

# Model Mis- specification

(c) Pongsa Pornchaiwiseskul, Faculty of Economics,  
Chulalongkorn University

1

## Covered Topics

- Functional forms
- Underfitting
- Overfitting
- linearity vs non-linear  
(Ramsey's RESET)

(c) Pongsa Pornchaiwiseskul, Faculty of Economics,  
Chulalongkorn University

2

# Functional Forms (1)

- **Linear linear**
- linear log
- linear reciprocal
- **quadratic (polynomial)**
- **interaction (cross terms)**
- **log linear**
- log reciprocal
- log quadratic
- **log log**
- **logistic**

# Functional Forms (2)

**log-log:**  $\ln Y_i = \beta_1 + \beta_2 \ln X_i + \varepsilon_i$

**log-linear:**  $\ln Y_i = \beta_1 + \beta_2 X_i + \varepsilon_i$

**interaction:**

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 (X_{2i} X_{3i}) + \varepsilon_i$$

**logistic:**  $\ln \frac{Y_i}{1 - Y_i} = \beta_1 + \beta_2 X_i + \varepsilon_i$

Note that all are linear in parameters.  
==>OLS applies.

## Functional Forms (3)

Assumption  $\implies$  form. For example,

- constant elasticity  $\implies$  double log (log log)
- constant rate of change  $\implies$  log linear
- Y between 0 and 1  $\implies$  logistic
- variable slope  $\implies$  interaction or polynomial
- combination, e.g., log-log+interaction

There could be more than one that fit.

## Box-Cox Transformation (1)

Box-Cox transformation for X

$$B(X, \lambda) = \frac{X^\lambda - 1}{\lambda}$$

Note that  $B(X, 1) = X - 1$  and  $B(X, 0) = \ln X$   $\Leftarrow$  Why?

Model 
$$\frac{Y_i^\lambda - 1}{\lambda} = \beta_1 + \beta_2 \frac{X_i^\lambda - 1}{\lambda} + \varepsilon_i$$

$\lambda = 1 \implies$  lin-lin model

$\lambda = 0 \implies$  double log model

Otherwise, non-linear model (need NLS or MLE)

## Box-Cox Transformation (2)

$$\frac{Y_i^\lambda - 1}{\lambda} = \beta_1 + \beta_2 \frac{X_i^\mu - 1}{\mu} + \varepsilon_i$$

- Use NLS or MLE to select the best value of  $(\lambda, \mu)$ .
- No need to pre-choose the specific functional form of the model.
- Require more computational effort. No big deal.

## Explanatory Variables

The complete list of X's is purely based on theoretical reasons.

Underfitting = exclusion of relevant variable X's

Overfitting = inclusion of irrelevant variable X's

What are their effects?

# Under-fitting (1)

Let  $X_K$  be the omitted variable with  $\beta_K \neq 0$

$$Y_i = \gamma_1 X_{1i} + \gamma_2 X_{2i} + \dots + \gamma_{K-1} X_{K-1,i} + v_i$$

## Effects

OLS estimator of  $\gamma$  will be a biased estimator of  $\beta$ . if the omitted variable is related to the remaining  $X$ 's. True?

# Under-fitting (2)

Let  $X_{Ki} = \theta_1 X_{1i} + \theta_2 X_{2i} + \dots + \theta_{K-1} X_{K-1,i} + \xi_i$

with  $\theta_k \neq 0$  for some  $k=1, \dots, K-1$

If  $\theta_k \neq 0$ ,  $\hat{\gamma}_k$  will be a biased estimator of

$\beta_k$  because  $\gamma_k = \beta_k + \beta_K \theta_k$ . Note that

$\gamma_k$  includes not only the effect of  $X_k$  but also that of the omitted variable ( $X_K$ ).

# Under-fitting (3)

Substitute into the exact model

$$\begin{aligned} Y_i &= \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_{K-1} X_{K-1,i} \\ &+ \beta_K (\theta_1 X_{1i} + \theta_2 X_{2i} + \dots + \theta_{K-1} X_{K-1,i} + \xi_i) + \varepsilon_i \\ &= (\beta_1 + \beta_K \theta_1) X_{1i} + (\beta_2 + \beta_K \theta_2) X_{2i} \\ &+ \dots + (\beta_{K-1} + \beta_K \theta_{K-1}) X_{K-1,i} \\ &+ (\varepsilon_i + \beta_K \xi_i) \end{aligned}$$

# Over-fitting (1)

Define  $Z =$  irrelevant variable

$$Y_i = \gamma_1 X_{1i} + \gamma_2 X_{2i} + \dots + \gamma_K X_{Ki} + \delta Z_i + v_i$$

Since  $Z$  is known to be irrelevant, the real  $\delta = 0$ .

