# Model Building with Variable Selection – Part I

## Johns Hopkins Engineering

## 625.461 Statistical Models and Regression

Module 9 – Lecture 9B

JOHNS HOPKINS

WHITING SCHOOL
*of* ENGINEERING

# Basics behind Variable Selection

$y$:  response variable

$x_1, \ldots, x_K$:  regressors

$$n \geq K + 1$$

**Assume: intercept is always in the model**

Fit the model containing all $K$ regressors

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$   $\mathbf{X}$:  $n \times (K+1)$

Suppose that we delete $r$ regressors and retain $p = K - r + 1$ regressors

$$\mathbf{y} = \mathbf{X}_p \boldsymbol{\beta}_p + \mathbf{X}_r \boldsymbol{\beta}_r + \boldsymbol{\varepsilon}$$

For the full model, the LS estimator of $\beta$

$$\hat{\boldsymbol{\beta}}^* = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

An estimator of the residual variance

$$\hat{\sigma}^2_* = \frac{\mathbf{y'y} - \hat{\boldsymbol{\beta}}*'\mathbf{X'y}}{n - K - 1} = \frac{\mathbf{y'}\left[\mathbf{I} - \mathbf{X}(\mathbf{X'X})^{-1}\mathbf{X'}\right]\mathbf{y}}{n - K - 1}$$

The components of $\hat{\boldsymbol{\beta}}*$ are $\hat{\boldsymbol{\beta}}_p^*$ and $\hat{\boldsymbol{\beta}}_r^*$ .

For the subset model containing *K-r* regressors,

$$\mathbf{y} = \mathbf{X}_p\boldsymbol{\beta}_p + \boldsymbol{\varepsilon}$$

$$\hat{\boldsymbol{\beta}}_p = \left(\mathbf{X}_p'\mathbf{X}_p\right)^{-1}\mathbf{X}_p'\mathbf{y}$$

$$\hat{\sigma}^2 = \frac{\mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}_p'\mathbf{X}_p'\mathbf{y}}{n-p} = \frac{\mathbf{y}'\left[\mathbf{I} - \mathbf{X}_p\left(\mathbf{X}_p'\mathbf{X}_p\right)^{-1}\mathbf{X}_p'\right]\mathbf{y}}{n-p}$$

# Basics behind Variable Selection

$$E\left(\hat{\boldsymbol{\beta}}_p\right) = \boldsymbol{\beta}_p + \left(\mathbf{X}'_p\mathbf{X}_p\right)^{-1}\mathbf{X}'_p\mathbf{X}_r\boldsymbol{\beta}_r = \hat{\boldsymbol{\beta}}_p + \mathbf{A}\boldsymbol{\beta}_r$$

Alias matrix: $\quad \mathbf{A} = (\mathbf{X}'_p\mathbf{X}_p)^{-1}\mathbf{X}'_p\mathbf{X}_r$

$\hat{\boldsymbol{\beta}}_p$ is biased for $\boldsymbol{\beta}_p$ unless $\mathbf{X}'_p\mathbf{X}_r = 0$

$$\mathrm{Var}\left(\hat{\boldsymbol{\beta}}_p\right) = \sigma^2\left(\mathbf{X}'_p\mathbf{X}_p\right)^{-1} \qquad \mathrm{Var}\left(\hat{\boldsymbol{\beta}}*\right) = \sigma^2\left(\mathbf{X}'\mathbf{X}\right)^{-1}$$

$\mathrm{Var}\left(\hat{\boldsymbol{\beta}}_p^*\right) - \mathrm{Var}\left(\hat{\boldsymbol{\beta}}_p\right)$ is positive definite

# Basics behind Variable Selection

$$E(\hat{\sigma}^2) = \sigma^2 + \frac{\beta_r' X_r' \left[ I - X_p (X_p' X_p)^{-1} X_p' \right] X_r \beta_r}{n - p}$$

The subset-model estimator of $\sigma^2$ is biased upward for $\sigma^2$

# Basics behind Variable Selection

Suppose that we wish to predict the response at $\mathbf{x}' = \left[ \mathbf{x}'_p, \mathbf{x}'_r \right]$.

If we use the full model, the predicted value is $\hat{y}* = \mathbf{x}'\hat{\boldsymbol{\beta}}*$, with mean $\mathbf{x}'\boldsymbol{\beta}$ and predicted variance

$$\text{Var}(\hat{y}*) = \sigma^2 \left[ 1 + \mathbf{x}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x} \right]$$

If the subset model is used, $\quad \hat{y} = \mathbf{x}_p' \hat{\boldsymbol{\beta}}_p$

$$E(\hat{y}) = \mathbf{x}_p' \boldsymbol{\beta}_p + \mathbf{x}_p' \mathbf{A} \boldsymbol{\beta}_r$$

$$\text{MSE}(\hat{y}) = \sigma^2 \left[ 1 + \mathbf{x}_p' (\mathbf{X}_p' \mathbf{X}_p)^{-1} \mathbf{x}_p \right] + \left( \mathbf{x}_p' \mathbf{A} \boldsymbol{\beta}_r - \mathbf{x}_r' \boldsymbol{\beta}_r \right)^2$$

$$Var(\hat{y}*) \geq Var(\hat{y})$$

JOHNS HOPKINS

WHITING SCHOOL
*of* ENGINEERING