

21.02.2017

L11

Hypothesis testing about β_j

Previously: $Y = \bar{X}\beta + \epsilon, \epsilon \sim N_n(0, \sigma^2 I)$

$\hat{\beta}_j$ is the jth element of $\hat{\beta} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T Y$

$$\hat{\beta} \sim N_p(\beta, (\bar{X}^T \bar{X})^{-1} \sigma^2)$$

$$\hat{\sigma}^2 = \frac{1}{n-p} \hat{\epsilon}^T \hat{\epsilon} = \frac{SSE}{n-p}$$

$$\hat{\epsilon} = Y - \hat{Y} = Y - \bar{X}\hat{\beta}$$

$$T_j = \frac{\hat{\beta}_j - \beta_j}{\sqrt{c_{jj}} \hat{\sigma}} \sim t_{n-p}$$

/ ↗
obs # param. we estimate

jth diagonal element of $(\bar{X}^T \bar{X})^{-1}$

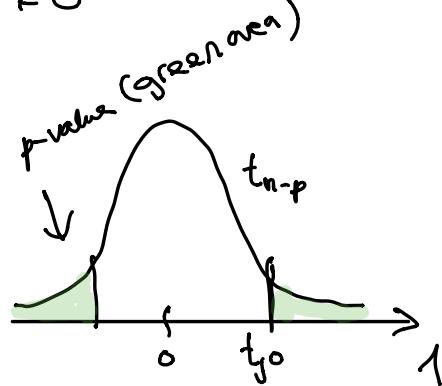
Test for association (linear) between response Y and X_j :

$$H_0: \beta_j = 0 \quad \text{vs} \quad H_1: \beta_j \neq 0$$

When H_0 is true:

$$T_{j0} = \frac{\hat{\beta}_j - 0}{\sqrt{c_{jj}} \hat{\sigma}} \sim t_{n-p}$$

test statistic



$$\begin{aligned}
 p\text{-value: } & P(|T_{j_0}| \geq |t_{j_0}|, H_0 \text{ true}) \\
 & = 2 \cdot P(T_{j_0} \geq |t_{j_0}|) \quad t_{j_0} \\
 & \quad \uparrow \quad t_{j_0} \\
 & \quad n-p
 \end{aligned}$$

Reject H_0 when $|t_{j_0}| > t_{\frac{\alpha}{2}, n-p}$ sign.level α .

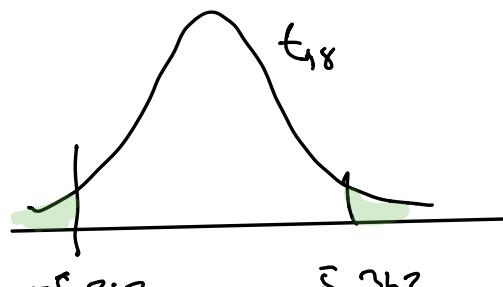
Ex: Acid rain : linear association between SO₄ and pH

$$H_0: \beta_1 = 0 \Rightarrow H_1: \beta_1 \neq 0$$

From summary of lm in R (slide)

$$t_{10} = -5.362$$

$$n-p = 18$$



$$p\text{-value} =$$

$$2 \cdot P(T_{18} > 5.362) \underset{\uparrow}{=} 4.3 \cdot 10^{-5}$$

$$R: 2 * (1 - pt(5.362, 18))$$

\uparrow
 area to left

Reject H_0 for all
 $\alpha > 4.3 \cdot 10^{-5}$.

Residuals (again)

$$\hat{\varepsilon} = y - \hat{y}, \quad \hat{\varepsilon} \sim N_n(0, \sigma^2(I - H))$$

R: residuals(fit)

The residuals have heteroscedastic variances

$$\text{Var}(\hat{\varepsilon}_i) = \sigma^2(1 - h_{ii}) \text{ and } \text{Cov}(\hat{\varepsilon}_i, \hat{\varepsilon}_j) = \sigma^2(0 - h_{ij})$$

can in general be
 $\neq 0$, but in most
 cases experience shows that
 ≈ 0

Standardized residuals:

$$r_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma} \sqrt{1 - h_{ii}}} \quad \text{will be (approx.) homoscedastic.}$$

R: rstandard(fit)

Studentized residuals: fitting the model to all obs. except i to make r_i^* . ← see slide

↓
use studentized!

R: rstudent(fit)

see example on slide
 for r_i vs $\hat{\varepsilon}_i$

Analysis of variance decomposition and R^2

$$y_1, \dots, y_n \text{ and } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 = \dots = \\ &\quad \text{sums of squared total} \\ &= \sum_{i=1}^n \underbrace{(y_i - \hat{y}_i)^2}_{\text{SSE}} + \sum_{i=1}^n \underbrace{(\hat{y}_i - \bar{y})^2}_{\text{SSR}} \\ &\quad \text{sums of squared error} \qquad \text{sums of sq regression} \\ &\quad \text{SSE} \qquad \text{SSR} \\ &\quad \uparrow \\ &\quad \text{explained by regression} \end{aligned}$$

With vectors and matrices.

$$\begin{array}{c} Y^T(I - \frac{1}{n}11^T)Y \\ \text{SST} \end{array} = \begin{array}{c} Y^T(I - H)Y \\ \text{SSE} \end{array} + \begin{array}{c} Y^T(H - \frac{1}{n}11^T)Y \\ \text{SSR} \end{array}$$

This is used to define:

$$R^2 = \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}}$$

↗ relative proportion of
total variability explained
by the regression

↗ coefficient of determination

Ex: Acid rain, full model (all covariates available)
 $R^2 = 0.73, 73\%$

Is the regression significant?

