

Mon 2019-11-4

Chapter 19 → Term paper

Structural Multicollinearity Data collinearity

Imperfect (near) Multicollinearity

$$X_{1i} = \alpha_0 + \alpha_1 X_{2i} + u_i$$

Where u_i can be regarded as an approximation. Imperfect multicollinearity ~~can be due to~~ is a strong linear relationship between variables. Multicollinearity is a sample phenomenon.

- i) If the goal is simply to predict y from a set of X variables, then multicollinearity is not a problem
- ii) If the goal is to understand how X influences y , it is a problem
- iii) P-values can be misleading
- iv) Confidence intervals are off

Correlation - linear strength between variables

Collinearity - One of the explanations is linearly dependent of another

Sources of Multicollinearity

- 1) Data collection method
- 2) Constraints on model or population
- 3) Model specification
- 4) Over defined model

Consequences of Multicollinearity

• Estimates will remain unbiased

• Computed & will fall

• High R^2 value but less significant β s
Wrong sign for regression coefficients

Auxiliary Regression

Regres each explanatory variable on the other explanatory

Variance Inflation Factor

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R_j^2}$$

Ready

• Drop one of the variables

• Acquire more data

• Rethinking the model: level-level, level-log, log-log

Wed 2019 - 11 - 06

Chapter 8 - Heteroscedasticity

HetSc is a violation of CLRM. u_i must have constant variance σ^2 :

$$u \sim N(0, \sigma^2)$$
$$\text{Var}(u_i) = \sigma^2$$

If that is not the case:

$$\text{Var}(u_i) = \sigma_i^2 \quad (\text{non-constant variance})$$

Types of HetSc:

- Conditional - nonconstant volatility when periods of high & low volatility cannot be predicted
- Unconditional - these periods can be predicted

Conditional \rightarrow not predictable by nature

What is the nature of HetSc?

Omitted Independent Variables

How HetSc can be tested

What are consequences

Formal Method tests for HetSc

How do we detect it \rightarrow in a given situation

What do we do if Heteroscedasticity is observed

Homoscedasticity is when:

$$E(u_i^2) = \sigma^2 \text{ for all } i$$

As opposed to Heteroscedasticity:

$$E(u_i^2) = \sigma_i^2$$

i.e., conditional variances are not constant.

$$\text{OR } E(u_i^2 | x_i) \neq E(u_j^2 | x_j) \text{ for some } i, j$$

Causes of Heteroscedasticity

Unconditional Heteroscedasticity can be used when variables have identifiable seasonal variability, such as electricity usage.

Why Heteroscedasticity happens

- i) Following error learning models: As people learn their errors of behavior become smaller, so σ_i^2 is expected to decrease
- ii) Growth of income: In a regression of savings, an income is likely to find σ_i^2 increasing with income
- iii) Data-collecting technique: As techniques improve σ_i^2 is likely to decrease. e.g. larger banks with better equipment are likely to commit fewer statement errors
- iv) presence of outliers: including or excluding outliers is up to the discretion of the analyst
- v) Misspecification of model: the wrong variables are included, or important ones excluded.
- vi) Skewness in distn of one or more regressors
- vii) Other sources:
 - 1) incorrect data transformation form
 - 2) incorrect functional form (e.g. linear vs log)

Pure vs impure Heteroscedasticity

- i) Pure is caused by the error term of the correctly specified equation
- ii) Impure is caused by specification errors such as an omitted variable

Omitted variables: impure Heteroscedasticity

Effects:

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

$$E(u_i) = 0$$

$$\text{var}(u_i) = \sigma^2$$

$$\text{cov}(u_i, u_j) = 0 \quad \text{if } i \neq j$$

Special considerations

Heteroscedasticity is frequent in cross-sectional data

Examples:

- i) Operating income versus advertising expenses of 30 firms

How it can be tested:

looking at the plot

(calculating σ^2 under Hct Sc
OLS is unbiased under Hct Sc

$$\text{Var}(\hat{\beta}_0 | X_i) = \sum \text{Var}:$$

$$\downarrow \\ \text{var}(\hat{\beta}_1 | X_i) = \sum (x_i - \bar{x})^2 \sigma_e^2$$

$$\downarrow \\ \text{var}(\hat{\beta}_1) = \sum (x - \bar{x})$$

$$t = \frac{\text{estimate} - \text{hypothesized value}}{\text{standard error}}$$