## Effects

- OLS estimator of  $\gamma$  is an unbiased estimator of  $\beta$  because  $\gamma = \beta$ .
- loss of accuracy. Why?

# Mis-fitting (1)

## Rules

- Include all the explanatory variables suggested by the underlying theories.
- Excluding them requires theoretical explanation.
- Even if the test statistic indicates insignificance, leave them in the model to avoid unbiasedness. Low statistic does not imply irrelevance. Data just can't reveal it.

# Mis-fitting (2)

- Add variables to test their relevancy.  
Remove the insignificant variables.  
Further theories could be developed if  
the test indicates their significance.

# Ramsey's RESET (1)

## REgression Specification Error Test

- test the linear (in X) model against unspecified non-linear model

## Concept

Using Taylor series expansion, a non-linear model can be expressed as a polynomial model. If the exact model is non-linear, using a linear model is equivalent to omitting variables (high order terms)

# Ramsey's RESET (2)

## Test Equation

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} \\ + \delta_2 \hat{Y}_i^2 + \delta_3 \hat{Y}_i^3 + \dots + \delta_M \hat{Y}_i^M + \varepsilon_i$$

Perform an F-test or Chi-square test

$$H_0 : \delta_2 = \delta_3 = \dots = \delta_M = 0$$

$$H_1 : \delta_2 \neq \delta_3 \neq \dots \neq \delta_M \neq 0$$

EViews can do RESET.



# Ramsey's RESET (3)

## Notes

- Start with a square term
- The highest order  $M$  must be pre-selected.  $M-1$  terms added.
- Reject  $H_0 \Rightarrow$  need a new model
- The test does not suggest the form of a new model. Try a polynomial order  $M$ .

## Model

## Selection Criteria

# Covered Topics

- Criteria
- Adding/dropping variables
- Double linear against double log
- Linear with different X's

## Forecasting Error

Forecasting error =  $Y_i - \hat{Y}_i$

Linear model  $\hat{Y}_i = \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_K X_{Ki}$

Double log

$$\widehat{\ln Y}_i = \hat{\beta}_1 \ln X_{1i} + \hat{\beta}_2 \ln X_{2i} + \dots + \hat{\beta}_K \ln X_{Ki}$$

$$\hat{Y}_i = \exp(\widehat{\ln Y}_i)$$

# General Selection Criteria (1)

Root Mean Square Error

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$$

Mean Absolute Error

$$MAE = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|$$

# General Selection Criteria (2)

Mean Absolute Percentage Error

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{Y_i - \hat{Y}_i}{Y_i} \right|$$

Theil's inequality coefficient

$$U = \frac{\sqrt{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}}{\sqrt{\sum_{i=1}^n Y_i^2} + \sqrt{\sum_{i=1}^n \hat{Y}_i^2}}$$

# Relevancy of Variables (1)

Should  $Z_1, Z_2, \dots, Z_M$  be added?

Test Equation

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} \\ + \delta_1 Z_{1i} + \delta_2 Z_{2i} + \dots + \delta_M Z_{Mi} + \varepsilon_i$$

Perform an F-test or  $\chi^2$ -test on

$$H_0 : \delta_1 = \delta_2 = \dots = \delta_M = 0$$

$$H_1 : \delta_1 \neq \delta_2 \neq \dots \neq \delta_M \neq 0$$

# Relevancy of Variables (2)

Note that  $F_{cal} \sim F(M, n - (K + M))$

$$\chi_{cal}^2 \sim \chi^2(M)$$

**EViews** gives the 2-run F-test and the LR test.

Reject  $H_0$  implies that the variables may be relevant. It needs explanation. No harm to the unbiasedness of the estimator of parameters  $\beta$  which have been already included.

Another test equation is the old residuals against all the X's (old & new) w/ a constant term.

# Redundancy of Variables (1)

Should  $X_{K-M+1}, X_{K-M+2}, \dots, X_K$  be dropped?

