

# STATYSTYCZNE WŁASNOŚCI ESTYMATORA MNK & HETEROSKEDASTYCZNOŚĆ

EKONOMETRIA WNE

Sebastian Zalas

University of Warsaw

[s.zalas@uw.edu.pl](mailto:s.zalas@uw.edu.pl)

# KLASYCZNY MODEL REGRESJI LINIOWEJ

1.  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  - model jest liniowy
2. Zmienne losowe  $\{(y_1, x_1), \dots, (y_i, x_i), \dots, (y_n, x_n)\}$  są niezależne oraz wylosowane z tego samego rozkładu (*independently and identically distributed - iid.*)
3.  $\text{rz}[\mathbf{X}_{n \times k}] = k$  - rząd kolumnowy  $\mathbf{X}$  jest pełny
4.  $\mathbb{E}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbf{0}$  - wartość oczekiwana składnika losowego jest równa 0
5.  $\mathbb{E}[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'|\mathbf{X}] = \mathbf{I}\sigma^2$  - sferyczność wariancji
6.  $\boldsymbol{\varepsilon}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}\sigma^2)$  składnik losowy ma rozkład normalny

# WŁASNOŚCI ESTYMATORA

- ▶ Estymator  $\hat{\theta}$  nieznanego parametru  $\theta$  jest funkcją danych, a więc jest on **zmienną losową**.

$$\hat{\theta} = g(X_1, \dots, X_n)$$

- ▶ Istnieje zatem **wartość oczekiwana** estymatora -  $\mathbb{E}[\hat{\theta}]$  oraz jego **wariancja** -  $\mathbb{V}[\hat{\theta}]$  (*sampling variance*)

- ▶ Estymator  $\hat{\theta}$  parametru  $\theta$  nazywamy **nieobciążonym** gdy:

$$\mathbb{E}[\hat{\theta}] = \theta$$

- ▶ **Obciążenie** estymatora

$$bias(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta$$

# WŁASNOŚCI ESTYMATORA

- ▶ Estymator  $\hat{\theta}$  nieznanego parametru  $\theta$  jest funkcją danych, a więc jest on **zmienną losową**.

$$\hat{\theta} = g(X_1, \dots, X_n)$$

- ▶ Istnieje zatem **wartość oczekiwana** estymatora -  $\mathbb{E}[\hat{\theta}]$  oraz jego **wariancja** -  $\mathbb{V}[\hat{\theta}]$  (*sampling variance*)
- ▶ Estymator  $\hat{\theta}$  parametru  $\theta$  nazywamy **nieobciążonym** gdy:

$$\mathbb{E}[\hat{\theta}] = \theta$$

- ▶ **Obciążenie** estymatora

$$bias(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta$$

# NIEOBCIĄŻONOŚĆ

- Założenia (1) oraz (4):

$$\begin{aligned}\mathbb{E}[\varepsilon|\mathbf{X}] &= 0 \\ &= \mathbb{E}[\mathbf{y} - \mathbf{X}\beta|\mathbf{X}] \\ &= \mathbb{E}[\mathbf{y}|\mathbf{X}] - \mathbb{E}[\mathbf{X}\beta|\mathbf{X}] \\ &= \mathbb{E}[\mathbf{y}|\mathbf{X}] - \mathbf{X}\beta \\ \Rightarrow \mathbb{E}[\mathbf{y}|\mathbf{X}] &= \mathbf{X}\beta\end{aligned}$$

# NIEOBCIĄŻONOŚĆ

- Estymator uzyskany MNK  $\hat{\beta}$  wektora parametrów  $\beta$ :

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

- Pokażmy, że estymator  $\hat{\beta}$  jest nieobciążony, pw. że założenia (1) - (4) są spełnione:

$$\begin{aligned}\mathbb{E}[\hat{\beta}|\mathbf{X}] &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}|\mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \underbrace{\mathbb{E}[\mathbf{y}|\mathbf{X}]}_{\mathbf{X}\beta} \\ &= \underbrace{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}}_{=I} \beta \\ &= \beta\end{aligned}$$