Only difference between t and the usual OLS t is how the SE is computed

(consequences of Hct Sc)

- a) OLS estimators are still linear
- b) " " " still unbiased
- c) " " no longer have min variance
- d) Usual OLS formulas to estimate are biased
- e) This bias arises from the fact that $\hat{\sigma}^2$ is no longer unbiased for σ^2
- f) CIs and hypothesis on t and f dists are unreliable
- g) Loss of efficiency

There are several tests for H₀: Sc
 We will do White test and Breusch-Pagan test
 It is

Chapter 4 questions 2 and 3

2(b) i) $H_0: \beta_3 = 0 \quad H_1: \beta_3 > 0$

ii) log-level $\% \Delta y = (100 \beta_3) \Delta x$

$$100 \cdot 0.00024 \cdot 50 \\ = 1.2\%$$

iii) $t_{\hat{\beta}_3} = \frac{0.00024}{0.00054} = 0.44 \quad t_c = 11.282$

At the 90% significance level we fail to reject H₀

iv) No, it fails to be significant and has only a marginal effect on the model

6) i) $H_0: \beta_0 = 0 \quad \alpha = 0.05 \quad df = 86 \quad \text{---}$
 $H_1: \beta_0 \neq 0 \quad t_c = 1.988 \quad t_{\hat{\beta}_0} = \frac{\hat{\beta}_0 - \beta_0}{Sc(\hat{\beta}_0)} = \frac{-14.47}{16.67} = -0.868$

Fail to reject H₀

$$H_0: \beta_1 = 1 \quad \alpha = 0.05 \quad t_c = 1.988$$

$$H_1: \beta_1 \neq 1 \quad df = 86 \quad t_{\hat{\beta}_1} = \frac{\hat{\beta}_1 - \beta_1}{Sc(\hat{\beta}_1)} = \frac{-0.976 - 1}{0.049} = -4.90$$

Fail to reject

Brensch-Pagan test for condition Heteroskedasticity

If the χ^2 value is significant with p-value below an appropriate threshold (e.g. $p < 0.05$), H_0 : homoskedasticity is rejected, and Heteroskedasticity is present

$$Y_i = \beta_1 + \beta_2 X_{2i} + \dots + \beta_n X_{ni} + u_i$$

$$\sigma_i^2 = \delta(\alpha_1 + \alpha_2 Z_{2i} + \dots + \alpha_m Z_{mi})$$

Assume

$$\sigma_i^2 = \alpha_1 + \alpha_2 Z_{2i} + \dots + \alpha_m Z_{mi}$$

If $\alpha_2 = \alpha_3 = \dots = \alpha_m = 0$,

this implies $\sigma_i^2 = \alpha_1$, which is constant and the model is Homoskedastic, so,

$$H_0: \alpha_2 = \alpha_3 = \dots = \alpha_m = 0, \quad \sigma_i^2 = \alpha_1$$

In Stata: `rvfplot, yline(0)`

~~reg i.~~ "hettest" → gives p-value. If $< 5\%$, there is Heteroskedasticity

White test

$$y = \beta_1 + \beta_2 X_2 - \beta_3 X_3 + u;$$

$$\tilde{y} = \hat{\beta}_1 + \hat{\beta}_2 X_2 - \hat{\beta}_3 X_3$$

$$\tilde{u}^2 = \alpha_1 + \alpha_2 X_{2i} + \alpha_3 X_{3i} + \alpha_4 X_{4i}^2 + \alpha_5 X_{5i}^2 + \alpha_6 X_{6i} X_{3i} + v_i^2$$

$$H_0: \alpha_2 = \alpha_3 = \dots = \alpha_6 = 0$$

$$H_1: \alpha_2 \neq 0, \alpha_3 \neq 0, \dots, \alpha_6 \neq 0$$

Lagrange Multiplier (λ_M) is R^2 times sample size n

$$LM = nR^2$$

If $nR^2 > \chi_{df}^2$

With $df = k - 1$

If $(nR^2 > \chi_{df}^2) \in$
there is heteroskedasticity

3

Cross-sectional data: Hetero is very frequent

Panel data: combined cross-sectional and time-series data

Recreational measures

If variance is known

Divide both sides by σ_i

$$\frac{y_i}{\sigma_i} = \beta_0 \left(\frac{1}{\sigma_i} \right) + \beta_1 \left(\frac{x_i}{\sigma_i} \right) + \frac{u_i}{\sigma_i}$$

