

BŁĄD ZMIENNEJ POMINIĘTEJ

EKONOMETRIA WNE

Sebastian Zalas

University of Warsaw

s.zalas@uw.edu.pl

BŁĄD ZMIENNEJ POMINIĘTEJ

- ▶ Jesteśmy zainteresowani oszacowaniem jak $x \longrightarrow y$

$$y = \beta_0 + \beta_1 x + \varepsilon = \mathbf{X}\beta + \varepsilon$$

- ▶ Prawdziwy model:

$$y = \beta_0 + \beta_1 x + \gamma z + \varepsilon = \mathbf{X}\beta + \mathbf{z}\gamma + \varepsilon \quad (1)$$

z - pominięta zmienna; niech z będzie skorelowane z x

- ▶ Co stanie się z $\hat{\beta}$?

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \mathbf{z}\gamma + \varepsilon) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\gamma + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon \\ &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\gamma + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon\end{aligned}$$

BŁĄD ZMIENNEJ POMINIĘTEJ

- ▶ Jesteśmy zainteresowani oszacowaniem jak $x \longrightarrow y$

$$y = \beta_0 + \beta_1 x + \varepsilon = \mathbf{X}\beta + \varepsilon$$

- ▶ Prawdziwy model:

$$y = \beta_0 + \beta_1 x + \gamma z + \varepsilon = \mathbf{X}\beta + \mathbf{z}\gamma + \varepsilon \quad (1)$$

z - pominięta zmienna; niech z będzie skorelowane z x

- ▶ Co stanie się z $\hat{\beta}$?

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \mathbf{z}\gamma + \varepsilon) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\gamma + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon \\ &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\gamma + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon\end{aligned}$$

BŁĄD ZMIENNEJ POMINIĘTEJ

- obciążenie:

$$\begin{aligned}\mathbb{E}[\hat{\beta} | \mathbf{X}] - \beta &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\gamma | \mathbf{X}] + \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon | \mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\gamma \mathbb{E}[\mathbf{z} | \mathbf{X}] + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \underbrace{\mathbb{E}[\varepsilon | \mathbf{X}]}_{=0} \\ &= \gamma(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbb{E}[\mathbf{z} | \mathbf{X}] \neq 0\end{aligned}\tag{2}$$

- Szacując model bez $\gamma\mathbf{z}$ pozostanie on w składniku losowym:

$$y = \beta_0 + \beta_1 x + \underbrace{(\gamma\mathbf{z} + \varepsilon)}_{=v}$$

- wtedy $\mathbb{E}[v | \mathbf{X}] = \mathbb{E}[\gamma\mathbf{z} | \mathbf{X}] + \mathbb{E}[\varepsilon | \mathbf{X}] = \mathbb{E}[\gamma\mathbf{z} | \mathbf{X}] \neq 0 \longrightarrow$ złamanie założenia KMRL prowadzi do **obciążenia** estymatora MNK

BŁĄD ZMIENNEJ POMINIĘTEJ

- obciążenie:

$$\begin{aligned}\mathbb{E}[\hat{\beta} | \mathbf{X}] - \beta &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\gamma | \mathbf{X}] + \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon | \mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\gamma \mathbb{E}[\mathbf{z} | \mathbf{X}] + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \underbrace{\mathbb{E}[\varepsilon | \mathbf{X}]}_{=0} \\ &= \gamma(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbb{E}[\mathbf{z} | \mathbf{X}] \neq 0\end{aligned}\tag{2}$$

- Szacując model bez $\gamma\mathbf{z}$ pozostanie on w składniku losowym:

$$y = \beta_0 + \beta_1 x + \underbrace{(\gamma\mathbf{z} + \varepsilon)}_{=v}$$

- wtedy $\mathbb{E}[v | \mathbf{X}] = \mathbb{E}[\gamma\mathbf{z} | \mathbf{X}] + \mathbb{E}[\varepsilon | \mathbf{X}] = \mathbb{E}[\gamma\mathbf{z} | \mathbf{X}] \neq 0 \longrightarrow$ złamanie założenia KMRL prowadzi do **obciążenia** estymatora MNK

BŁĄD ZMIENNEJ POMINIĘTEJ

- obciążenie:

$$\begin{aligned}\mathbb{E}[\hat{\beta} | \mathbf{X}] - \beta &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\gamma | \mathbf{X}] + \underbrace{\mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon | \mathbf{X}]}_{=0} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\gamma \mathbb{E}[\mathbf{z} | \mathbf{X}] + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \underbrace{\mathbb{E}[\varepsilon | \mathbf{X}]}_{=0} \\ &= \gamma(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbb{E}[\mathbf{z} | \mathbf{X}] \neq 0\end{aligned}\tag{2}$$

- Szacując model bez $\gamma\mathbf{z}$ pozostanie on w składniku losowym:

$$y = \beta_0 + \beta_1 x + \underbrace{(\gamma\mathbf{z} + \varepsilon)}_{=v}$$

- wtedy $\mathbb{E}[v | \mathbf{X}] = \mathbb{E}[\gamma\mathbf{z} | \mathbf{X}] + \mathbb{E}[\varepsilon | \mathbf{X}] = \mathbb{E}[\gamma\mathbf{z} | \mathbf{X}] \neq 0 \longrightarrow$ złamanie założenia KMRL prowadzi do **obciążenia** estymatora MNK

