# **Bootstrap Confidence Intervals**

ST552 Lecture 13

Charlotte Wickham 2019-02-11

The inferences we've covered so far relied on our assumption of Normal errors:

$$\epsilon \sim N(0, \sigma^2 I_{n \times n})$$

For example, we've seen under this assumption, the least squares estimates are also Normally distributed:

$$\hat{\boldsymbol{\beta}} \sim \boldsymbol{N} \left( \boldsymbol{\beta}, \, \sigma^2 \left( \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X} \right)^{-1} \right)$$

If the errors aren't truly Normally distributed, what distribution do the estimates have?

Imagine the errors are in fact  $t_3$  distributed?

With your neighbour: **design a simulation to understand the distribution of the least squares estimates.** 

a particular B

Some model ENtz 1) Assure some things about model : y= XB+ E Decide on X and B, make something up. (2) Compute \$\$'s 100,000 times a) Simulate y : simulate E frølig y= XB+E b) Using y, fit model:  $\beta = (xTx)^{-1} X^{T} Y$ (3) Make (3) Histogram Examine many  $\beta$ 's  $\beta_{0}$   $\beta_{0}$   $\beta_{0}$   $\beta_{0}$   $\beta_{0}$   $\beta_{0}$   $\beta_{0}$   $\beta_{1}$   $\beta_{2}$   $\beta_{3}$   $\beta_{2}$   $\beta_{3}$   $\beta_{3}$   $\beta_{2}$   $\beta_{3}$   $\beta_{3}$ Red: We da't know this with real data 4

# Example: 1. Fix n, fix X



# Example: 1. Fix $\beta$ , find $\hat{\chi} \in (\mathbf{u})$



# **Example: 2. Simulate errors, find** *y*



## Example: 3. Find least squares line



8

# Example: 4. Repeat #2. and #3. many times



9

# Example: Examine distribution of estimates



#### Example: Compared to theory



Think of our estimates like linear combinations of the errors. I.e. a sort of average of i.i.d random variables. Some version of the Central Limit Theorem will apply.

For large samples, even when the errors aren't Normal,

$$\hat{\beta} \sim N(\beta, \sigma^2(X^T X)^{-1})$$
  
approximately  
as  $h \rightarrow \infty$  that approximation  
will improve.

If we knew the error distribution and true parameters we could use simulation to understand the sampling distribution the least squares estimates.

Simulation can also be used to demonstrate the CLT at work in regression.

In practice, with data in front of us, we don't know the distribution of the errors (nor the true parameter values).

The bootstrap is one approach to estimate the sampling distribution of  $\hat{\beta}$ , by using the simulation idea, and substituting in our best guesses for the things we don't know.

(Model based resampling)

0. Fit model and find estimates,  $\hat{\beta}$ , and residuals,  $e_i$ 

1. Fix X. residual 2. For k = 1,..., B Repeat many times 2.1 Generate errors,  $\epsilon_i^*$  sampled with replacement from  $e_i$ 2.2 Construct y, using the model,  $y = \hat{y} + \epsilon^* = \times \hat{\beta} + \epsilon^*$ 2.3 Use least squares to find  $\hat{\beta}^*_{(k)}$ 3. Examine the distribution of  $\hat{\beta}^*$  and compare to  $\hat{\beta}$ 95% One confidence interval for  $\beta_j$  is the 2.5% and 97.5% quantiles of the distribution of  $\hat{\beta}_i^*$ . (Known as the Percentile method, there are other (better?) methods).

# Example: Faraway Galapagos Islands

(I'll illustrate with simple linear regression, Faraway does multiple case in 3.6)

#### Observed data



# **Bootstrap:** 1. Find $\hat{\beta}$ , $\hat{y}$ , and $e_i$ .



# Bootstrap: Using fixed X, beta from observed data

Observed data



# Bootstrap: 2. Resample residuals to construct bootstrapped response



# Bootstrap: 3. Fit regression model to bootstrapped response

Bootstrapped data



### Bootstrap: 3. Repeat #2. and #3. many times



## Examine distribution of estimates



# High level: bootstrap idea

We don't know the distribution of the errors, but our best guess is probably the empirical c.d.f on the residuals.

Sampling from a random variable with a c.d.f. defined as the empirical c.d.f. of the residuals, boils down to sampling with replacement from residuals.

dist .t E, 0/0/5

(esidual)

Experiment:

We might rely on bootstrap confidence intervals when we are worried about the assumption of Normal errors. But, there are limitations.

- We still rely on the assumption that the errors are independent and identically distributed.
- Generally scaled residuals are used (residuals don't have the same variance, more later)
- An alternative bootstrap resamples the  $(y_i, x_{i1}, \ldots, x_{ip})$ vectors, i.e. resamples the rows of the data, a.k.a resampling cases bootstrap. sesaple residuals Saphi study

We might rely on bootstrap confidence intervals when we are worried about the assumption of Normal errors. But, there are limitations.

- We still rely on the assumption that the errors are independent and identically distributed.
- Generally scaled residuals are used (residuals don't have the same variance, more later)
- An alternative bootstrap resamples the (y<sub>i</sub>, x<sub>i1</sub>, ..., x<sub>ip</sub>) vectors, i.e. resamples the rows of the data, a.k.a *resampling cases bootstrap*.