Let $v_i = \frac{u_i}{\sigma_i}$, the "transformed" error term

Now square v_i

$$v_i^2 = \frac{u_i^2}{\sigma_{\epsilon_i}^2}$$

$$E(v_i^2) = E\left(\frac{u_i^2}{\sigma_{\epsilon_i}^2}\right) = \frac{1}{\sigma_{\epsilon_i}^2} E(u_i^2)$$

$$\rightarrow E(u_i^2) = \sigma_{\epsilon_i}^2$$

$$E(v_i^2) = \left(\frac{1}{\sigma_{\epsilon_i}^2}\right) \sigma_{\epsilon_i}^2 = 1$$

↳ constant, therefore
transformed model
is homoskedastic

This is called method of weighted Least squares,
or WLS

olsrr package in R provides:

- o Barlett Test
- * o Bresch-Pagan Test
- o Score test
- * o F test

If σ^2 is unknown
Two subcases:

i) Error variance is proportional x_i
Square root transformation
so

$$E(u_i^2) = \sigma^2 x_i$$

consider

$$Y_i = \beta_1 + \beta_2 x_i + u_i$$

$$\downarrow \frac{C}{\sqrt{x_i}}$$

$$\frac{Y_i}{\sqrt{x_i}} = \frac{\beta_1}{\sqrt{x_i}} + \beta_2 \cdot \frac{x_i}{\sqrt{x_i}} + \frac{u_i}{\sqrt{x_i}}$$

↓

$$\frac{Y_i}{\sqrt{x_i}} = \beta_1 \frac{1}{\sqrt{x_i}} + \beta_2 \sqrt{x_i} + \frac{u_i}{\sqrt{x_i}}$$

let $V_i = \frac{u_i}{\sqrt{x_i}}$ $V_i^2 = \frac{u_i^2}{x_i}$

$$E(u_i^2) = \sigma^2 x_i \quad E(V_i^2) = \frac{\sigma^2 x_i}{x_i} = \sigma^2$$

so it's homoscedastic

Situation 2: Error variance is proportional to mean of x_i

Wed 2019 - 11 - 13

$$E(u^2) = \sigma^2 x_i^2 = \sigma^2 [E(x_i)]^2$$

Assume the model: mean transform.,

$$y_i = \beta_1 + \beta_2 x_i + u_i$$

↓ divide by x_i

$$\frac{y_i}{x_i} = \frac{\beta_1}{x_i} + \frac{\beta_2 x_i}{x_i} + \frac{u_i}{x_i} \approx \frac{y_i}{x_i} = \frac{\beta_1}{x_i} + \beta_2 + \frac{u_i}{x_i}$$

Let $v_i = \frac{u_i}{x_i}$

$$\frac{y_i}{x_i} = \beta_1 \frac{1}{x_i} + \beta_2 + v_i \quad v_i^2 = \frac{u_i^2}{x_i^2}$$

$$E(u_i^2) = E\left(\frac{u_i^2}{x_i^2}\right) = \frac{1}{x_i^2} E(u_i^2)$$

$$E(u_i^2) = \frac{1}{x_i^2} \sigma^2 (x_i^2) = \sigma^2 \leftarrow \text{As a result, transformed model does not differ from } \text{homoscedastic.}\text{E.I.}\text{netto}$$

Chapter 7 - Dummy variables

Qualitative variables

Dummy variables are binary 1 or 0

Used in regression to represent subgroups

Regression models containing only dummy explanatory variables are called analysis of variance (ANOVA)