# NIEOBCIĄŻONOŚĆ

- Estymator uzyskany MNK  $\hat{\beta}$  wektora parametrów  $\beta$ :

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

- Pokażmy, że estymator  $\hat{\beta}$  jest nieobciążony, pw. że założenia (1) - (4) są spełnione:

$$\begin{aligned}\mathbb{E}[\hat{\beta}|\mathbf{X}] &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}|\mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \underbrace{\mathbb{E}[\mathbf{y}|\mathbf{X}]}_{\mathbf{X}\beta} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X} \underbrace{\beta}_{=I} \\ &= \beta\end{aligned}$$

# NIEOBCIĄŻONOŚĆ

- Zdekomponujmy estymator MNK  $\hat{\beta}$ :

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

- Pokażmy nieobciążony korzystając z def. obciążenia:

$$\begin{aligned}\mathbb{E}[\hat{\beta} - \beta | \mathbf{X}] &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon | \mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbb{E}[\varepsilon | \mathbf{X}] = 0\end{aligned}$$



# NIEOBCIĄŻONOŚĆ

- Zdekomponujmy estymator MNK  $\hat{\beta}$ :

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

- Pokażmy nieobciążony korzystając z def. obciążenia:

$$\begin{aligned}\mathbb{E}[\hat{\beta} - \beta | \mathbf{X}] &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon | \mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbb{E}[\varepsilon | \mathbf{X}] = 0\end{aligned}$$

# NIEOBCIĄŻONOŚĆ

## **Twierdzenie. Nieobciążoność estymatora MNK.**

W Klasycznym Modelu Regresji Liniowej (założenia (1) - (4)), estymator uzyskany MNK jest nieobciążony:

$$\mathbb{E}[\hat{\beta} \mid \mathbf{X}] = \beta$$

- ▶  $\mathbb{E}[\varepsilon \mid \mathbf{X}] = 0 \Rightarrow$  nieuwzględnienie ważnego czynnika w modelu  $\Rightarrow$  jest obecny w  $\varepsilon \Rightarrow$  oszacowania będą obciążone
- ▶ Rozkład warunkowy względem  $\mathbf{X}$  - estymator jest nieobciążony dla każdej realizacji macierzy regresorów  $\mathbf{X}$
- ▶ Warunkowy rozkład  $\hat{\beta}$  jest skoncentrowany wokół  $\beta$

# NIEOBCIĄŻONOŚĆ

## **Twierdzenie. Nieobciążoność estymatora MNK.**

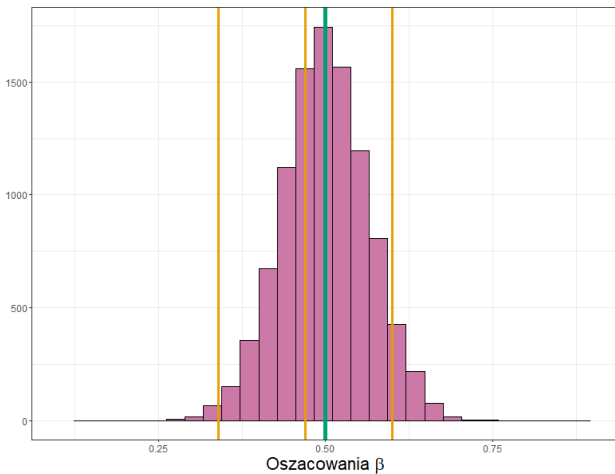
W Klasycznym Modelu Regresji Liniowej (założenia (1) - (4)), estymator uzyskany MNK jest nieobciążony:

$$\mathbb{E}[\hat{\beta} \mid \mathbf{X}] = \beta$$

- ▶  $\mathbb{E}[\varepsilon \mid \mathbf{X}] = 0 \Rightarrow$  nieuwzględnienie ważnego czynnika w modelu  $\Rightarrow$  jest obecny w  $\varepsilon \Rightarrow$  oszacowania będą obciążone
- ▶ Rozkład warunkowy względem  $\mathbf{X}$  - estymator jest nieobciążony dla każdej realizacji macierzy regresorów  $\mathbf{X}$
- ▶ Warunkowy rozkład  $\hat{\beta}$  jest skoncentrowany wokół  $\beta$

