

제 10 강

이분산성

Heteroskedasticity

이분산성의 본질

- 이분산성(Heteroskedasticity) 은 오차항의 분산이 일정하지 않게 되는 오차항에 있어서의 체계적인 패턴임
- 통상적 최소제곱추정(Ordinary least squares, OLS)는 모든 관측치들이 동등한 정도로 신뢰할 만하다고 가정함
- 유효한(efficient) 추정을 위해서는 관측치에 대해 가중치를 줌으로써 동일한 오차항 분산을 갖도록 해주어야 함

단순선형모형

$$y_t = \beta_1 + \beta_2 x_t + \varepsilon_t$$

0의 기대값(zero mean): $E(\varepsilon_t) = 0$

동분산성(homoskedasticity): $\text{var}(\varepsilon_t) = \sigma^2$

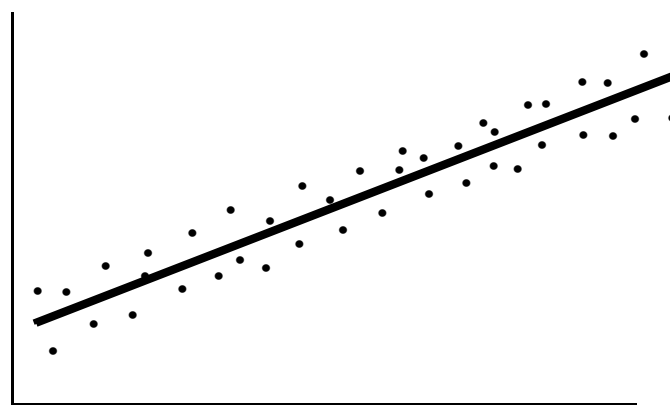
비자기상관(nonautocorrelation):
 $\text{cov}(\varepsilon_t, \varepsilon_s) = 0 \quad t \neq s$

