

ZMIENNE BINARNE I KATEGORYCZNE

EKONOMETRIA WNE

Sebastian Zalas

University of Warsaw

s.zalas@uw.edu.pl

WPROWADZENIE

- Dotychczas zmienne miały interpretację *ilościową*
przykłady: wynagrodzenie, średnia ocen, cena domów etc.
- w pracy empirycznej należy również uwzględniać czynniki *jakościowe*
 - płeć lub rasa
 - sektor w którym operuje firma (przemysł, usługi etc.)
 - region
- dziś zajmiemy się **zależnymi** zmiennymi jakościowymi

WPROWADZENIE

- Dotychczas zmienne miały interpretację *ilościową*
przykłady: wynagrodzenie, średnia ocen, cena domów etc.
- w pracy empirycznej należy również uwzględniać czynniki *jakościowe*
 - płeć lub rasa
 - sektor w którym operuje firma (przemysł, usługi etc.)
 - region
- dziś zajmiemy się **zależnymi** zmiennymi jakościowymi

WPROWADZENIE

- Dotychczas zmienne miały interpretację *ilościową*
przykłady: wynagrodzenie, średnia ocen, cena domów etc.
- w pracy empirycznej należy również uwzględniać czynniki *jakościowe*
 - płeć lub rasa
 - sektor w którym operuje firma (przemysł, usługi etc.)
 - region
- dziś zajmiemy się **zależnymi** zmiennymi jakościowymi

ZMIENNE BINARNE

- czynniki jakościowe można wyrazić jako **zmienną zero-jedynkową (zmienną binarną)**.
- w takich przypadkach należy zdecydować jaka cecha ma przyjąć wartość 1 (a jaka 0)
- Przykład: analizujemy zależność między wynagrodzeniem a płcią. Płeć możemy zakodować jako:
 1. gender (1 kobieta) (0 mężczyzna)
 2. female (1 kobieta) (0 mężczyzna)
 3. male (0 kobieta) (1 mężczyzna)
- Który sposób jest najlepszy?

ZMIENNE BINARNE

- *gender* jest nieintuicyjne: nie domyślamy się co oznacza wartość 1
- *male* lub *female* to odpowiedni wybór, w zależności od pytania które stawiamy.
- Dlaczego zmienne zerojedynkowe przyjmują wartości 0 i 1?
 - wartości te są arbitralne, jakiegokolwiek dwie wartości mogłyby opisać cechę jakościową.
 - stosowanie zera i jedynki opłaca się, ponieważ w regresji taka zmienna zyskuje intuicyjną interpretację

ZMIENNE BINARNE

- *gender* jest nieintuicyjne: nie domyślamy się co oznacza wartość 1
- *male* lub *female* to odpowiedni wybór, w zależności od pytania które stawiamy.
- Dlaczego zmienne zerojedynekowe przyjmują wartości 0 i 1?
 - wartości te są arbitralne, jakiegokolwiek dwie wartości mogłyby opisać cechę jakościową.
 - stosowanie zera i jedynki opłaca się, ponieważ w regresji taka zmienna zyskuje intuicyjną interpretację

ZMIENNE BINARNE

Rozważmy model objaśniający płace:

$$wage = \beta_0 + \delta female + \beta_1 educ + \varepsilon \quad (1)$$

gdzie $female = 1 \Rightarrow$ kobieta, $female = 0 \Rightarrow$ mężczyzna

Co mierzy δ ?

$$\mathbb{E}[wage \mid female = 1, educ] = \beta_0 + \delta \times 1 + \beta_1 educ + \varepsilon$$

$$\mathbb{E}[wage \mid female = 0, educ] = \beta_0 + \delta \times 0 + \beta_1 educ + \varepsilon$$

Po odjęciu stronami otrzymujemy:

$$\delta = \mathbb{E}[wage \mid female = 1, educ] - \mathbb{E}[wage \mid female = 0, educ]$$

δ = różnica średnich płac między kobietami i mężczyznami, przy takim samym poziomie edukacji

ZMIENNE BINARNE

Używając danych wage1 oszacowano model objaśniający płace (1):

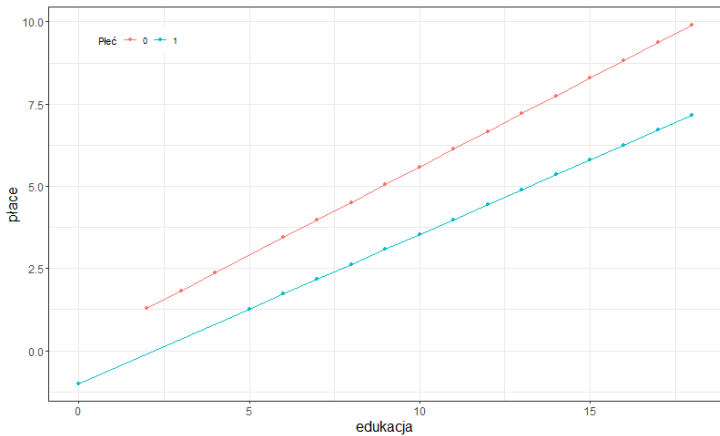
$$wage = 0.62 - 2.27 \text{ female} + 0.51 \text{ educ} + e$$

- interpretacja $\delta = 2.27$:

Kobiety zarabiają średnio o \$2.27 na godzinę mniej niż mężczyźni przy niezmiennym poziomie wykształcenia

