
Building the Regression Model





C. Andy Tsao

Department of Applied Math, National Dong Hwa University





Outline







Learning from Data

-  Controlled Experiments ($F = ma, PV = nRT$)
-  Controlled Experiments with Supplemental variables
-  **Confirmatory** observational studies
-  **Exploratory** observational studies


Data Collection and preparation


-  How the data are collected? (Design, Nature of the study)
-  Data consistency. (Plots and numerical summaries, logical relations)
-  Is this data analysis-ready? (Format checking, file conversion, etc.)
-  **GIGO** (Garbage In; Garbage Out.)


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 vs. Parsimony
Criteria and how to use them? "Good" models?

Selection-I.1

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

 $R_a^2 = 1 - \left(\frac{n-1}{n-p}\right) \frac{SSE_p}{SSTO} = 1 - \frac{MSE_p}{SSTO/(n-1)}$
 $MSE_p = SSE_p/(n-p)$.
ID those with smaller/smallest MSE_p .


 $\Gamma_p = \frac{E(SSE_p)}{SSTO}$

Selection-I.2

 AIC: Akaike's information criterion

$$AIC_p = -2 \ln \text{likelihood} + 2p \propto n \ln SSE_p - n \ln n + 2p.$$

ID models with smaller AIC.

 BIC (or SBC in Text): Schwartz' Bayesian information criterion.





$$BIC_p = -2 \ln \text{likelihood} + p \ln n \propto n \ln SSE_p - n \ln n + (\ln n)p.$$

ID models with smaller BIC.

 Does these criteria make sense?

Increasing/Decreasing in ..

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 < loer order terms
($X^4 < X^1$)
 - Interaction terms < main effect terms
($X_1X_2 < X_1$ 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 Y vs X_q while X_1, \dots, X_{q-1} in the model.

 Partial Regression Plots:

$e(Y|X_1, \dots, X_{q-1})$ vs. $e(X_q|X_1, \dots, X_{q-1})$.

$e(Y|X_1, \dots, X_{q-1})$: residual of Y regresses on X_1, \dots, X_{q-1}

$e(X_q|X_1, \dots, X_{q-1})$: residual of X_q regresses on X_1, \dots, X_{q-1}

 Why bother?

Outliers-I

The model (fitted) shouldn't be affected by just a few points.

Outliers-II

 Studentized Deleted Residual:

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


 Hat matrix Leverage values ! Outlying X

$$0 < h_{ii} < 1; \quad \sum_{i=1}^n h_{ii} = p:$$

$$h_{ii} = \frac{p}{n}: \text{ } 2p/n, \text{ extreme } h_{ii}, \text{ outside } (0.2; 0.5)$$

$$h_{\text{new}} = X_{\text{new}}^0 (X^0 X)^{-1} X_{\text{new}} \text{ for hidden extrapolation.}$$

Influential obs

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

 Cook's Distance

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



$$(DFBETAS)_i = \frac{\hat{\beta}_k - \hat{\beta}_k(i)}{\sqrt{MSE_{(i)}c_{kk}}}$$

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

Flag: $DFBETAS > 1$ for small/medium data; $> 2/\sqrt{n}$.

Change of signs.

 DFINF

 One vs many trouble makers.


Multicollinearity: VIF

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

Remedial Measure

For unequal error variances, high multicollinearity,
influential obs



Model Validation

