

# Lecture 10

## Model Selection

*STAT 8020 Statistical Methods II*  
September 11, 2019

Whitney Huang  
Clemson University

## 1 Variable Selection Criteria

- What is the appropriate subset size?
- What is the best model for a fixed size?

$$\begin{aligned}(\hat{Y}_i - \mu_i)^2 &= (\hat{Y}_i - E(\hat{Y}_i) + E(\hat{Y}_i) - \mu_i)^2 \\&= \underbrace{(\hat{Y}_i - E(\hat{Y}_i))^2}_{\text{Variance}} + \underbrace{(E(\hat{Y}_i) - \mu_i)^2}_{\text{Bias}^2},\end{aligned}$$

where  $\mu_i = E(Y_i|X_i = x_i)$

- Mean squared prediction error (MSPE):

$$\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^n (E(\hat{Y}_i) - \mu_i)^2$$

- $C_p$  criterion measure:

$$\begin{aligned}\Gamma_p &= \frac{\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^n (E(\hat{Y}_i) - \mu_i)^2}{\sigma^2} \\&= \frac{\sum \text{Var}_{\text{pred}} + \sum \text{Bias}^2}{\text{Var}_{\text{error}}}\end{aligned}$$

- Do not know  $\sigma^2$  nor numerator
- Use  $\text{MSE}_{X_1, \dots, X_{p-1}} = \text{MSE}_F$  as the estimate for  $\sigma$
- For numerator:
  - Can show  $\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 = p\sigma^2$
  - Can also show  $\sum_{i=1}^n (\text{E}(\hat{Y}_i) - \mu_i)^2 = \text{E}(\text{SSE}_F) - (n - p)\sigma^2$

$$\Rightarrow C_p = \frac{\text{SSE} - (n-p)\text{MSE}_F + p\text{MSE}_F}{\text{MSE}_F}$$

### Recall

$$\Gamma_p = \frac{\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^n (E(\hat{Y}_i) - \mu_i)^2}{\sigma^2}$$

- When model is correct  $E(C_p) \approx p$
- When plotting models against  $p$ 
  - Biased models will fall above  $C_p = p$
  - Unbiased models will fall around line  $C_p = p$
  - By definition:  $C_p$  for full model equals  $p$