BŁĄD ZMIENNEJ POMINIĘTEJ

- obciążenie:

$$\begin{aligned}\mathbb{E}[\hat{\beta} | \mathbf{X}] - \beta &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\gamma | \mathbf{X}] + \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon | \mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\gamma \mathbb{E}[\mathbf{z} | \mathbf{X}] + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \underbrace{\mathbb{E}[\varepsilon | \mathbf{X}]}_{=0} \\ &= \gamma(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbb{E}[\mathbf{z} | \mathbf{X}] \neq 0\end{aligned}\tag{2}$$

- Szacując model bez $\gamma\mathbf{z}$ pozostanie on w składniku losowym:

$$y = \beta_0 + \beta_1 x + \underbrace{(\gamma\mathbf{z} + \varepsilon)}_{=\nu}$$

- wtedy $\mathbb{E}[\nu | \mathbf{X}] = \mathbb{E}[\gamma\mathbf{z} | \mathbf{X}] + \mathbb{E}[\varepsilon | \mathbf{X}] = \mathbb{E}[\gamma\mathbf{z} | \mathbf{X}] \neq 0 \longrightarrow$ złamanie założenia KMRL prowadzi do **obciążenia** estymatora MNK

BŁĄD ZMIENNEJ POMINIĘTEJ

- Kiedy błąd zmiennej pominiętej zachodzi? Przekształćmy (2):

$$\begin{aligned}\mathbb{E}[\hat{\beta} | \mathbf{X}] - \beta &= \gamma \frac{\text{Cov}[z, x]}{\mathbb{V}[x]} \\ &= \gamma \frac{\sigma_z}{\sigma_x} \rho_{x,z}\end{aligned}\tag{3}$$

- x jest skorelowana z pominiętą zmienną z , czyli

$$\rho_{x,z} \neq 0$$

- pominięta zmienna z jest skorelowana ze zmienną zależną y , czyli:

$$\gamma \neq 0$$

BŁĄD ZMIENNEJ POMINIĘTEJ

- Kiedy błąd zmiennej pominiętej zachodzi? Przekształćmy (2):

$$\begin{aligned}\mathbb{E}[\hat{\beta} | \mathbf{X}] - \beta &= \gamma \frac{\text{Cov}[z, x]}{\mathbb{V}[x]} \\ &= \gamma \frac{\sigma_z}{\sigma_x} \rho_{x,z}\end{aligned}\tag{3}$$

- x jest skorelowana z pominiętą zmienną z , czyli

$$\rho_{x,z} \neq 0$$

- pominięta zmienna z jest skorelowana ze zmienną zależną y , czyli:

$$\gamma \neq 0$$

BŁĄD ZMIENNEJ POMINIĘTEJ

- Kiedy błąd zmiennej pominiętej zachodzi? Przekształćmy (2):

$$\begin{aligned}\mathbb{E}[\hat{\beta} | \mathbf{X}] - \beta &= \gamma \frac{\text{Cov}[z, x]}{\mathbb{V}[x]} \\ &= \gamma \frac{\sigma_z}{\sigma_x} \rho_{x,z}\end{aligned}\tag{3}$$

- x jest skorelowana z pominiętą zmienną z , czyli

$$\rho_{x,z} \neq 0$$

- pominięta zmienna z jest skorelowana ze zmienną zależną y , czyli:

$$\gamma \neq 0$$

BŁĄD ZMIENNEJ POMINIĘTEJ - WNIOSKI

- ▶ Błąd zmiennej pominiętej zachodzi gdy oba warunki są spełnione (via (3)):
 1. zmienna w modelu jest skorelowana z pominiętą zmienną
 2. pominięta zmienna jest skorelowana ze zmienną zależną
- ▶ Pominięcie zmiennej w modelu prowadzi do obciążenia estymatora MNK z powodu złamania założenia KMRL: $\mathbb{E}[\varepsilon \mid \mathbf{X}] = 0$

BŁĄD ZMIENNEJ POMINIĘTEJ - WNIOSKI

- ▶ Błąd zmiennej pominiętej zachodzi gdy oba warunki są spełnione (via (3)):
 1. zmienna w modelu jest skorelowana z pominiętą zmienną
 2. pominięta zmienna jest skorelowana ze zmienną zależną
- ▶ Pominięcie zmiennej w modelu prowadzi do obciążenia estymatora MNK z powodu złamania założenia KMRL: $\mathbb{E}[\varepsilon \mid \mathbf{X}] = 0$

BŁĄD ZMIENNEJ POMINIĘTEJ - PRZYKŁAD

- ▶ wpływ *profocc* na płace
 - *profocc* = 1 gdy osoba pracuje w wysoko wyspecjalizowanym zawodzie

- ▶ oszacujmy dwa modele:

$$\log(\text{wage}) = \text{profocc} + \text{exper} + \text{expersq} + \text{female} + \varepsilon$$

$$\log(\text{wage}) = \text{profocc} + \text{educ} + \text{exper} + \text{expersq} + \text{female} + \varepsilon$$

- ▶ Czy edukacja może być pominięta? Czy są inne możliwości?