heteroskedasticity: $\text{var}(\varepsilon_t) = \sigma_t^2$

동분산성

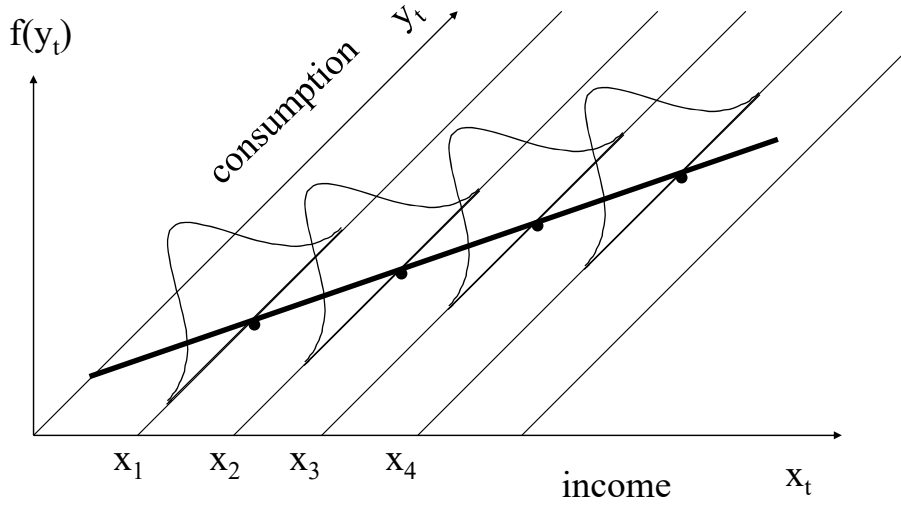
consumption

y_t



income x_t

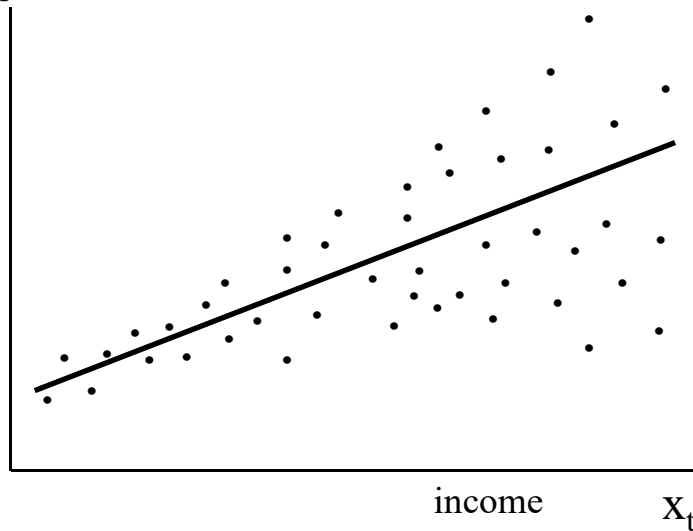
동분산성



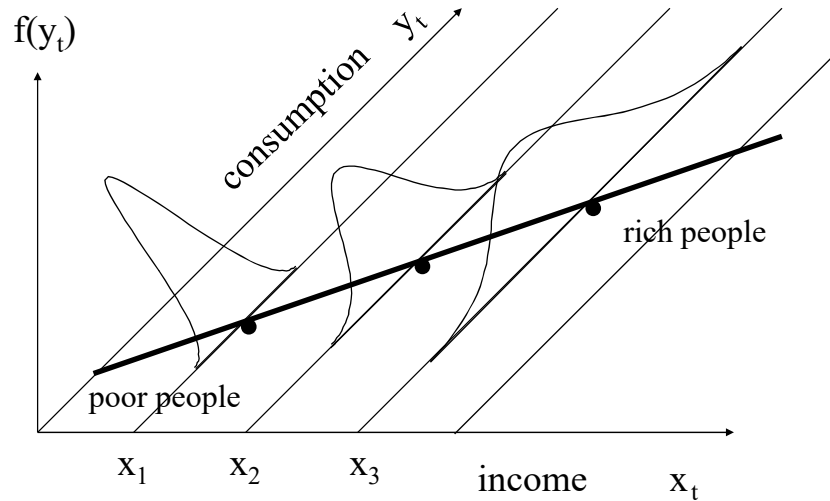
이분산성

consumption

y_t



이분산성



- 최소제곱추정은 일치 추정이며, 여전히 선형이고 불편 추정
- 최소제곱 추정은 유효 추정이 아님, 즉 BLUE가 아님
- 통상적인 공식은 최소제곱 추정량에 대한 잘못된 표준 오차를 제공함
- 이러한 잘못된 표준오차에 근거한 신뢰구간이나 가설검정은 모두 잘못된 것임

$$y_t = \beta_1 + \beta_2 x_t + \varepsilon_t$$

heteroskedasticity: $\text{var}(\varepsilon_t) = \sigma_t^2$

최소제곱추정량의 분산에 대한 잘못된 공식

$$\text{var}(b_2) = \frac{\sigma^2}{\sum (x_t - \bar{x})^2}$$

최소제곱추정량의 분산에 대한 올은 공식

$$\text{var}(b_2) = \frac{\sum \sigma_t^2 (x_t - \bar{x})^2}{[\sum (x_t - \bar{x})^2]^2}$$

이분산성

화이트의 최소제곱 추정량의 분산에 대한 추정량

$$\hat{\text{var}}(b_2) = \frac{\sum \hat{\varepsilon}_t^2 (x_t - \bar{x})^2}{[\sum (x_t - \bar{x})^2]^2}$$

대표본에서, 화이트의 표준오차는 적절함 (i.e. 일치추정량)

1. 비례적(Proportional) 이분산성
(continuous function of x_t , for example)
2. 분할된(Partitioned) 이분산성
(discrete categories/groups)

비례적 이분산성

$$y_t = \beta_1 + \beta_2 x_t + \varepsilon_t$$

$E(\varepsilon_t) = 0$	$\text{var}(\varepsilon_t) = \sigma_t^2$	$\text{cov}(\varepsilon_t, \varepsilon_s) = 0 \quad t \neq s$
------------------------	--	---

where $\sigma_t^2 = \sigma^2 x_t$

The variance is assumed to be proportional to the value of x_t

비례적 이분산성

$$y_t = \beta_1 + \beta_2 x_t + \varepsilon_t$$

variance: $\text{var}(\varepsilon_t) = \sigma_t^2 \longrightarrow \sigma_t^2 = \sigma^2 x_t$

standard deviation: $\sigma_t = \sigma \sqrt{x_t}$

To correct for heteroskedasticity divide the model by $\sqrt{x_t}$

$$\frac{y_t}{\sqrt{x_t}} = \beta_1 \frac{1}{\sqrt{x_t}} + \beta_2 \frac{x_t}{\sqrt{x_t}} + \frac{\varepsilon_t}{\sqrt{x_t}}$$

비례적 이분산성

$$\frac{y_t}{\sqrt{x_t}} = \beta_1 \frac{1}{\sqrt{x_t}} + \beta_2 \frac{x_t}{\sqrt{x_t}} + \frac{\varepsilon_t}{\sqrt{x_t}}$$

$$y_t^* = \beta_1 x_{1t}^* + \beta_2 x_{2t}^* + \varepsilon_t^*$$

$$\text{var}(\varepsilon_t^*) = \text{var}\left(\frac{\varepsilon_t}{\sqrt{x_t}}\right) = \frac{1}{x_t} \text{var}(\varepsilon_t) = \frac{1}{x_t} \sigma^2 x_t$$

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

ε_t is **heteroskedastic**, but ε_t^* is **homoskedastic**.