Dummies are proxies for qualitative ~~other~~ factors

Dummy variable trap
When dummies are multicollinear.

$$\text{Wage} = \beta_0 + \delta_0 \text{female} + \beta_1 \text{education} + \epsilon$$

if $\delta_0 > 0$, women make more money
else, if $\delta_0 < 0$, women make less money

Intercept for men is β_0
for women is $\beta_0 + \delta_0$

Mon 2019 - 11 - 18

Interaction:

Interaction variables: formed by multiplication of two variables

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon_{ui}$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$$

ANOVA: used to analyze if there is a significant difference between groups (categorical variable)

ANOVA can show if survey or experimental results are significant

Look at p-value of β_3

$$Y = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

ANCOVA - Analysis of covariance

Wed 2019-11-20

- Chapter 12 - Autocorrelation / serial correlation
another violation of OLS
Frequent in time series data, as opposed to cross-sectional

The relationship between a given variable and itself over various time intervals, often found repeating patterns which the term "auto-correlation" reflects. (e.g. seasonal effects such as stock prices, GDP)

Occurs when the error term observations in a regression are correlated
Can be a problem when analyzing historical data

Pure auto.: Model is correctly specified but errors are ~~no~~ serial correlation

Most common is 1st order: in which error at time t is related to error @ previous time t-1:

$$\epsilon_t = \rho \epsilon_{t-1} + u_t \quad -1 < \rho < 1$$

↑ population serial correlation

Occurs when assumption assuming uncorrelated observations of the error term is violated.

No autocorrelation means

$$E(u_i, u_j) = 0$$

pure serial correlation is the "default"

Impose serial correlation
Model is not well specified

Reasons for serial correlation

- 1) Inertia: (or sluggishness) frequent in Macro time series data.
- 2) Model specification error.
- 3) Cobweb phenomenon: economic model that explain why prices might be subject to periodic fluctuations - based on time lag between the supply and demand decisions.

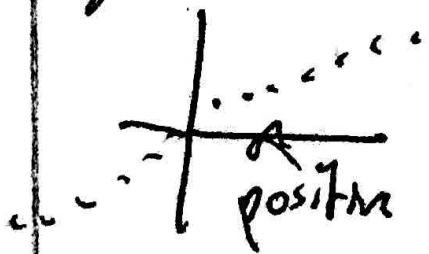
Data manipulation

Time series reg's are involving quarterly data are often derived by averaging the data from 3 months

Spatial auto correlation:

occurs when two errors are spatially/geographically related

Negative auto correlation: switching of positive and negative errors:



$$E_t = \rho E_{t-1} + u_t \quad \leftarrow \text{markos 1st order scheme}$$

↑
error of
eq in question desired
error term

$\rho = 0 \rightarrow$ no error term

$\rho \approx 0 \pm 1 \rightarrow$ strong auto correlate

Most econometricians focus on positive serial correlation

- 1) Pure serial correlation does not cause bias in the coefficient estimates
- 2) They are not efficient i.e. not minimum bias
- 3) Estimated ~~variance~~ of OLS estimates are biased

d : Durbin Watson statistic
always between 0 and 4
detects autocorrelation at lag 1

- 2 - no autocorrelation
- $0 < d < 2$ - positive serial correlation
- $2 < d < 4$ - negative " "

- 1) Compute OLS
- 2) " $u_t = Y_t - \hat{\beta}_1 X_1 - \hat{\beta}_2 X_2 - \dots - \hat{\beta}_k X_k$

- 3)
$$d = \frac{\sum_{t=2}^N (\hat{u}_t - \bar{\hat{u}}_{t-1})^2}{\sum \hat{u}_t^2}$$
 d is bounded by 2 distributions:
 d_L and d_U , given by table

ϵ for error term

$$\epsilon_t = \rho \epsilon_{t-1} + u_t$$

$$\hat{\epsilon}_t = \delta_0 + \delta_1 \hat{\epsilon}_{t-1} + \text{sample error term}$$

δ_1 corresponds to ρ

Hence the test is

$$H_0: \rho = 0$$

$$H_1: \rho \neq 0$$

Another way to test for auto correlation

$$DW = \sum$$