a surprisingly good estimator in this case. However, bootstrap confidence interval may be worth using even for larger sample sizes.
- ▶ When conditional uncorrelatedness (OLSA.4) is **NOT** assumed, even for larger sample sizes it is worth using wild bootstrap. This is called **wild cluster bootstrap**.

## Framework of the Simulations

- ▶ How bootstrap distributions and confidence intervals of regression coefficients change in wild bootstrap simulations when increasing the sample size and bootstrap replication numbers.
- ▶ CEU MicroData's Hungarian balance sheet data in the commerce sector from 2014.
- ▶ Python package [ols\\_bootstrap](#)
- ▶ The model:  $\ln(w_i/L_i) = \beta_0 + \beta_1 \ln(L_i) + \beta_2 D_i + \epsilon_i \quad \forall i = 1, \dots, n$  where
- ▶  $\beta_0 = 6.7833$  for the constant
- ▶  $\beta_1 = 0.2105$  for the log\_employment ( $\ln(L_i)$ )
- ▶  $\beta_2 = 0.3226$  for the is\_exporter ( $D_i$ )

# Confidence Interval Test

sample size	100	600	1100
BCa	0: 60	0: 46	0: 56
	1: 940	1: 954	1: 944
BC	0: 59	0: 44	0: 55
	1: 941	1: 956	1: 945
Percentile	0: 58	0: 44	0: 55
	1: 942	1: 956	1: 945
Reverse Percentile	0: 58	0: 45	0: 56
	1: 942	1: 955	1: 944

TABLE 4.1: The 95% Confidence Interval Test for  $\hat{\beta}_1$ 

sample size	100	600	1100
BCa	0: 67	0: 58	0: 46
	1: 933	1: 942	1: 954
BC	0: 59	0: 51	0: 42
	1: 941	1: 949	1: 958
Percentile	0: 60	0: 51	0: 46
	1: 940	1: 949	1: 954
Reverse Percentile	0: 57	0: 50	0: 47
	1: 943	1: 950	1: 953






TABLE 4.2: The 95% Confidence Interval Test for  $\hat{\beta}_2$

# Bootstrap Distribution Test

[https://github.com/pvh95/ols\\_bootstrap/blob/masterthesis2022/bca\\_application.ipynb](https://github.com/pvh95/ols_bootstrap/blob/masterthesis2022/bca_application.ipynb)

Thank you very much for your attention!





## Reference 1

-  Rao, C. R. (1989). *Statistics and Truth. Putting Chance to Work*. International Co-operative Publishing House, Burtonsville, Md.
-  Shao, J. & Tu, D. (1995). *The Jackknife and Bootstrap*. New York, NY, USA: Springer Verlag.
-  Tukey, J. W. (1958). *Bias and Confidence in not-quite Large Samples (Abstract)*. In: The Annals of Mathematical Statistics 29, p. 614
-  Efron, B. (1979). *Bootstrap Methods: Another Look at the Jackknife: Annual Statistics*. In.
-  Efron, B. & Tibshirani, R. (1993). *An Introduction to the Bootstrap*. New York: Chapman and Hall.

## Reference 2

-  Efron, B. & Balasubramanian N. (2020). *The Automatic Construction of Bootstrap Confidence Intervals*. In: Journal of Computational and Graphical Statistics 29.3, pp. 608–619.
-  Hansen, B. (2022b). *A Modern Gauss-Markov Theorem*. In: Econometrica 90.3, pp. 1283–1294
-  White, H. (1980). *A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity*. Econometrica: journal of the Econometric Society, 817-838.
-  Davidson, R. & Flachaire, E. (2008). *The Wild Bootstrap, Tamed at Last*. In: Journal of Econometrics 146.1, pp. 162–169.

## Reference 3

-  MacKinnon, J.G. and White, H., 1985. *Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties*. Journal of econometrics, 29(3), pp.305-325.
-  MacKinnon, J. G. (2013). *Thirty Years of Heteroskedasticity-Robust Inference*. In: Recent Advances and Future Firections in Causality, Prediction, and Specification Analysis. Springer, pp. 437–461.
-  Liu, R.Y. (1988). *Bootstrap procedures under some non-I.I.D. models*. Annals of Statistics, 16, 1696–1708.
-  Mammen, E. (1993). *Bootstrap and wild bootstrap for high dimensional linear models*. Annals of Statistics 21, 255–285.