

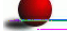



Building the Regression Model

C. Andy Tsao

Department of Applied Math, National Dong Hwa University

Outline

-  Overview
-  Selection of Predictors
-  Diagnosis
-  Remedial Measure
-  Validation


Learning from Data

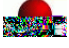


Objectives

- ▄▄▄▄ Reduction of explanatory or predictor variables
Find parsimonious model with good explanatory prediction power. Trade-off.
- Model refinement and selection
Choosing from many "good" models, checking the adequacy of the models, sensitivity of the models, fixing the weak spots.
- ▄▄▄ Model Validation
Ready to explain what's going on? Ready to predict what the future will be?
- ▄▄▄ Trade-off
Best explanatory prediction power v. Parsimony
Criteria and how to use them? "Good" models

Selection-I

 $R_p^2 = 1 - \frac{SSE_p}{SSTO}$. ID those with substantial increases.
NOT the biggest one.

 $R_a^2 = 1 - ($

Comments

- 🌐 No easy, clear-cut way to ID the best model
- Usually, many "good" models rather than one best model
- Respect the hierarchy of models
 - 🌐 Higher order terms < lower order terms
($X^4 < X^1$)
 - 🌐 Interaction terms < main effect terms
($X_1 X_2 < X_1$ or X_2)
- Chapter 10 Variable Selection of Faraway, J. (2002).
Also his Chapter 11 is highly recommended

Diagnosis

Checking the adequacy of a regression model

- 🇺🇸 Improper functional form of a predictor

- 🇺🇸 Outliers

 - Influential observation

- 🇺🇸 Multicollinearity

Improper functional form of a predictor

- Goal: Detect the suitable form of



Outliers-I

🏠 Studentized Deleted Residual:

$$t_i = \frac{d_i}{s(d_i)} \text{ where } s(d_i) = \frac{\text{MSE}_{(i)}(1 - h_{ii})}{n}$$


■ Hat matrix Leverage values → Outlying \mathbf{X}

🟢 $0 \leq h_{ii} \leq 1, \quad \sum_{i=1}^n h_{ii} = p.$

🏠 $\bar{h}_{ii} = \frac{p}{n}. \quad 2p/n, \text{ extreme } h_{ii}, \text{ outside } (0.2, 0.5)$

🟢 $\mathbf{h}_{new} = \mathbf{X}_{new}(\mathbf{X} \mathbf{X})^{-1} \mathbf{X}_{new}$ for hidden extrapolation.

Influential obs


(DFFITS)_i = $\frac{\hat{Y}_i - \hat{Y}_{i(i)}}{\sqrt{\text{MSE}_{(i)} h_{ii}}}$ **Flag:** If |DFFITS| > 1 for small/medium data set or $> 2\sqrt{p/n}$, large data set.


Cook's Distance

$$D_i = \frac{\sum_{j \neq i}^n (\hat{Y}_j - \hat{Y}_{j(i)})^2}{p \text{MSE}} = \frac{e_i^2}{p \text{MSE}} \frac{h_{ii}}{(1 - h_{ii})^2} \sim F_{p, n-p}$$



$$(\text{DFBETAS})_i = \frac{\hat{\beta}_k - \hat{\beta}_{k(i)}}{\sqrt{\text{MSE}_{(i)} c_{kk}}}$$

where c_{kk} is the diagonal entries of $(\mathbf{X} \mathbf{X})^{-1}$

Flag: DFBETAS > 1 for small/medium data; $> 2\sqrt{\frac{p}{n}}$ for large data.
 Change of signs.


DFIN


 One vs many trouble makers.

Multicollinearity: VIF

📍 Problems of MLCL: X , Extra SSR, $s(\hat{\cdot})$, nonsignificance

📍 Informal Diagnosis

- 📍 Sensitive incl/exclud of X or data

- Nonsignificance on important predictors

- 📍 Wrong sign of estimated


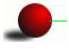


- 📍 Large coefficient in r_{XX} , Large R^2 among X

- Wide confidence intervals of

📍 Variat68 -156.36on I6 -0 1647.61Vat68 Factor (I6 -0 1647.6try

r_{XX} . R

Model Validation

-  Estimation/Fit the past; Predict the future
-  Consistency with New Data
-  Comparison with theoretical expectation, earlier empirical and simulation results
-  Cross-Validation: Use of a holdout sample to check the model and predictive ability.

What's next?