비례적 이분산성 - 추정방법

다음 3단계로 이루어지는 가중최소제곱 (weighted least squares) 추정 :

1. 어떤 변수가 이분산성에 비례적인가를 결정함 (x_t in previous example).
2. 원래의 모형의 모든 항들을 그 변수의 제곱근으로 나누어 줌 (divide by $\sqrt{x_t}$).
3. 새로운 종속변수와 설명변수 (y_t^* , x_{t1}^* and x_{t2}^* but **no intercept**) 들을 가지는 변환된 모형에 대해 최소제곱 추정을 적용

분할된 이분산성

$$y_t = \beta_1 + \beta_2 x_t + \varepsilon_t$$

y_t = 1에이커 당 옥수수 수확량 $t = 1, \dots, 100$

x_t = 1에이커 당 공급된 물 (rain or other)

구 옥수수 종자의 오차항 분산: $\text{var}(\varepsilon_t) = \sigma_1^2$
 $t = 1, \dots, 80$

신 옥수수 종자의 오차항 분산: $\text{var}(\varepsilon_t) = \sigma_2^2$
 $t = 81, \dots, 100$

분할된 이분산성 - 추정방법

구 종자:	$y_t = \beta_1 + \beta_2 x_t + \varepsilon_t$	$\text{var}(\varepsilon_t) = \sigma_1^2$
-------	---	--

$$\frac{y_t}{\sigma_1} = \beta_1 \frac{1}{\sigma_1} + \beta_2 \frac{x_t}{\sigma_1} + \frac{\varepsilon_t}{\sigma_1} \quad t = 1, \dots, 80$$

신 종자:	$y_t = \beta_1 + \beta_2 x_t + \varepsilon_t$	$\text{var}(\varepsilon_t) = \sigma_2^2$
-------	---	--

$$\frac{y_t}{\sigma_2} = \beta_1 \frac{1}{\sigma_2} + \beta_2 \frac{x_t}{\sigma_2} + \frac{\varepsilon_t}{\sigma_2} \quad t = 81, \dots, 100$$

분할된 이분산성 - 추정방법

$$\frac{y_t}{\sigma_1} = \beta_1 \frac{1}{\sigma_1} + \beta_2 \frac{x_t}{\sigma_1} + \frac{\varepsilon_t}{\sigma_1} \quad t = 1, \dots, 80$$

$$\frac{y_t}{\sigma_2} = \beta_1 \frac{1}{\sigma_2} + \beta_2 \frac{x_t}{\sigma_2} + \frac{\varepsilon_t}{\sigma_2} \quad t = 81, \dots, 100$$

↓

$$y_t^* = \beta_1 x_{1t}^* + \beta_2 x_{2t}^* + \varepsilon_t^* \quad t = 1, \dots, 100$$

$$\text{var}(\varepsilon_t^*) = 1 : \text{Homoskedasticity}$$

However, σ_1 and σ_2 are unknown. Need to estimate them

그룹별 이분산성 - 추정방법

각 그룹의 자료에 대해 최소제곱추정을 적용.

$\hat{\sigma}_1^2$ provides estimator of σ_1^2 using the 80 observations.

$\hat{\sigma}_2^2$ provides estimator of σ_2^2 using the 20 observations.

일반화된 최소제곱추정(Generalized Least Square):
변수들을 표준적 가정에 부합되게 변환하고, 변환된 변수들에 대해 최소제곱 추정을 적용하여 모수들에 대한 유효한 추정을 얻음

