

# Testowanie hipotez statystycznych

## Wykład 8

Natalia Nehrebecka Stanisław Cichocki

29 listopada 2015

# Plan zajęć

- 1 Rozkład estymatorów MNK w KMRL
  - Rozkład estymatora  $b$
  - Rozkład sumy kwadratów reszt
- 2 Testowanie hipotez prostych
  - Hipotezy proste - test t
  - Badanie istotności zmiennych w modelu
- 3 Przedziały ufności
  - dla parametrów
- 4 Testowanie hipotez łącznych
  - Hipotezy łączne - test F
- 5 Pytania teoretyczne

## Dodatkowe założenie

Oprócz założeń o:

- braku autokorelacji i homoskedastyczności:  $Var(\varepsilon) = \sigma^2 \mathbf{I}$
- zerowej wartości oczekiwanej:  $E(\varepsilon) = \mathbf{0}$

Dochodzi założenie o:

- normalności rozkładu błędów losowych.

Reasumując:

$$\varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

Wiemy już że:

- $\mathbf{b} = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon$
- $E(b) = \beta$  oraz  $\text{Var}(b) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$

Stąd:

$$b \sim N(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$$

Wiemy już, że:

- $\mathbf{e}'\mathbf{e} = \varepsilon'\mathbf{M}_X\varepsilon$
- macierz  $\mathbf{M}_X$  - symetryczna i idempotentna
- rząd macierzy  $\mathbf{M}_X = N - K$

Stąd:

$$\frac{\sum_{i=1}^N e_i^2}{\sigma^2} = \frac{\mathbf{e}'\mathbf{e}}{\sigma^2} = \frac{\varepsilon'\mathbf{M}_X\varepsilon}{\sigma^2} \sim \chi_{N-K}^2$$

## Brak korelacji między $b$ a $e$

$$\text{cov}(\mathbf{b}, \mathbf{e}) = 0$$

co implikuje, że:

$$\text{cov}(\mathbf{b}, \mathbf{e}'\mathbf{e}) = 0$$

# Testowanie hipotez

**Hipotezy proste** dotyczą pojedynczego parametru modelu lub kombinacji liniowej parametrów

## Rozkład statystyki t

$$t = \frac{b_k - \beta_k}{\hat{se}(b_k)} = \dots = \frac{\frac{b_k - \beta_k}{\sqrt{\sigma^2(\mathbf{X}'\mathbf{X})_{kk}^{-1}}}}{\sqrt{\frac{\mathbf{e}'\mathbf{e}}{N-K}}}$$

Ponieważ:  $\frac{b_k - \beta_k}{\sqrt{\sigma^2(\mathbf{X}'\mathbf{X})_{kk}^{-1}}} \sim N(0, 1)$  oraz  $\frac{\mathbf{e}'\mathbf{e}}{\sigma^2} \sim \chi_{N-K}^2$ , stad

$$t \sim t_{N-K}$$



## Przykład (1/2)

Założmy, że teoria mówi, że pewien parametr modelu,  $\beta_k$ , jest równy określonej wartości,  $\beta_k^*$ ,  $\beta_k = \beta_k^*$

Jeżeli:

- spełnione są założenia KMRL
- błąd losowy ma rozkład normalny
- teoria jest słuszna / hipoteza zerowa  $H_0$  jest prawdziwa

## Przykład (2/2)

Wtedy:

- statystyka testowa:

$$t = \frac{b_k - \beta_k^*}{\hat{se}(b_k)} \sim t_{N-K}$$

- statystyka krytyczna (odczytujemy z tablic rozkładu *t-Studenta*):

$$t^* = t \left( \underbrace{N - K}_{\text{Stopni swobody}}, \underbrace{1 - \frac{\alpha}{2}}_{\text{Rzad kwantyla}} \right)$$

gdzie:  $\alpha$ - poziom istotności

⇒ Jeśli  $|t| \geq t^*$  - odrzucamy  $H_0$

⇒ Jeśli  $|t| < t^*$  - nie ma podstaw do odrzucenia  $H_0$

## Hipotezy dwustronne

$$\begin{cases} H_0 : \beta_k = 0 \\ H_1 : \beta_k \neq 0 \end{cases}$$

Jeśli brak podstaw do odrzucenia  $H_0$ , wówczas model ma postać:

$$\mathbf{y} = \beta_0 + \cdots + \underbrace{\beta_k}_{0} \mathbf{x}_k + \cdots + \beta_K \mathbf{x}_K + \varepsilon$$

zmienna  $\mathbf{x}_k$  nie ma znaczenia dla wyjaśnienia zmienności  $\mathbf{y}$

