

# Resampling

## MC & Bootstrap

↳ Statistical model



unbiased estimator when  $E[\hat{\theta}] = \theta$

↳ so when bias  $b(\hat{\theta}) = E[\hat{\theta}] - \theta$

→ Data when CLT doesn't apply

## Importance Sampling

- Choosing an adequate sampling dist<sup>n</sup> is order to approximate the EV with good control on the estimation variance

MC @ 18-19 @3

$$y_i = \theta^* x_i + z_i$$

$$z_i \sim \chi^2_1$$

what could have a negative impact on the estimation of  $\theta^*$

- noise is heavy tailed & skewed positively

↳ could make observations larger than the true data.

↳ not an issue for slope but more outliers make skew the distribution

A data analyst is implementing a Monte Carlo simulation of  $M = 1,000$  random samples of realisations of the model

$$Y_i = \theta^* X_i + Z_i, \quad i = 1, \dots, n \quad (1)$$

with  $n = 100$ ,  $\theta^* = 8$  and a sequence of i.i.d. realizations  $Z_i \stackrel{iid}{\sim} t_d$  with  $d = 3$  degrees of freedom, using a single sample  $\{X_i\}_{i=1}^n$  from  $X \sim \mathcal{U}(1, 2)$  to generate all  $M$  Monte Carlo samples, and computes and stores the Monte Carlo least squares estimates of  $\theta^*$  for analysis. Note the analyst is making sure to not include an intercept in the regression model when fitting it to the simulated data.

- Quote the theoretic expected value of  $Y$ , i.e. the true value of  $E(Y)$ , showing your calculation.
- Quote the theoretic expected value of  $\hat{\theta}$ , i.e. the true value of  $E(\hat{\theta})$ , justifying your answer with a brief statement.

$$\begin{aligned} E[Y] &= E[\theta X_i] + E[Z_i] && Z_i \text{ \& } X_i \text{ independent!} \\ &= \theta E[X_i] + E[Z_i] \\ &= 8 \cdot (1.5) + 0 = 12 \end{aligned}$$

mean of a  $t$ -dist<sup>n</sup> = 0  $\forall df > 1$

$E(\theta)$

$\rightarrow 1000$  samples

$\rightarrow$  we find the mean of each sample where we have 100 observations.

$$\hat{\theta} = OLS = \frac{\sum X_i Y_i}{\sum X_i^2}$$

Ans

- The OLS being consistent under this model  $E[\theta^*] = 8$
- The MC simulation demonstrates this consistency.

c) How would MC estimate differ if  $Z \sim N(0, 1)$   
 $t$  dist<sup>n</sup> variance =  $\frac{df}{df-2}$

So good at modelling low population samples

- So if  $Z \sim N(0, 1)$ , the mean would be closer to the true value.

- There would be less variance in the model.

Actual Answer:

- mean should remain the same since the OLS is still unbiased
- The variance would decrease, since there is less outliers in our data

d) - Bias decrease  $E(\hat{\theta}_n) \xrightarrow[M \rightarrow \infty]{LLN} \theta$

- Higher proportion of outliers  
↳ more symmetric distribution as  $n$  increases.

Confidence Intervals for BS

- naive: just take quantiles of the estimated value
- Appropriate:  $2 \left( \text{true model parameter} \right)$  - quantile of the estimated value.

## MC estimation of standard normal CDF

$M \rightarrow$  sample size of each estimate.

$G \rightarrow$  number of grid pts where the CDF is evaluated

$x \rightarrow$  values where CDF is evaluated.

CDF  $\rightarrow$  store the value at each  $x$

$g \rightarrow$  func to be evaluated  $(e^{-\frac{x^2}{2}})$   
 $\hookrightarrow$  take the mean of it

$u \rightarrow$