1. 잔차 도표 ⇒ Park Test

오차항의 분산의 이분산성에 있어서 체계적이고 유의한 영향이 의심되는 설명변수를 확인하는데 유용함

2. Goldfeld-Quandt Test

두 그룹의 관측치들 간에 그 신뢰성에 체계적이고 유의한 차이를 확인

3. Breusch-Pagan test

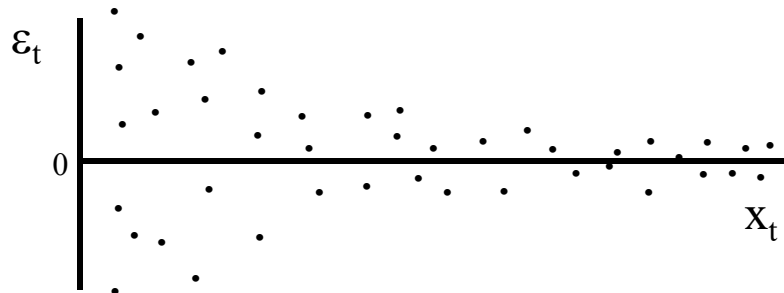
오차항의 분산이 의존할 것으로 의심되는 설명변수들이 특정될 때, 이들의 분산에 대한 체계적이고 유의한 영향을 확인

4. White's heteroscedasticity test

전체적인 설명변수들의 오차항의 분산에 대한 체계적이고 유의한 영향을 확인함 (설명변수들을 특정할 필요 없음)

잔차 도표

의심되는 설명변수 하나를 크기에 따라 Sorting하고 그에 대해 잔차를 그려봄으로써 자료에 이분산적 패턴이 있는가 확인



Park 검정

H_0 : No heteroscedasticity exists i.e., $\text{Var}(\varepsilon_t) = \sigma^2$
(homoscedasticity)

H_1 : Yes, heteroscedasticity exists i.e., $\text{Var}(\varepsilon_t) = \sigma_t^2$

Park test procedures:

1. Run OLS on regression: $Y_t = \beta_1 + \beta_2 X_t + \varepsilon_t$, obtain $\hat{\varepsilon}_t$
2. Take square and take log : $\ln(\hat{\varepsilon}_t^2)$
3. Run OLS on regression: $\ln(\hat{\varepsilon}_t^2) = \alpha_1 + \alpha_2 X_t + v_t$

4. Use t-test to test $H_0 : \alpha_2 = 0$ (Homoscedasticity)

$$H_1 : \alpha_2 \neq 0$$

Suspected variable
that causes
heteroscedasticity

Park 검정- 일반화

이분산성의 구조는 훨씬 복잡할 수 있음:

$$\sigma_t^2 = \sigma^2 \exp\{\alpha_1 z_{t1} + \alpha_2 z_{t2}\}$$

z_{t1} 와 z_{t2} 는 분산이 의존하는 것으로 의심되는 임의의 관측가능한 변수들임

Note: The function $\exp\{\cdot\}$ ensures that σ_t^2 is positive.

Park 검정- 일반화

$$\sigma_t^2 = \sigma^2 \exp\{\alpha_1 z_{t1} + \alpha_2 z_{t2}\}$$

$$\ln(\sigma_t^2) = \ln(\sigma^2) + \alpha_1 z_{t1} + \alpha_2 z_{t2}$$

$$\ln(\sigma_t^2) = \alpha_0 + \alpha_1 z_{t1} + \alpha_2 z_{t2}$$

$$\text{where } \alpha_0 = \ln(\sigma^2)$$

$$H_0: \alpha_1 = 0, \alpha_2 = 0$$

$$H_1: H_0 \text{ not true}$$

Least squares residuals, $\hat{\varepsilon}_t$

$$\ln(\hat{\varepsilon}_t^2) = \alpha_0 + \alpha_1 z_{t1} + \alpha_2 z_{t2} + v_t$$

the usual F test

GQ 검정

골드펠드-퀀트 검정(The Goldfeld-Quandt test)은 비례적 이분산성 혹은 분할된 이분산성의 경우 모두 이분산성을 확인하는데 이용될 수 있음

- 비례적 이분산성의 경우, 먼저 어떤 변수가 오차항의 분산에 비례적인가를 확인하는 것이 필요함
- 그리고 나서 자료를 그 변수에 대해 큰 값에서 작은 값으로 sorting함

GQ 검정

- 통상적으로는 가운데 r 개의 관측치는 제외함 (for example, $r \approx T/6$),
- 그리고 처음 T_1 관측치들 그리고 마지막 T_2 관측치들에 대해 최소제곱 추정을 적용.
- 이 경우, 어느 그룹의 관측치들이 잠재적으로 큰 분산을 가지고 있는가를 알기 때문에 다음과 같은 단측 가설 검정을 하게 됨

$$H_0: \sigma_1^2 = \sigma_2^2 \quad H_1: \sigma_1^2 > \sigma_2^2$$

GQ 검정

**Goldfeld-Quandt
Test Statistic**

$$\mathbf{GQ} = \frac{\hat{\sigma}_1^2}{\hat{\sigma}_2^2} \sim F_{[T_1-K, T_2-K]}$$

Small values of **GQ** support H_0 while large values support H_1 .