Test Equation

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_{K-M} X_{K-M,i} \\ + \beta_{K-M+1} X_{K-M+1,i} + \dots + \beta_K X_{Ki} + \varepsilon_i$$

Perform an F-test or  $\chi^2$ -test on

$$H_0 : \beta_{K-M+1} = \beta_{K-M+2} = \dots = \beta_K = 0$$

$$H_1 : \beta_{K-M+1} \neq \beta_{K-M+2} \neq \dots \neq \beta_K \neq 0$$

# Redundancy of Variables (2)

Note that  $F_{cal} \sim F(M, n - K)$

$$\chi_{cal}^2 \sim \chi^2(M)$$

EViews also gives 2-run F-test and LR test results.

Accept  $H_0$  does not imply that the variables are not relevant. It just says that they are redundant.

With the other  $X$ 's in the model, they are not needed to explain the variation of  $Y$ . If they are dropped, there might be estimation bias but no harm to prediction of  $Y$ .

# MWD Test (1)

MacKinnon-White-Davidson Test

Choose between lin-lin and log-log with same X and Y

Model A

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} + \varepsilon_i$$

Model B

$$\ln Y_i = \gamma_1 \ln X_{1i} + \gamma_2 \ln X_{2i} + \dots + \gamma_K \ln X_{Ki} + \varepsilon_i$$

# MWD Test (2)

Note that  $R^2$  cannot be used to judge.

Given  $\hat{Y}_i$  the fitted value from Model A

$\widehat{\ln Y}_i$  the fitted value from Model B

Define

$$Z_{Ai} = \ln \hat{Y}_i - \widehat{\ln Y}_i$$

$$Z_{Bi} = \hat{Y}_i - \exp(\widehat{\ln Y}_i)$$

# MWD Test (3)

## Model A'

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} \\ + \delta_A Z_{Ai} + \varepsilon_i$$

t-Test on  $H_0 : \delta_A = 0$

$$H_1 : \delta_A \neq 0$$

Accept  $H_0 \Rightarrow$  double lin “encompasses”  
double log

# MWD Test (4)

## Model B'

$$\ln Y_i = \gamma_1 \ln X_{1i} + \gamma_2 \ln X_{2i} + \dots + \gamma_K \ln X_{Ki} \\ + \delta_B Z_{Bi} + \varepsilon_i$$

t-Test on  $H_0 : \delta_B = 0$

$$H_1 : \delta_B \neq 0$$

Accept  $H_0 \Rightarrow$  double log “encompasses”  
double line

# MWD Test (5)

## Conclusion

	$\delta_B = 0$	$\delta_B \neq 0$
$\delta_A = 0$	No clear preference	Double linear is chosen
$\delta_A \neq 0$	Double log is chosen	Neither model good enough

# J-test (1)

## Davidson-MacKinnon's J-test

Choose between two linear models with different set of explanatory variables but same dependent variable.

$R^2$  cannot be used either.

Model C  $Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} + \varepsilon_i$

Model D  $Y_i = \gamma_1 Z_{1i} + \gamma_2 Z_{2i} + \dots + \gamma_L Z_{Li} + \varepsilon_i$



# J-test (2)

## Model C'

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} + \delta_C \hat{Y}_{Di} + \varepsilon_i$$

where  $\hat{Y}_i^D$  is the fitted value from Model D

t-Test on  $H_0 : \delta_C = 0$

$$H_1 : \delta_C \neq 0$$

Accept  $H_0 \Rightarrow$  model C “encompasses”  
model D

# J-test (3)

## Model D'

$$Y_i = \gamma_1 Z_{1i} + \gamma_2 Z_{2i} + \dots + \gamma_L Z_{Li} + \delta_D \hat{Y}_{Ci} + \varepsilon_i$$

where  $\hat{Y}_{Ci}$  is the fitted value from Model C

t-Test on  $H_0 : \delta_D = 0$

$$H_1 : \delta_D \neq 0$$

Accept  $H_0 \Rightarrow$  model D “encompasses”  
model C

# J-test (4)

## Conclusion

	$\delta_D = 0$	$\delta_D \neq 0$
$\delta_C = 0$	No clear preference	Model C is chosen
$\delta_C \neq 0$	Model D is chosen	Neither model good enough