

Model Adequacy checking

Wednesday, 21 August 2019 11:06 AM

$$E(\mathbf{y}) = \mathbf{X}\beta \quad \beta \text{ is model parameters}$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I}) \quad \sigma^2 \text{ unknown}$$

- 1) Errors are uncorrelated
- 2) Errors are independent.
- 3) Errors are normally distributed

Coefficient of determination (R^2)

$$R^2 = \frac{\text{Variation of } y \text{ explained by model}}{\text{Total Variation}}$$

$$= \frac{SS_{\text{Model}}}{SS_T}$$

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

⊗ Larger value of R^2 is considered better for the model

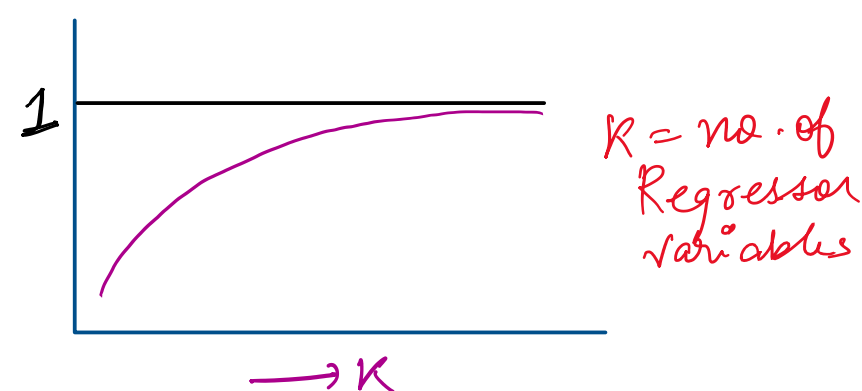
⊗ $R^2 \in [0, 1]$

$$SS_T = SS_{\text{Model}} + SS_{\text{Error}}$$

$$1 = \frac{SS_{\text{Model}}}{SS_T} + \frac{SS_{\text{Error}}}{SS_T} \Rightarrow R^2 = 1 - \frac{SSE}{SST}$$

$$\text{As } R^2 = 1 - (SSE/SS_T)$$

If we increase the number of regression variables such that $(\mathbf{X}^T \mathbf{X})$ remains invertible then R^2 will increase with the number of regressor variable.



Adjusted R^2

$$R_{\text{adj}}^2 = 1 - \frac{SSE/df(SSE)}{SST/df(SST)}$$

$$= 1 - \frac{SSE/(n-k-1)}{SST/(n-1)}$$

when R_{adj}^2 attains maxima then we can stop and build the model.

ERROR ANALYSIS

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{P}_X) \mathbf{y} \sim \mathcal{N}(\mathbf{0}, \sigma^2 (\mathbf{I} - \mathbf{P}_X))$$

$$\frac{\mathbf{e}^T \mathbf{e}}{\sigma^2} \sim \chi_{n-k-1}^2$$

$$e_i = y_i - \hat{y}_i \sim \mathcal{N}(0, \sigma^2 (1 - h_{ii}))$$

Assuming $\mathbf{H} = \mathbf{P}_X$

$$\text{cov}(e_i, e_j) = \text{cov}(y_i - \hat{y}_i, y_j - \hat{y}_j)$$

$$= \sigma^2 (-h_{ij}) = \sigma^2 (\mathbf{I} - \mathbf{H})_{ij}$$

⊗ estimated errors are correlated in general

Define

$$r_i = \frac{e_i}{\sqrt{\sigma^2 (1 - h_{ii})}}$$

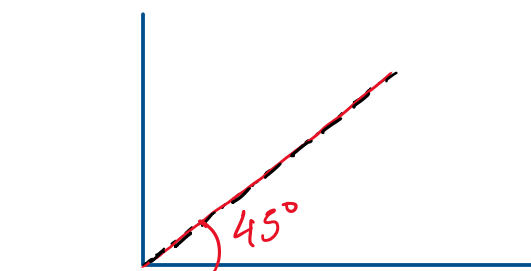
We can estimate

$$\hat{\sigma}^2 = \frac{SSE}{n-k-1}$$

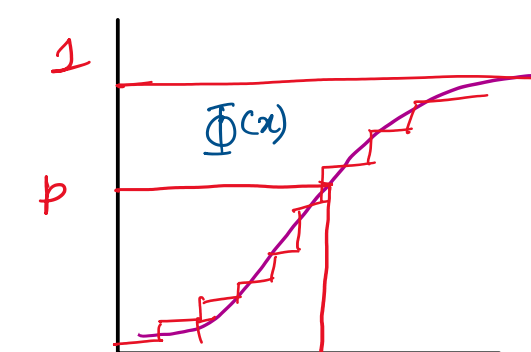
r_i will follow $\mathcal{N}(0, 1)$ when $n \uparrow \infty$
because $\hat{\sigma}^2 \rightarrow \sigma^2$ with Prob 1 (Chebyshev)

If we plot the histogram of empirical density then we are supposed to get the pdf of $\mathcal{N}(0, 1)$.

We also can do $q-q$ plot for r_i 's. If the errors are normally distributed then only qq plot will give a straight line passing through origin with slope 1.



$q = x$ such that



$$P(Z \leq a) = \Phi$$

$$\frac{\#(r_i \leq x)}{n} = \Phi$$

$$\frac{\sum \mathbb{1}_{\{r_i \leq x\}}}{n} = \Phi$$

$p-p$ plot

$e_i = y_i - \hat{y}_i$ prediction error of y_i but the same y_i has been used to get $\hat{\beta}$ and \hat{y}_i .

⊗ If we remove (y_i, x_i) from the dataset and predict y_i based on other $(n-1)$ data then what will be the change corresponding

$$\frac{e_i}{\sqrt{v(e_i)}}$$

where $e_i = y_i - \hat{y}_i$

$\hat{y}_i =$ predicted value of y_i based on rest $(n-1)$ observations.

$$\sum_{i=1}^n e_i^2 = \text{Predicted error sum of squares (Residual)}$$

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

Leave out/Jack-knife

$$e_i = \frac{e_i}{1 - h_{ii}} \sim \mathcal{N}(0, \frac{\sigma^2}{1 - h_{ii}})$$

This can be shown

Standardize e_i as

$$\frac{e_i}{\sqrt{\sigma^2 (1 - h_{ii})}} = \frac{e_i / (1 - h_{ii})}{\sqrt{\sigma^2 / (1 - h_{ii})}}$$

$$= \frac{e_i}{\sqrt{\sigma^2 (1 - h_{ii})}}$$

↳ $(n-1)$ observations

Here we need to estimate σ^2 based on $(n-1)$ observations

①

②

③