# Lecture 10
## Model Selection

*STAT 8020 Statistical Methods II*
September 11, 2019

Whitney Huang
Clemson University

Model Selection

Variable Selection
Criteria

10.1

---

**Agenda**

Model Selection

Variable Selection
Criteria

**1** **Variable Selection Criteria**

10.2

---

**Variable Selection**

Model Selection

Variable Selection
Criteria

- What is the appropriate subset size?

- What is the best model for a fixed size?

10.3

## Mallows' $C_p$ Criterion

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

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

- Mean squared prediction error (MSPE):
  $\sum_{i=1}^{n} \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^{n} (\mathrm{E}(\hat{Y}_i) - \mu_i)^2$

- $C_p$ criterion measure:

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

Model Selection

CLEMS#N
UNIVERSITY

Variable Selection
Criteria

10.4

---

## $C_p$ Criterion

- Do not know $\sigma^2$ nor numerator

- Use $\mathsf{MSE}_{X_1,\cdots,X_{p-1}} = \mathsf{MSE}_\mathsf{F}$ as the estimate for $\sigma$

- For numerator:

  - Can show $\sum_{i=1}^{n} \sigma_{\hat{Y}_i}^2 = p\sigma^2$

  - Can also show
    $\sum_{i=1}^{n} (\mathrm{E}(\hat{Y}_i) - \mu_i)^2 = \mathrm{E}(\mathsf{SSE}_\mathsf{F}) - (n-p)\sigma^2$

  $\Rightarrow C_p = \frac{\mathsf{SSE} - (n-p)\mathsf{MSE}_\mathsf{F} + p\mathsf{MSE}_\mathsf{F}}{\mathsf{MSE}_\mathsf{F}}$

Model Selection

CLEMS#N
UNIVERSITY

Variable Selection
Criteria

10.5

---

## $C_p$ Criterion Cont'd

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

- When model is correct $\mathrm{E}(C_p) \approx p$

- 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 $p$

Model Selection

CLEMS#N
UNIVERSITY

Variable Selection
Criteria

10.6