

# Lecture 4: Heteroskedasticity and Dependence

OLS estimator:

$$\hat{\beta}_1 = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

$$\hat{\beta}_1 = \beta_1 + \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) u_i$$

$E[u_i] = 0 \rightarrow \hat{\beta}_1$  unbiased  
 $\rightarrow 0 \rightarrow \hat{\beta}_1$  consistent

$$v[\hat{\beta}_1] = \cancel{v[\beta_1]} + v\left[\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) u_i\right]$$

$$v[\hat{\beta}_1] = \left(\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)^2 \left[ \sum_{i=1}^n (x_i - \bar{x})^2 v[u_i] + \sum_{i=1}^n \sum_{j \neq i} (x_i - \bar{x})(x_j - \bar{x}) \text{cov}(u_i, u_j) \right]$$

# Classical model

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

- GM assumptions
- ①  $E[U_i] = 0$
  - ②  $V[U_i] = \sigma^2$
  - ③  $\text{cov}(U_i, U_j) = 0 \quad i \neq j$

$$V[\hat{\beta}_1] = \left( \frac{1}{\sum_i (X_i - \bar{X})^2} \right)^2 \left[ \sum_i (X_i - \bar{X})^2 \overset{\sigma^2}{V[U_i]} + \sum_i \sum_{j \neq i} \cancel{(X_i - \bar{X})(X_j - \bar{X}) \text{cov}(U_i, U_j)} \right]$$

homoskedasticity  $\rightarrow$  simple math!

no assumptions about shape of  $U_i$  distribution

$\rightarrow SE[\hat{\beta}_1]$  and then apply CLT

$\rightarrow$  large sample inference!

## Normal linear model

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

$$U_i \sim N(0, \sigma^2)$$

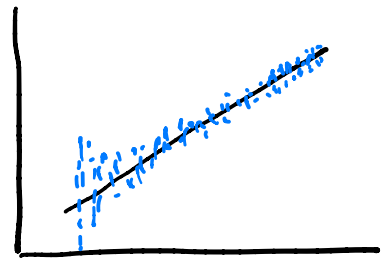
→ then  $t = \frac{\hat{\beta}_1 - \beta_1^0}{SE(\hat{\beta}_1)}$  has  $t(N-2)$  dist.

## Random X's

$$E[U_i] = 0 \Rightarrow E[U_i | X_1, X_2, \dots, X_N] = 0$$

$$E[\hat{\beta}_1] = \beta_1 \Rightarrow E[\hat{\beta}_1 | X_1, \dots, X_N] = \beta_1$$

# Heteroskedasticity



①  $E[U_i | X_i] = 0$

②  $(X_i, Y_i)$  are iid

③ outliers are unlikely

$$V[\hat{\beta}] = \left( \frac{1}{\sum_i (x_i - \bar{x})^2} \right)^2 \left[ \sum_i (x_i - \bar{x})^2 V[U_i] + \cancel{\sum_i \sum_{j \neq i} (x_i - \bar{x})(x_j - \bar{x}) \text{cov}(U_i, U_j)} \right]$$

→ still set  $\text{cov}(\ ) = 0$ , but more complicated formula

## Dependence

→ observations are dependent within clusters, but incl. across clusters.

- ① clustered sample design
- ② group-level treatment

$$v[\hat{\beta}_1] = \left( \frac{1}{\sum_i (x_i - \bar{x})^2} \right)^2 \left[ \sum_i (x_i - \bar{x})^2 v(u_i) + \sum_i \sum_{j \neq i} (x_i - \bar{x})(x_j - \bar{x}) \text{cov}(u_i, u_j) \right]$$

↑  $u_i^2$   
↑ still allow hetero  
↑  $u_i u_j$   
↑  $j \neq i$ 's cluster

$v_{\text{cov}} = \text{'hetero'}$

$v_{\text{cov}} = \sim \text{cluster error}$

→ Another option: group data  $\bar{Y}_K = \beta_0 + \beta_1 X_K + \bar{U}_K$

## Back WLS

- Recall:  $\hat{\beta}_i^{WLS} = \frac{\sum_i w_i (y_i - \bar{y})(x_i - \bar{x})}{\sum_i w_i (x_i - \bar{x})^2}$

- If  $U_i$  is heteroskedastic,  $w_i = \frac{1}{V[U_i]}$

- Grouped data example:

individual:

$$y_{ig} = \beta_0 + \beta_1 x_{ig} + U_{ig}$$

$$V[U_{ig} | x_{ig}] = \sigma^2$$

group:

$$\bar{y}_g = \beta_0 + \beta_1 \bar{x}_g + \bar{U}_g$$

$$V[\bar{U}_g | \bar{x}_g] = \frac{\sigma^2}{N_g}$$

$$- w_g = \frac{1}{V[\bar{U}_g]} = \frac{N_g}{\sigma^2}$$

← constant → can set  $w_g = N_g$