# NIEOBCIĄŻONOŚĆ: $\hat{\beta}$ vs. $\beta$

Symulacja modelu  $y = 1 + 0.5x + \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, 4)$



# WARIANCJA ESTYMATORA MNK

- Definicja wariancji. Niech  $\mathbf{x}$  będzie wektorem losowym:

$$\mathbb{V}[\mathbf{X}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])']$$

- Wariancja warunkowa:

$$\mathbb{V}[\mathbf{X} \mid \mathbf{Z}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])' \mid \mathbf{Z}]$$

- Wariancja składnika losowego:

$$\mathbb{V}[\varepsilon \mid \mathbf{X}] = \mathbb{E}[\varepsilon \varepsilon' \mid \mathbf{X}]$$

# WARIANCJA ESTYMATORA MNK

- ▶ Przywołajmy fakt:  $\hat{\beta} = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon$
- ▶ Wyprowadzimy formułę wariancji estymatora MNK, czyli  $\mathbb{V}[\hat{\beta}]$ :

$$\begin{aligned}\mathbb{V}[\hat{\beta} \mid \mathbf{X}] &= \mathbb{E}[(\beta - \hat{\beta})(\beta - \hat{\beta})' \mid \mathbf{X}] \\&= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon)(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon)' \mid \mathbf{X}] \\&= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon\varepsilon'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mid \mathbf{X}] \\&= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbb{E}[\varepsilon\varepsilon' \mid \mathbf{X}]\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\end{aligned}$$

- ▶ wariancja estymatora  $\hat{\beta}$  zależy od  $\mathbb{E}[\varepsilon\varepsilon' \mid \mathbf{X}] = \mathbb{V}[\varepsilon \mid \mathbf{X}]$

# WARIANCJA ESTYMATORA MNK

- W KMRL zakładamy sferyczność wariancji -  $\mathbb{E}[\varepsilon\varepsilon'|\mathbf{X}] = \mathbf{I}\sigma^2$ , czyli:

$$\mathbb{V}[\varepsilon | \mathbf{X}] = \mathbb{E}[\varepsilon\varepsilon' | \mathbf{X}] = \mathbf{\Omega} = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \mathbf{I}\sigma^2$$

- W takim przypadku, wariancja **estymatora MNK** przyjmuje postać:

$$\mathbb{V}[\hat{\beta} | \mathbf{X}] = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\sigma^2\mathbf{I}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$$

# WARIANCJA ESTYMATORA MNK

- Sferyczność wariancji - co to oznacza?

$$\mathbb{V}[\boldsymbol{\varepsilon} \mid \mathbf{X}] = \mathbb{E}[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}' \mid \mathbf{X}] = \boldsymbol{\Omega} = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \mathbf{I}\sigma^2$$

- **homoskedastyczność**  $\Rightarrow$  stałość wariancji  $\Rightarrow$  takie same elementy na diagonalu macierzy  $\boldsymbol{\Omega}$
- brak autokorelacji  $\Rightarrow$  zera poza diagonalą macierzy  $\boldsymbol{\Omega}$



## Twierdzenie Gaussa Markova.

W Klasycznym Modelu Regresji Liniowej (założenia (1) - (5)), nieobciążony estymator uzyskany MNK ma najniższą wariancję spośród liniowych, nieobciążonych estymatorów.

$$\mathbb{V}[\tilde{\beta}^1 | \mathbf{X}] \geq \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$$

---

<sup>1</sup> oznacza dowolny liniowy nieobciążony estymator

- ▶ Żaden inny nieobciążony estymator nie może mieć niższej wariancji niż  $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$
- ▶ Estymator MNK jest najefektywniejszy w klasie liniowych nieobciążonych estymatorów - **BLUE** - *Best Linear Unbiased Estimator*