두 관측치 그룹간에 분산이 다르다는 것이 의심되지
만 어느 쪽 분산이 잠재적으로 크지 모른다면 양측 가
설검정이 적절할 것임

GQ의 분모 분자도 바뀔 수 있으며, 그에 따른 적절한
기각역의 설정이 요구됨

Breusch-Pagan test

$$Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + \dots + \beta_K X_{Kt} + \varepsilon_t,$$

$$H_0 : \text{homoscedasticity} \quad V(\varepsilon_t) = \sigma^2$$

$$H_1 : \text{heteroscedasticity} \quad V(\varepsilon_t) = \sigma_t^2 = E(\varepsilon_t^2 | Z_t) = h(\delta_0 + \delta_1 Z_{1t} + \dots + \delta_S Z_{St})$$

$Z_t = (Z_{1t}, \dots, Z_{St})$: 설명변수 들 중 분산이 의존하는 변수들

$h(\cdot)$: Some smooth function

오차항의 분산이 설명변수들 일부에 의존하는 것이 의심될 때,
해당 설명변수들을 특정하여 수행할 수 있음

Breusch-Pagan test**Test procedures:**

(1) Run OLS on regression: $Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + \dots + \beta_K X_{Kt} + \varepsilon_t$,
obtain the residuals.

(2) Run the auxiliary regression (보조적 회귀):

$$\hat{\varepsilon}_t^2 = \delta_0 + \delta_1 Z_{1t} + \dots + \delta_S Z_{St} + v_t$$

$$H_0: \delta_1 = \dots = \delta_S = 0$$

(3) Compute $W = T \cdot R^2 \sim_{asy} \chi^2_{df}$ (Lagrange Multiplier Test)

(4) Compare the W and $\chi^2_{\alpha, df}$ (α : significance level)

(where the df is # of slope coefs in the auxiliary regression: S

if $W > \chi^2_{\alpha, df} \implies$ reject the H_0

White 이분산성 검정(교차항 있음)-참고

H_0 : homoscedasticity $\text{Var}(\varepsilon_t) = \sigma^2$

H_1 : heteroscedasticity $\text{Var}(\varepsilon_t) = \sigma_t^2$

Test procedures:

(1) Run OLS on regression: $Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + \varepsilon_t$,
obtain the residuals, $\hat{\varepsilon}_t$

(2) Run the auxiliary regression:

$$\hat{\varepsilon}_t^2 = \alpha_1 + \alpha_2 X_{2t} + \alpha_3 X_{3t} + \alpha_4 X_{2t}^2 + \alpha_5 X_{3t}^2 + \alpha_6 X_{2t} X_{3t} + v_t$$

(3) Compute W (or LM) = $T \cdot R^2 \sim_{asy} \chi^2_{df}$ (F검정도 이용가능)

(4) Compare the W and $\chi^2_{\alpha, df}$

(where the df is # of slope coefs in eqn (2) = 5)

if $W > \chi^2_{df} \implies$ reject the H_0

White 이분산성 검정(교차항 없음)- 참고

H_0 : homoscedasticity $\text{Var}(\varepsilon_t) = \sigma^2$

H_1 : heteroscedasticity $\text{Var}(\varepsilon_t) = \sigma_t^2$

Test procedures:

(1) Run OLS on regression: $Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + \dots + \beta_K X_{Kt} + \varepsilon_t$,
obtain the residuals, $\hat{\varepsilon}_i$

(2) Run the auxiliary regression (보조적 회귀):

$$\hat{\varepsilon}_t^2 = \delta_1 + \delta_2 X_{2t} + \dots + \delta_K X_{Kt} + \delta_{K+1} X_{2t}^2 + \dots + \delta_{2K-1} X_{Kt}^2 + v_t$$

$$H_0: \delta_2 = \delta_3 = \dots = \delta_{2K-1} = 0$$

(3) Compute $\mathbf{W} = \mathbf{T} \cdot \mathbf{R}^2 \underset{asy}{\sim} \chi^2_{df}$ (F검정도 이용가능)

(4) Compare the \mathbf{W} and $\chi^2_{\alpha, df}$ (α : significance level)

(where the **df** is # of slope coefs in (2) = $2K-2$)

if $\mathbf{W} > \chi^2_{df} \implies$ reject the H_0