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0 \quad \text{vs}$$

$H_1:$ at least one $\beta_j \neq 0 \quad j=1, \dots, k$

Test statistic: $F = \frac{\text{SSR}/k}{\text{SSE}/(n-p)} \sim F_{k, n-p}$

\uparrow
 prove this in Part 3 in
 a general setting

Ex: Acid rain:

$$\left. \begin{array}{l} F\text{-observed: } 34.15 \\ k=7, n-p=18 \end{array} \right\} \underbrace{P(F_{7,18} > 34.15)}_{p\text{-value}} = 3.9 \cdot 10^{-7}$$

Ex: Volume of tree and the lumberjack

Big model: $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \varepsilon_i$

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

$R^2_{\text{Big}} \geq R^2_{\text{Small}}$ since $\hat{\beta}_3$ is found

to minimize $SSE = \hat{\varepsilon}^T \hat{\varepsilon}$, thus maximize

$$R^2 = 1 - \frac{SSE}{SST}$$

$R^2_{\text{Big}} = R^2_{\text{Small}}$ if $\hat{\beta}_3 = 0$

if $\hat{\beta}_3 \neq 0$ then $SSE_{\text{Big}} < SSE_{\text{Small}}$ and

$R^2_{\text{Big}} > R^2_{\text{Small}}$.

R^2 will always increase (or stay unchanged) when a new covariate is added to the model.

Next: more on choosing a good model, and then

$$R^2_{\text{Adj}} = 1 - \frac{SSE/(n-p)}{SST/(n-1)} \leftarrow \begin{array}{l} \text{penalizing} \\ \text{adding} \\ \text{many} \\ \text{covariates} \end{array}$$

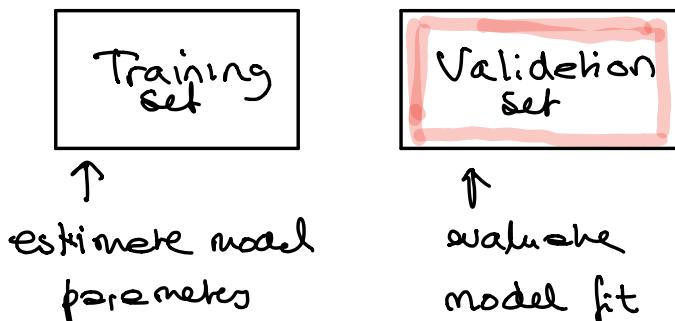
is one criterion to use instead of R^2 for model selection.

Model choice and variable selection [F3.4]

Question 1: Is a full model (all available covariates fitted) the best model?

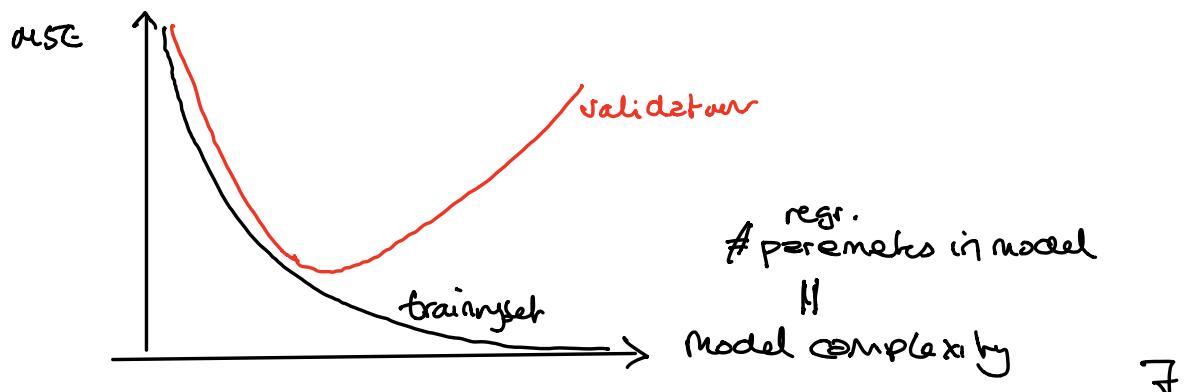
good interpretability good for future predictions

Data is divided into



$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Now calculate the MSE on the training and on the validation set and plot as a function of model complexity:



Answer 1: No, this may lead to overfitting =
fitting the trend + the noise!

⇒ so, what can we do instead?