BŁĄD ZMIENNEJ POMINIĘTEJ - PRZYKŁAD

- ▶ wpływ *profocc* na płace
 - *profocc* = 1 gdy osoba pracuje w wysoko wyspecjalizowanym zawodzie
- ▶ oszacujmy dwa modele:

$$\log(\text{wage}) = \text{profocc} + \text{exper} + \text{expersq} + \text{female} + \varepsilon$$

$$\log(\text{wage}) = \text{profocc} + \text{educ} + \text{exper} + \text{expersq} + \text{female} + \varepsilon$$

- ▶ Czy edukacja może być pominięta? Czy są inne możliwości?

BŁĄD ZMIENNEJ POMINIĘTEJ - PRZYKŁAD

- ▶ wpływ *profocc* na płace
 - *profocc* = 1 gdy osoba pracuje w wysoko wyspecjalizowanym zawodzie
- ▶ oszacujmy dwa modele:

$$\log(\text{wage}) = \text{profocc} + \text{exper} + \text{expersq} + \text{female} + \varepsilon$$

$$\log(\text{wage}) = \text{profocc} + \text{educ} + \text{exper} + \text{expersq} + \text{female} + \varepsilon$$

- ▶ Czy edukacja może być pominięta? Czy są inne możliwości?

BŁĄD ZMIENNEJ POMINIĘTEJ - PRZYKŁAD

	<i>Zmienna zależna:</i>	
	<i>log(wage)</i>	
	(1)	(2)
profocc	0.402*** (0.040)	0.227*** (0.043)
educ		0.063*** (0.008)
exper	0.038*** (0.005)	0.037*** (0.005)
expersq	-0.001*** (0.0001)	-0.001*** (0.0001)
female	-0.312*** (0.038)	-0.310*** (0.036)
N	526	526
R ²	0.359	0.430
Note:	*p<0.1; **p<0.05; ***p<0.01	

BŁĄD ZMIENNEJ POMINIĘTEJ - PRZYKŁAD

Edukacja jest pominiętym czynnikiem:

- ▶ edukacja jest skorelowana z płacą
- ▶ *profocc* jest skorelowane z edukacją

Współczynnik przy *profocc* w kol. (1):

- ▶ jest obciążony (przeszacowany):

$$\hat{\beta}_{profocc} - \beta_{profocc} = \beta_{educ} \frac{\sigma_{educ}}{\sigma_{profocc}} \rho_{educ,profocc} > 0$$

$$\hat{\beta}_{profocc} > \beta_{profocc}$$

- ▶ mierzy wpływ pracy w wysoko wyspecjalizowanym zawodzie oraz edukacji

BŁĄD ZMIENNEJ POMINIĘTEJ - PRZYKŁAD

Edukacja jest pominiętym czynnikiem:

- ▶ edukacja jest skorelowana z płacą
- ▶ *profocc* jest skorelowane z edukacją

Współczynnik przy *profocc* w kol. (1):

- ▶ jest obciążony (przeszacowany):

$$\hat{\beta}_{profocc} - \beta_{profocc} = \beta_{educ} \frac{\sigma_{educ}}{\sigma_{profocc}} \rho_{educ,profocc} > 0$$

$$\hat{\beta}_{profocc} > \beta_{profocc}$$

- ▶ mierzy wpływ pracy w wysoko wyspecjalizowanym zawodzie oraz edukacji

BŁĄD ZMIENNEJ POMINIĘTEJ - CO ZROBIĆ?

Jak go wykryć?

- ▶ nie da się przetestować czy pominięto ważną zmienną
- ▶ należy bazować na własnej ocenie/ teorii/ literaturze

Co robić?

- ▶ warto wyznaczyć, który współczynnik jest interesujący
- ▶ włączyć do modelu zmienną *proxy* (zastępczą) dla pominiętej
- ▶ można wyznaczyć kierunek obciążenia korzystając z (3) oraz ocenić jego siłę

BŁĄD ZMIENNEJ POMINIĘTEJ - CO ZROBIĆ?

Jak go wykryć?

- ▶ nie da się przetestować czy pominięto ważną zmienną
- ▶ należy bazować na własnej ocenie/ teorii/ literaturze

Co robić?

- ▶ warto wyznaczyć, który współczynnik jest interesujący
- ▶ włączyć do modelu zmienną *proxy* (zastępczą) dla pominiętej
- ▶ można wyznaczyć kierunek obciążenia korzystając z (3) oraz ocenić jego siłę

BŁĄD ZMIENNEJ POMINIĘTEJ - CO ZROBIĆ?

Przykład zmiennej *proxy* zastępczej:

- ▶ cena samochodu zależy od wieku
- ▶ wieku samochodu nie obserwujemy
- ▶ może obserwujemy jak długo posiada go obecny właściciel ← zmienna proxy

Pytania? Wątpliwości?
Dziękuję!

e: s.zalas@uw.edu.pl