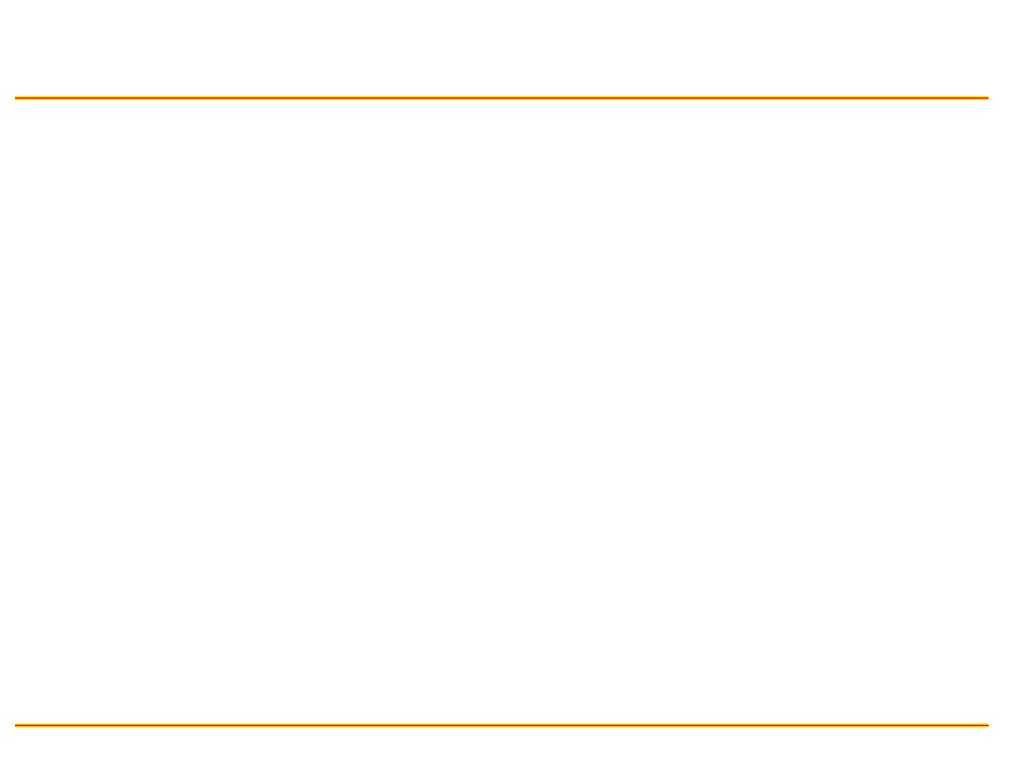# Building the Regression Model

C. Andy Tsao

Department of Applied Math, National Dong Hwa University

# Outline

# 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

# Selection-I.1

- $\mathbf{R}_p^2 = 1 - \frac{SSE_p}{SSTO}.$

# Selection-I.2

- AIC: Akaike's information criterion
  $$\text{AIC}_p = -2 \ln \textbf{likelihood} + 2\textbf{p} \propto \textbf{n} \ln \textbf{SSE}_p - \textbf{n} \ln \textbf{n} + 2\textbf{p}.$$
  ID models with smaller AIC.

-

# Selection-III

**General Linear Test Approach** Given $Y\,|\,X_1, \cdots, X_2$
Should $X_3, \cdots, X_5$ be added?.
Testing $H_0 : \mathbf{Reduced}$ vs $H_1 : \mathbf{Full}$

- Fit the *full* model $(Y\,|\,X_1, \cdots, X_5)$ and get $\mathbf{SSE(F)}, \mathbf{df}_F$

- Fit the *reduced* model $(Y\,|\,X_1, \cdots, X_2)$ and get $\mathbf{SSE(R)}, \mathbf{df}_R$

- Calculate
$$\mathbf{F} = \frac{\frac{SSE(R)}{df_R} \; \frac{SSE(F)}{df_F}}{} \Big/ \frac{SSE(F)}{df_F} \sim \mathbf{F}_{df_R \; df_F, df_F}$$
Then perform a level test.