- mężczyźni są poziomem odniesienia dla kobiet.
- w przypadku modelu z większą liczbą zmiennych objaśniających, interpretacja jest analogiczna

ZMIENNE BINARNE



ZMIENNE BINARNE - $\log(y)$

Model objaśniający *logarytm* płac:

$$\log(wage) = \beta_0 + \delta female + \beta_1 educ + u$$

Co mierzy δ ? Zapiszmy modele dla kobiet i mężczyzn:

$$\mathbb{E}[\log(wage)|female = 1, educ] = \beta_0 + \delta \times 1 + \beta_1 educ + \varepsilon$$

$$\mathbb{E}[\log(wage)|female = 0, educ] = \beta_0 + \delta \times 0 + \beta_1 educ + \varepsilon$$

odejmijmy stronami

$$\delta = \mathbb{E}[\log(wage)|female = 1, educ] - \mathbb{E}[\log(wage)|female = 0, educ]$$

teraz mamy logarytmy!

$$e^{\delta} = \frac{wage_{female=1}}{wage_{female=0}}$$

ZMIENNE BINARNE - $\log(y)$

Zmiana procentowa:

$$\frac{wage_{female=1}}{wage_{female=0}} - 1 = e^{\delta} - 1$$
$$\frac{wage_{female=1} - wage_{female=0}}{wage_{female=0}} \times 100 = 100 \times (e^{\delta} - 1)$$

Interpretacja $\delta \Rightarrow$ kobiety zarabiają średnio mniej/więcej niż mężczyźni o $100(e^{\delta} - 1)$ procent, przy innych czynnikach niezmiennych

Do interpretacji możemy użyć także przybliżenia $e^{\delta} \approx \delta + 1$ (szeregu Taylora):

$$\frac{wage_{female=1} - wage_{female=0}}{wage_{female=0}} \times 100 = 100 \delta$$

ZMIENNE BINARNE - $\log(y)$

Oszacowano model objaśniający *logarytm* płac:

$$\log(wage) = 0.8263 - 0.3609female + 0.0772educ + u$$

Interpretacja:

- $100 * (e^{\delta} - 1) = -30.29272$
- $100\delta \approx -36.08654$
- kobiety zarabiają średnio mniej niż mężczyźni o 30 procent (36 procent), przy innych czynnikach niezmiennych
- przybliżone oszacowanie efektu pozostanie takie samo przy zamianie poziomu odniesienia

ZMIENNE KATEGORYCZNE

- zmienne jakościowe mogą opisywać więcej niż dwa poziomy:
 - sektory: przemysł, usługi, rolnictwo..
 - zawody
 - wykształcenie: podstawowe, średnie, wyższe
- aby te cechy uwzględnić w regresji, każdy poziom należy zakodować jako oddzielne zmienne binarne
- należy wybrać jedną grupę jako poziom odniesienia i nie uwzględniać go regresji

ZMIENNE KATEGORYCZNE

- ▶ model objaśniający płace za pomocą statusu matrymonialnego (dla kobiet i mężczyzn) oraz edukacji i doświadczenia:

$$\begin{aligned}\log(wage) = & 0.32 + 0.21married_male - 0.19married_fem \\ & - 0.11single_fem + 0.079educ + 0.027exper \\ & - 0.0005exper^2 + 0.02tenure - 0.0005tenure^2 + e\end{aligned}$$

- ▶ poziom odniesienia? nieżonaci mężczyźni (single)
 - ▶ żonaci mężczyźni zarabiają średnio o $\approx 21\%$ więcej niż single, przy pozostałych czynnikach niezmiennych
 - ▶ zamężne kobiety zarabiają średnio o $\approx 20\%$ mniej niż single, przy pozostałych czynnikach niezmiennych

ZMIENNE KATEGORYCZNE UPORZĄDKOWANE

- szczególny przypadek zmiennej kategorycznej: tu kolejność poziomów zmiennej ma znaczenie
- jako poziom odniesienia można wybrać np.: najniższy poziom
- przykłady:
 - ▶ edukacja (1=podstawowe, 2=średnie, 3=wyższe)
 - ▶ ranking szkół

ZMIENNE KATEGORYCZNE UPORZĄDKOWANE

- ▶ medianą początkowych zarobków absolwentów szkół prawnych a ranking szkoły
- ▶ można podzielić ranking na grupy: *top10*, *r11_25*, *r26_40*, *r41_60*, *r61_100*

$$\begin{aligned}\log(\text{salary}) = & 9.17 + 0.700\text{top10} + 0.594r11_25 + 0.375r26_40 \\ & + 0.263r41_60 + 1.132r61_100 + 0.0057LSAT \\ & + 0.014GPA + 0.036 \log(\text{libvol}) + 0.0008 \log(\text{cost}) + e\end{aligned}$$

- ▶ interpretacja: absolwenci topowych dziesięciu szkół mają $(e^{0.7} - 1 = 1.014)$ medianę początkowych zarobków wyższą o ponad 100% w porównaniu z absolwentami szkół o rankingu poniżej setnego, przy pozostałych czynnikach niezmiennych

PODSUMOWANIE

- ▶ zmienne binarne/kategoryczne pozwalają uwzględnić cechy jakościowe
- ▶ zmienne kategoryczne sprowadzamy do zbiory zmiennych binarnych
- ▶ uważamy na poziom odniesienia i logarytmy

Pytania? Wątpliwości?
Dziękuję!

e: s.zalas@uw.edu.pl