

## Model Building / Variable Selection

### • PRIOR INFORMATION

#### • Data-driven techniques

{ Forward  
Backward Algorithms  
FB Stepwise

If there are  $K$  possible predictors:  $K$  steps

All possible models:  $2^K$

### Criterion-driven

AIC Akaike information Criterion

BIC Bayesian information Criterion

Let  $M$  be the model:

$$Y_i = \mu(x_{1i}, x_{2i}, \dots, x_{pi}) + \varepsilon_i$$

$$\mu(x_i) = \beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi}$$

where  $\varepsilon_i \text{ iid } N(0, \sigma^2)$

The maximum likelihood estimator of  $\sigma^2$  computed from the data  $\{(x_i, y_i), i = 1, \dots, n\}$  is

$$\hat{\sigma}_{MLE}^2 = \frac{\|\underline{R}\|^2}{n} \quad \text{where} \quad \underline{R} = \begin{pmatrix} R_1 \\ \vdots \\ R_n \end{pmatrix}, \quad R_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip})$$

$\uparrow$   
 $n - (p+1)$

the AIC & BIC values for this model are, respectively:

$$AIC(x_1 \dots x_p) = \log \hat{\sigma}_{MLE}^2(p) + \frac{n + (p+1)}{n - (p+1) - 2}$$

$$BIC(x_1 \dots x_p) = \log \hat{\sigma}_{MLE}^2(p) + \frac{(p+1) \log n}{n}$$

	$\underbrace{\hspace{10em}}$ FIT	$\underbrace{\hspace{10em}}$ PENALTY
As $p \uparrow$	$\downarrow$ <hr style="width: 50%; margin: 0 auto;"/>	$\uparrow$ <hr style="width: 50%; margin: 0 auto;"/>

Squared Prediction error criterion

$$y_i = \mu(x_i) + \varepsilon_i$$

$$\mu(x_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

$$\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

# DATA

Cross Validation

Training

$$\{(x_1, y_1), (x_3, y_3), \dots\}$$

Fit the model

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi} + \varepsilon_i$$

$$i = 1, 3, 5, \dots$$

$$\Rightarrow \hat{\beta}_{\text{TRAINING}} = \begin{pmatrix} \hat{\beta}_{0, \text{TR}} \\ \vdots \\ \hat{\beta}_{p, \text{TR}} \end{pmatrix}$$

Testing

$$\{(x_2, y_2), (x_4, y_4), \dots\}$$

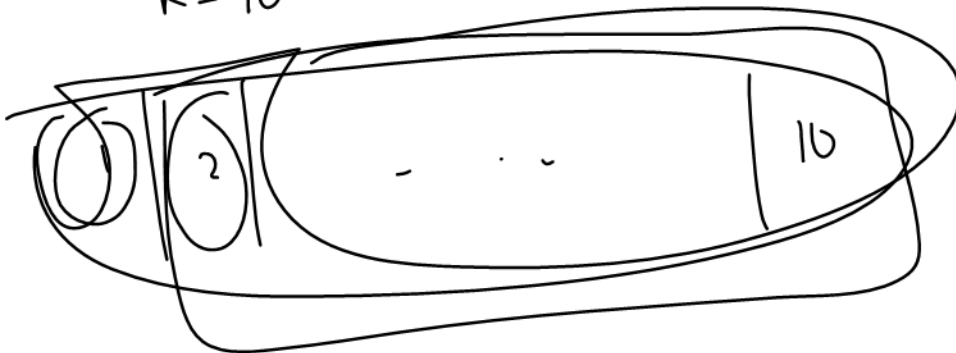
$$PE(i) = y_i - (\hat{\beta}_{0, \text{TR}} + \hat{\beta}_{1, \text{TR}} x_{1i} + \dots + \hat{\beta}_{p, \text{TR}} x_{pi})$$

$$i = 2, 4, 6, \dots$$

$$SPE = \sum_{i=2,4,\dots} (PE(i))^2$$

"k-fold cross-validation"

k = 10



$$\text{DATA} = \bigcup_{g=1}^{10} \text{Data}_g$$

For  $g=1:10$

$$\text{Test}_g = \text{Data}_g$$

$$\text{Train}_g = \text{Data} \setminus \text{Data}_g$$

$$\text{SPE}_g$$

END

$$\text{SPE} = \sum_{g=1}^{10} \text{SPE}_g$$

Multi collinearity  $\rightarrow$

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \varepsilon_i$$

$$Y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, I \otimes \sigma^2)$$

$$X = \begin{pmatrix} 1 & X_{11} & X_{21} \\ \vdots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} \end{pmatrix}_{n \times 3}$$

$$\hat{\beta} = (X'X)^{-1}X'Y$$

$$\text{cov}(\hat{\beta}) = (X'X)^{-1} \otimes \sigma^2 = V_{\hat{\beta}} \quad (3 \times 3)$$

$$\text{Var}(\hat{\beta}_1) = V_{\hat{\beta}}(2,2)$$

$$= \left( \frac{1}{1 - r_{12}^2} \right) S_{X_1 X_2} \sigma^2$$

$$\text{Where } S_{X_1 X_2} = \sum_{i=1}^n (X_{1i} - \bar{X}_1)(X_{2i} - \bar{X}_2)$$

$$r_{12} = \frac{\sum (X_{1i} - \bar{X}_1)(X_{2i} - \bar{X}_2)}{\sqrt{\sum (X_{1i} - \bar{X}_1)^2 \cdot \sum (X_{2i} - \bar{X}_2)^2}}$$