# 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
  $(\mathbf{X}^4 < \mathbf{X}^1)$

  Interaction terms $<$ main effect terms
  $(\mathbf{X}_1\mathbf{X}_2 < \mathbf{X}_1 \text{ or } \mathbf{X}_2)$

- Chapter 10 Variable Selection of Faraway, J. (2002). Also his Chapter 11 is highly recommended

# Improper functional form of a predic

- Goal: Detect the suitable form of $\mathbf{Y}$ vs $\mathbf{X}_q$ while $\mathbf{X}_1, \cdots, \mathbf{X}_{q-1}$ in the model.

- Partial Regression Plots:
  $\mathbf{e}(\mathbf{Y}|\mathbf{X}_1, \cdots, \mathbf{X}_{q-1})$ vs. $\mathbf{e}(\mathbf{X}_q|\mathbf{X}_1, \cdots, \mathbf{X}_{q-1})$.

  $\mathbf{e}(\mathbf{Y}|\mathbf{X}_1, \cdots, \mathbf{X}_{q-1})$: residual of $\mathbf{Y}$ regresses on $\mathbf{X}_1, \cdots, \mathbf{X}_{q-1}$

  $\mathbf{e}(\mathbf{X}_q|\mathbf{X}_1, \cdots, \mathbf{X}_{q-1})$: residual of $\mathbf{X}_q$ regresses on $\mathbf{X}_1, \cdots, \mathbf{X}_{q-1}$

- Why bother?

# Outliers-I

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

- LSE is EXTREMELY sensitive to outliers. ExampX0.

- Detection: Residual-based tests and pXots towards outlying $\mathbf{Y}$. Why? What to expect?

    Semistudentized residual: Same scale (Naive).
    $$\mathbf{e}_i^* = \frac{e_i}{\sqrt{MSE}}, \mathbf{e}_i = \mathbf{Y}_i - \hat{\mathbf{Y}}_i$$

    Studentized residual: In the same scale (Refined).
    $$\mathbf{r}_i = \frac{e_i}{\sqrt{MSE(1-h_{ii})}} \text{ since}$$

    $$^2(\mathbf{e}) = {}^2(\mathbf{I} - \mathbf{H}), \quad \mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'.$$

    Deleted Residuals: With or Without You. Outlying $\mathbf{Y}$.
    $$\mathbf{d}_i = \mathbf{Y}_i - \hat{\mathbf{Y}_{i(i)}}$$

# Outliers-II

- Studentized Deleted Residual:
  $$\mathbf{t}_i = \frac{d_i}{s(d_i)} \text{ where } \mathbf{s}(\mathbf{d}_i) = \sqrt{\mathbf{MSE}_{(i)}(1 - \mathbf{h}_{ii})}$$

- Hat matrix Leverage values $\rightarrow$ Outlying $\mathbf{X}$
  $$0 \leq \mathbf{h}_{ii} \leq 1, \qquad \sum_{i=1}^{n} \mathbf{h}_{ii} = \mathbf{p}.$$

# Influential obs

- $$(\mathbf{DFFITS})_i = \frac{\widehat{\mathbf{Y}}_i - \widehat{\mathbf{Y}_{i(i)}}}{\sqrt{\mathbf{MSE}_{(i)h_{ii}}}}$$ Flag: If $|\mathbf{DFFITS}| > 1$ for

  small/medium 13433] and $>1$

# Multicollinearity: VIF

- Problems of MLCL: X, Extra SSR, $\mathbf{s}(\hat{\phantom{o}})$, nonsignificance

- Informal Diagnosis

  Sensitive incl/exclud of $\mathbf{X}$ or data

  Nonsignificance on important predictors

  Wrong sign of estimated

  Large coefficient in $\mathbf{r}_{XX}$, Large $\mathbf{R}^2$ among $\mathbf{X}$

  Wide confidence intervals of

- Variation Inflation Factor $(TL^{-1})$ $\mathbf{V\,IF}_k$ diagonal entry of $\mathbf{r}_{XX}$.

$$(\mathbf{V\,IF})_k = (1 - \mathbf{R}_k^2)^{-1},$$

$\mathbf{R}_k^2$: $\mathbf{R}^2$ of $\mathbf{X}_k$ regressing on the other $\mathbf{X}'\mathbf{s}$.

Flag: Larger than 10 or $\gg \mathbf{V\bar{I}F}$

# 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?**