# Statystyka testowa

- statystyka testowa:

$$t = \frac{b_k}{\hat{se}(b_k)}$$

czyli jest to stosunek wielkości **estymatora parametru** przez estymator jego **odchylenia standardowego**

- statystyka krytyczna (odczytujemy z tablic rozkładu *t-Studenta*):

$$t^* = t\left(N - K, 1 - \frac{\alpha}{2}\right)$$

- Jeśli  $|t| \geq t^*$  - *odrzucaamy*  $H_0$
- Jeśli  $|t| < t^*$  - *nie ma podstaw do odrzucenia*  $H_0$

## Wnioskowanie statystyczne

### hipotezy dwustronne

$$P(|t| > t^*) = 2[1 - F_{t_{N-K}}(t^*)] = \alpha$$

gdzie:  $t^*$  – statystyka krytyczna.

Obecnie, zamiast stosować wartości krytyczne, oblicza się  $p$  – *value* (*policzony poziom istotności*):

$$2[1 - F_{t_{N-K}}(t)] = p - value$$

gdzie:  $t$  – statystyka testowa.

- Jeśli  $p$  – *value* poniżej określonego poziomu istotności (np. 0,05) - odrzucamy  $H_0$
- W przeciwnym przypadku - nie ma podstaw do odrzucenia  $H_0$

- Jaki jest przedział, w którym z określonym prawdopodobieństwem znajdzie się nieznaną wartość parametru  $\beta_k$ .
- Odpowiedź na to pytanie uzyskamy wyznaczając tak zwany przedział ufności.
- Przedział ufności dla nieznanego parametru  $\beta_k$  na poziomie ufności  $1 - \alpha$  można skonstruować następująco:

$$\begin{aligned} P(|t| < t^*) &= P\left(\left|\frac{b_k - \beta_k}{\hat{s}e(b_k)}\right| < t^*\right) = \\ &= P(b_k - \hat{s}e(b_k)t^* < \beta_k < b_k + \hat{s}e(b_k)t^*) = 1 - \alpha \end{aligned}$$

gdzie:

$$t^* = t\left(N - K, 1 - \frac{\alpha}{2}\right)$$

Hipotezy łączne są ważne z punktu widzenia:

- rozważań teoretycznych
- doboru zmiennych do modelu

**Uwaga:**

Hipotezy łączne nie są równoważne iloczynowi hipotez prostych!

## Typowa hipoteza łączna

dana jest układem równań:

$$H_0 : \mathbf{H}\beta = \mathbf{h}$$

gdzie:  $\mathbf{H}$  - macierz o pełnym rzędzie wierszowym =  $g$ .

Liczba równań w tym układzie nazywana jest **liczba ograniczeń**

Układ równań:

- zawiera równania liniowo niezależne
- nie jest sprzeczny



- 1 (\*) Udowodnić, że rozkład sumy kwadratów reszt jest rozkładem  $\chi^2_{N-K}$  niezależnym od rozkładu  $b$ .
- 2 Wyprowadzić rozkład małopróbkowy estymatora MNK. Jakie założenie, poza standardowymi KMRL, należy w tym przypadku przyjąć?
- 3 Jaka postać ma statystyka służąca do testowania hipotezy o tym, że  $\beta_k = \beta_k^*$  ?
- 4 Mając oszacowanie  $b_k$  oraz oszacowanie odchylenia standardowego tego oszacowania  $\hat{se}(b_k)$  wyjaśnić w jaki sposób należy zbudować przedział ufności dla  $\beta_k$ . Ilość obserwacji wynosi  $N$ , ilość szacowanych parametrów  $K$ , a poziom ufności  $1 - \alpha$ .