# WARIANCJA $\varepsilon$

►  $\mathbb{V}[\varepsilon] = \sigma^2$  - nie znamy  $\sigma^2 \Rightarrow$  nie znamy  $\mathbb{V}[\hat{\beta}]$

► Należy oszacować  $\sigma^2$ :

$$s^2 = \frac{1}{n-k} \sum_{i=1}^n e_i^2 = \frac{\mathbf{e}'\mathbf{e}}{n-k}$$

nieobciążony estymator  $\sigma^2$

► Estymator wariancji  $\hat{\beta}$ :

$$\widehat{\mathbb{V}[\hat{\beta} \mid \mathbf{X}]} = s^2(\mathbf{X}'\mathbf{X})^{-1} = \frac{\mathbf{e}'\mathbf{e}}{n-k}(\mathbf{X}'\mathbf{X})^{-1} = \hat{\mathbf{V}}_{\hat{\beta}}$$

# HETEROSKEDASTYCZNOŚĆ $\varepsilon$

- Analizowane dane mogą jednak nie spełniać założenia o sferyczności składnika losowego. Wtedy jego wariancja ma postać:

$$\mathbb{V}[\varepsilon \mid \mathbf{X}] = \mathbb{E}[\varepsilon \varepsilon' \mid \mathbf{X}] = \mathbf{\Omega} = \begin{bmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n^2 \end{bmatrix}$$

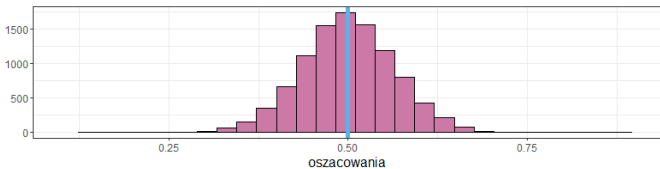
- heteroskedastyczność  $\Rightarrow$  brak stałości wariancji  $\Rightarrow$  różne elementy na diagonalu macierzy  $\mathbf{\Omega}$
  - utrzymujemy założenia o braku autokorelacji  $\Rightarrow$  zera poza diagonalą macierzy  $\mathbf{\Omega}$
- W takim przypadku, wariancja **estymatora MNK** przyjmuje postać:

$$\mathbb{V}[\hat{\beta} \mid \mathbf{X}] = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{\Omega} \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}$$

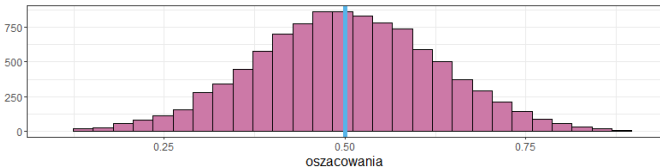
# HETEROSKEDASTYCZNOŚĆ $\varepsilon$ - KONSEKWENCJE

- ▶ Nieobciążoność  $\hat{\beta}$  pozostaje nienaruszona.
- ▶ Jednak estymator  $\hat{\beta}$  będzie **nieefektywny** tzn. możemy znaleźć estymator o niższej wariancji w klasie liniowych estymatorów
- ▶ Oznacza to, że precyzja estymatora się zmniejsza, co wyływa także na jakość wnioskowania statystycznego z oszacowanego modelu.

# HOMO- VS HETERO- SKEDASTYCZNOŚĆ



(a) Rozkład  $\hat{\beta}$  w modelu z **homoskedastycznym** składnikiem losowym.



(b) Rozkład  $\hat{\beta}$  w modelu z **heteroskedastycznym** składnikiem losowym.

