

# Model selection [F.3.4]

L12

24.02.2017

SPSE = Expected squared prediction error

$E(Y) = \mu$  is the truth, but we model  $\mu$  as  $\mu = X\beta_M$  and we assume

$$Y = X\beta_M + \varepsilon, \quad E(\varepsilon) = 0, \quad \text{Var}(\varepsilon) = \sigma^2 I$$

$M \subseteq \{0, 1, 2, \dots, k\}$  and  $|M| = \text{size of model}$   
 = #parameters  
 based on a subset of all the available covariates

TRAINING:  $i = 1, \dots, n$  our observations, available  $y_i, x_i^T$   
 and  $\hat{\beta}_M$  is the estimator from the training set  
 for our model  $M$ .

VALIDATION:  $j = 1, \dots, T$  new observations, available as  
 $y_j$  and  $x_j^T$ .

$$\sum_{j=1}^T E \left( (y_j - \hat{y}_{jM})^2 \right) = \text{SPSE}$$

$y_j$  → new obs       $\hat{y}_{jM}$  → predicted value based on  $\hat{\beta}_M$   
 and  $x_j^T$

$$\text{Expected value of } \hat{Y}_j = E(\hat{Y}_j) = \sum_{m=1}^M (\mu_{jm} - \mu_j)^2$$

Bias - Variance trade-off

$\text{Var}(\epsilon) = \sigma^2 I$

$= \dots =$

$n\sigma^2 + |M| \cdot \sigma^2 +$

Variance of model

||

Irreducible prediction error

Can be made smaller by choosing fewer variables

Verance of errors

Verance of model

True value of  $y_j$

Expected value of  $\hat{y}_j$

Want to minimize SPSE, but difficult since  $\sigma^2$  and  $\mu_j$  is unknown.

PLAN: estimate SPSE and choose the  $\gamma$  that minimizes this estimate.

Finding the best model ← all subsets method

↓  
smallest SSE

1) Have  $k$  covariates that might be used

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \varepsilon_i$$

How many possible models can we make  
(want to have intercept)?

$$2 \cdot 2 \cdot 2 \cdots 2 = 2^k \text{ possible models}$$

fit all possible  $2^k$  models.

2) For  $M = \{0, 1, 2, \dots, k\}$  choose the model with the smallest SSE.

Ex: Acid rain:  $k=7 \Rightarrow$  total  $2^7 = 128$  possible models

complexity $ M  = 0$	searched $\frac{1}{7}$	best model $x_4$ (Al)
$ M  = 1$	$\binom{7}{1} = 21$	$x_1, x_3$
$ M  = 2$	$\binom{7}{2} = 35$	$x_1, x_2, x_3$
$ M  = 3$	$\binom{7}{3} = 35$	$x_1, x_2, x_3, x_5$
$ M  = 4$	$\binom{7}{4} = 35$	$x_1, x_2, x_3, x_5, x_7$
$ M  = 5$	$\binom{7}{5} = 21$	$x_1, x_2, x_3, x_5, x_7, x_8$
$ M  = 6$	$\binom{7}{6} = 7$	$x_1, x_2, x_3, x_4, x_5, x_7$
$ M  = 7$	$\binom{7}{7} = 1$	$x_1, x_2, x_3, x_4, x_5, x_6, x_7$

3) Now we need to choose between these k+1 models found in 2). Which criterion should I use?

$$i) R^2_{\text{adj}} = 1 - \frac{\frac{SSE}{(n-p)}}{\frac{SST}{n-1}} \quad |M|=p=k+1$$

Ex: Happiness:

$|M|=1$ : only intercept

$|M|=2$ : love (3) : 60.5

$|M|=3$ : love+work (10) : 66.3

$|M|=4$ : love+work+money (13) : 68.4

$|M|=5$ : love+work+sex+money (15) : 67.6

$$ii) \text{ Mallows' } C_p = \frac{\hat{\sigma}_{\text{full}}^2}{\hat{\sigma}_{M}^2} - n + 2|M|$$

$$\text{vs } S^2 \hat{\sigma}_{\text{full}}^2 = SSE + 2|M| \cdot \hat{\sigma}_{\text{full}}^2 \leftarrow \begin{matrix} \uparrow \\ \text{give same result.} \end{matrix}$$

$$iii) AIC = n \cdot \ln(\hat{\sigma}^2) + 2(|M| + 1)$$

$$\frac{SSE}{n-p+k+1}$$

$$iv) BIC = n \cdot \ln(\hat{\sigma}^2) + \ln(n) (|M| + 1)$$

BIC gives more penalty than AIC to large models.

Ex: Happy: BIC best model = love + work

Homework: slide 24 Acid rain

Transformation of response and predictors might improve the fit of the regression model.

The BoxCox transform

$$g_\lambda(\lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \ln(y) & \lambda = 0 \end{cases}$$

Class of function

For  $Y = X\beta + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2 I)$  the best value of  $\lambda$  is based on maximizing the likelihood profile

$$l(\lambda) = -\frac{n}{2} \ln \left( \frac{\text{SSE}_\lambda}{n} \right) - (\lambda - 1) \sum_{i=1}^n \ln \hat{e}_i$$

$\text{SSE}_\lambda$  is the SSE when  $g_\lambda(Y)$  is the response

R: boxcox(fit), see plot.