**Best Subset Selection Methods**

**Cp statistic**

<u>Remarks</u>

1. If a model is underfitted, $\hat{\sigma}^2$, $\hat{\beta}$ and $\hat{y}$ at a new observation of $x_0$ are not unbiased estimators.

2. If a model is overfitted, $\hat{\sigma}^2$, $\hat{\beta}$ and $\hat{y}$ at a new observation of $x_0$ are still unbiased estimators but with larger variance.

3. Normally, we choose a larger value for $\alpha_{IN}$ than $\alpha_{OUT}$ to avoid underfitting a model.

4. The Lack of Fit test in Chapter 1 is a test for testing whether the model is underfitted if there are repeated measurements of $y$ for the same $x$.

Consider the mean square error of $\hat{y}(x_i)$

$$\sum_{i=1}^{n} \frac{\text{MSE}(\hat{y}(x_i))}{\sigma^2} = p' + \frac{\text{E}(\hat{\sigma}_{p'}^2) - \sigma^2}{\sigma^2}(n - p')$$

Then, Cp is defined as its estimate, i.e.,

$$Cp = 2p' - n + \frac{ResS.S._{p'}}{\hat{\sigma}_{\text{full model}}^2}$$

<u>Remarks</u>

1. For each $p'$, calculate $C_p$ for the model with smallest $ResS.S._{p'}$.

2. Find the model with the smallest $C_p$.

3. Find a model with the smallest mean square error on predicted values of $y$.