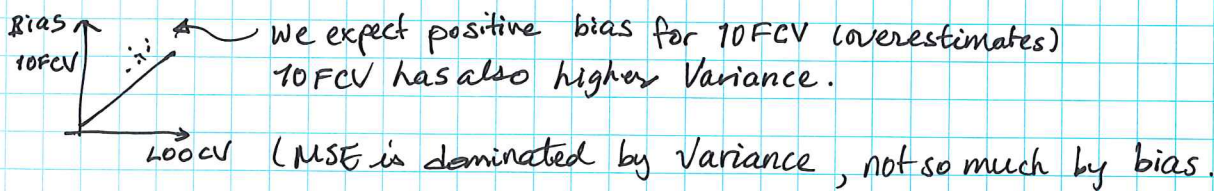


# this week Results of Cross Validation Bootstrap Monte Carlo tests (Panini)

- Recall  $Y = f(X) + \epsilon$  Goal: estimate  $\theta$  (expected test mSE) of an estimator  $\hat{f}$  of  $f$  using cross-validation  
ex.  $\hat{\theta}_{LOOCV}$ ,  $\hat{\theta}_{10FCV}$ . \* (fit on 90% of data, so it is more biased, overestimate)
- LOOCV: average the test error of each datum.  
→ more  $\hat{f}$  but less correlated
- 10FCV: get the test error over the (whole) dataset (ie. not individual data)  
→ fewer  $\hat{f}$  but more correlated
- Select / change  $X$ ,  $\epsilon$ , params

(Reviewing X Valid'n results of assignment survey - ref R code)

- ↪ assess / compare: test MSE vs. mean of 10FCV, LOOCV (results obtained via cv)
  - Variance of 10FCV, LOOCV



- Bootstrap - see lecture slides. Additional remarks:
  - $P^*$  is the conditional distribution based on
  - estimate bias using boot package in R
    - note that boot. is asymptotic
  - $\text{Var}(\hat{\theta}_n)$  and  $\text{Var}^*(\hat{\theta}_n^*)$  are similar for large sample sizes.

$$\hat{\theta}_n^* - \hat{\theta}_n \approx \hat{\theta}_n - \theta$$

- Consider the slightly weaker statement  $\frac{\hat{\theta}_n^* - \hat{\theta}_n}{\text{sd}(\hat{\theta}_n^*)} \approx \frac{\hat{\theta}_n - \theta}{\text{sd}(\hat{\theta}_n)}$   
how to obtain these?

- Given  $Z_1, \dots, Z_n$ ,

$$\begin{array}{lcl} \text{so } Z_1^{*1}, \dots, Z_n^{*1} & \rightarrow & \hat{\theta}_n^{*1} \\ Z_1^{*2}, \dots, Z_n^{*2} & \rightarrow & \hat{\theta}_n^{*2} \\ \vdots & & \vdots \\ Z_1^{*B}, \dots, Z_n^{*B} & \rightarrow & \hat{\theta}_n^{*B} \end{array} \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} \begin{array}{l} \\ \\ \\ \end{array} \begin{array}{l} \text{get the sample variance} = \hat{\text{sd}}(\hat{\theta}_n) \\ \\ B \text{ bootstraps} \end{array}$$

- How do we get  $\hat{\text{sd}}(\hat{\theta}_n^*)$ ?  
Do another level of bootstrap:  
ex. for  $\hat{\theta}_n^{*1}$ ,

$$\begin{array}{lcl} Z_1^{**}, \dots, Z_n^{**} & \rightarrow & \hat{\theta}_n^{**1} \\ Z_1^{**}, \dots, Z_n^{**} & \rightarrow & \hat{\theta}_n^{**2} \\ \vdots & & \vdots \\ Z_1^{**}, \dots, Z_n^{**} & \rightarrow & \hat{\theta}_n^{**M} \end{array} \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} \begin{array}{l} \text{Sample variance} = \hat{\text{sd}}(\hat{\theta}_n^{*1}) \\ \\ M \text{ bootstraps} \end{array}$$



- Bootstrap CIs: normal: 1<sup>st</sup> order asymptotically correct  
bootstrap T: 2<sup>nd</sup> order asymptotically correct  
→ good for location params but not for ratios.

- Is the ref. size ( $n$ ) enough? Assess via QQ plot: if it is roughly Gaussian distributed, then that's good ( $n$  is probably big enough) — can be confident about bootstrap results.  
If not, probably shouldn't trust.

- If the parametric model is bad, use nonparametric perhaps?  $\hat{P}_n$  is empirical.

$Z_1, \dots, Z_n \rightarrow \hat{\theta}$ . To get  $\hat{sd}(\hat{\theta}_n)$ , we generate  $B$  bootstrapped datasets and we take the sample variance of  $\hat{\theta}_n^{*1}, \dots, \hat{\theta}_n^{*B}$ .  
Next, we need the quantiles of  $\frac{\hat{\theta}_n^{*i} - \hat{\theta}_n}{\hat{sd}(\hat{\theta}_n^{*i})}$

$$\begin{aligned} Z_1^{*1}, \dots, Z_n^{*1} &\rightarrow \hat{\theta}_n^{*1} \\ &\vdots \\ Z_1^{*B}, \dots, Z_n^{*B} &\rightarrow \hat{\theta}_n^{*B} \end{aligned}$$

We approximate this by the empirical quantiles of  $\frac{\hat{\theta}_n^{*1} - \hat{\theta}_n}{\hat{sd}(\hat{\theta}_n^{*1})}, \dots, \frac{\hat{\theta}_n^{*B} - \hat{\theta}_n}{\hat{sd}(\hat{\theta}_n^{*B})}$

Hestenberg: check that bias, variance is consistent; else  $\hat{\theta}_n - \theta$  not so accurate?

$$\frac{\hat{\theta}_n - \theta}{\hat{sd}(\hat{\theta}_n)}$$

- ref R code (ex. airconditioning: aircond (from boot lib.)

- Parametric bootstrap: sample from param. data rather than from  $\hat{P}_n$

- `rexp(~)`: samples from exponential distribution using some params.  
arg `air.rg` function corresponds to params.

- Bootstrap for regression  $E[Y|X=x] = \mu(x)$

- Check that the bootstrapped samples  $n$  look like the data (sets)

- Why can't we fix the  $x$ 's and just resample the  $y$ 's?  
We lose any trends in the data



If, however, there are heteroskedastic errors, i.e. errors are not iid, then we lose trends in sampling the residuals.