# HETEROSKEDASTYCZNOŚĆ - ESTYMATOR ODPORNY

- ▶ Gdy występuje heteroskedastyczności, zwykły estymator wariancji może być obciążony  $\Rightarrow$  musimy skonstruować taki estymator wariancji, który będzie odporny na heteroskedastyczność

- ▶ Idealnie byłoby, gdybyśmy mogli mieć:

$$\begin{aligned}\hat{\mathbf{V}}_{\hat{\beta}}^{ideal} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbb{E}[\varepsilon\varepsilon' \mid \mathbf{X}]\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \begin{bmatrix} \varepsilon_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \varepsilon_n^2 \end{bmatrix} \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\end{aligned}$$

- ▶ Otrzymanie takiego estymatora nie jest możliwe...

# HETEROSKEDASTYCZNOŚĆ - ESTYMATOR WHITE'A

- ▶ White (1980) pokazał, że poniższy estymator wariancji estymatora MNK:

$$\hat{\mathbf{V}}_{\hat{\beta}}^{HC0} = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n x_i x_i' e_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1}$$

jest zgodnym estymatorem  $\mathbb{V}[\hat{\beta} \mid \mathbf{X}]$ , przez co jest odporny na heteroskedastyczność.

- ▶ HC - *Heteroscedasticity consistent*, czasami mówi się *heteroscedasticity robust*

# DLACZEGO $\mathbb{V}[\varepsilon \mid \mathbf{X}]$ JEST TAK WAŻNA?

- ▶ Wariancja składnika losowego,  $\mathbb{V}[\varepsilon]$  decyduje o wariancji estymatora,  $\mathbb{V}[\hat{\beta}]$
- ▶ Błąd standardowy  $\hat{\beta}$ :

$$se[\hat{\beta}] = \sqrt{\mathbb{V}[\hat{\beta}]} = \begin{bmatrix} \sqrt{\widehat{\mathbb{V}[\hat{\beta}_0]}} & \dots & \sqrt{\widehat{\text{Cov}[\hat{\beta}_0, \hat{\beta}_k]}} \\ \vdots & \ddots & \vdots \\ \sqrt{\widehat{\text{Cov}[\hat{\beta}_0, \hat{\beta}_k]}} & \dots & \sqrt{\widehat{\mathbb{V}[\hat{\beta}_k]}} \end{bmatrix}$$

- ▶ Na podstawie  $\hat{\beta}$  będziemy testować istotność oszacowań  $\Rightarrow$  niepoprawna wariancja  $\varepsilon$  prowadzi do niepoprawnego wnioskowania



# TEST BREUSCH'A - PAGAN'A

- Hipoteza zerowa

$$H_0 : \mathbb{E}[\varepsilon^2 \mid x_1, x_2, \dots, x_k] = \mathbb{E}[\varepsilon^2] = \sigma^2.$$

- Oszacuj model, uzyskaj kwadraty reszt,  $e^2$ :

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + e$$

- Oszacuj model, policz  $R_{e^2}^2$  z tej regresji:

$$e^2 = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \dots + \delta_k x_k + u$$

- Oblicz statystykę testową, uzyskaj  $p$ -value:

$$F = \frac{R_{e^2}^2 \frac{1}{k}}{(1 - R_{e^2}^2) \frac{1}{n-k-1}} \sim F_{k, n-k-1}$$

lub skorzystaj ze statystyki:

$$LM = n \times R_{\hat{u}^2}^2 \sim \chi_k^2$$

# TEST WHITE'A

- ▶ Hipoteza zerowa

$$H_0 : \mathbb{E}[\varepsilon^2 \mid x_1, x_2, \dots, x_k] = \mathbb{E}[\varepsilon^2] = \sigma^2.$$

- ▶ Oszacuj model, uzyskaj kwadraty reszt,  $e^2$ :

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + e$$

- ▶ Oszacuj model (z kwadratami i interakcjami), policz  $R_{e^2}^2$  z tej regresji:

$$e^2 = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \dots + \delta_k x_k + \sum_{i=1}^k \sum_{j=1}^k \delta_{ij} x_i x_j + u$$

- ▶ Oblicz statystykę i uzyskaj  $p$ -value tak jak w przypadku testu BP.

Pytania? Wątpliwości?  
Dziękuję!

**e:** [s.zalas@uw.edu.pl](mailto:s.zalas@uw.edu.pl)