

# Egzamin z ekonometrii - wersja II E, MSEMAT

7-02-2013

## Pytania teoretyczne

1. Porównać zastosowania znanych Ci kontrastów ze standardowym sposobem rozkodowania zmiennej dyskretnej.
2. Wyprowadzić estymator MNK dla modelu z wieloma zmiennymi objaśniającymi. Proszę również pokazać warunek dostateczny na istnienie minimum funkcji i udowodnić, że jest to rzeczywiście minimum tej funkcji.
3. Wyprowadzić postać macierzy wariancji kowariancji  $b$  i podać interpretację jej elementów.
4. Wyjaśnić, dlaczego  $R^2$  nie można używać do porównywania modeli.

**Zadanie 1**

Dla modelu:

$$(\star) y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I)$$

gdzie  $X$  jest nielosowe,

stworzono macierz  $X^* = XA$ , gdzie  $A$  jest pewną macierzą nieosobliwą, oraz wektor  $y^* = cy$ , gdzie  $c \in \mathbb{R}$ ,  $c \neq 0$ .

Następnie zdefiniowano następujący model:

$$(\star\star) y^* = X^*\beta^* + \eta, \quad \eta \sim N(0, \sigma^2 I)$$

1. Wyznaczyć estymator MNK dla regresji  $(\star\star)$ . Estymator  $b^*$  należy przedstawić jako funkcję estymatora  $b$ .
2. Wyznaczyć  $Var(b^*)$ .
3. Wyznaczyć wektor reszt dla regresji  $(\star\star)$  jako funkcję wektora reszt dla regresji  $(\star)$ .
4. Pokazać, że  $R^2$  w obu regresjach  $(\star)$  i  $(\star\star)$  jest takie same.

**Rozwiązanie Zadanie 1**

1. Wyznaczyć estymator MNK dla regresji  $(\star\star)$ ? Estymator  $b^*$  należy przedstawić jako funkcję estymatora  $b$ .

$$b = ((X^*)' X^*)^{-1} (X^*)' y^* = \{X^* = XA, y^* = cy\} = ((XA)'(XA))^{-1} (XA)'(cy) = c(A'X'XA)^{-1} A'X'y.$$

Dla dowolnych nieosobliwych macierzy  $A$  i  $B$  zachodzi  $(AB)^{-1} = B^{-1}A^{-1}$ .

Uogólniając ten wzór na przypadek trzech odwracalnych macierzy mamy  $(ABC)^{-1} = C^{-1}B^{-1}A^{-1}$ .

Czyli  $(A'X'XA)^{-1} = A^{-1}(X'X)^{-1}(A')^{-1}$ .

$$b^* = cA^{-1}(X'X)^{-1} \underbrace{(A')^{-1}A'}_{I} X'y = cA^{-1}(X'X)^{-1}X'y = cA^{-1}b$$

2. Wyznaczyć  $Var(b^*)$ .

$$Var(b^*) = Var(cA^{-1}b) = cA^{-1}Var(b)(cA^{-1})' = c^2A^{-1}Var(b)(A^{-1})' = c^2\sigma^2A^{-1}(X'X)^{-1}(A^{-1})'.$$

3. Wyznaczyć wektor reszt dla regresji  $(\star\star)$  jako funkcję wektora reszt dla regresji  $(\star)$ .

$$\hat{y}^* = X^*b^* = (XA)cA^{-1}b = cXAA^{-1}b = cXb = c\hat{y}.$$

$$e^* = y^* - \hat{y}^* = cy - c\hat{y} = c(y - \hat{y}) = ce$$

4. Pokazać, że  $R^2$  w obu regresjach (★) i (★★) jest takie samo.

$$R^{2*} = 1 - \frac{RSS^*}{TSS^*} = 1 - \frac{(e^*)' e^*}{(y^* - \bar{y}^*)'(y^* - \bar{y}^*)}.$$

$$(e^*)' e^* = (ce)'(ce) = e' c' ce = c^2 e' e$$

$$\bar{y}^* = \frac{1}{N} \sum cy_i = \frac{1}{N} c \sum y_i = c\bar{y}.$$

$$TSS^* = (y^* - \bar{y}^*)'(y^* - \bar{y}^*) = (cy - c\bar{y})'(cy - c\bar{y}) = c^2(y - \bar{y})'(y - \bar{y}).$$

$$R^{2*} = 1 - \frac{RSS^*}{TSS^*} = 1 - \frac{(e^*)' e^*}{(y^* - \bar{y}^*)'(y^* - \bar{y}^*)} = 1 - \frac{c^2 e' e}{c^2 (y - \bar{y})'(y - \bar{y})} = R^2$$

**Zadanie 2**

Na podstawie danych BAEL z 2010 roku oszacowano wysokość wynagrodzenia dla pracowników w mikroprzedsiębiorstwach ( $\ln\_placa$  - logarytm wynagrodzenia). Zmiennymi objaśniającymi są wiek ( $wiek$ ), wiek do kwadratu ( $wiek\_2$ ), miejsce zamieszkania ( $wies$ : 0 - miasto, 1 - wieś), płeć ( $plec$ : 0 - mężczyzna, 1 - kobieta), wykształcenie ( $wyksz$ : 0 - podstawowe, 1 - zawodowe, 2 - średnie, 3 - wyższe), wymiar zatrudnienia ( $etat$ : 0 - niepełny etat, 1 - cały etat), stan cywilny ( $stan$ : 0 - kawaler / panna, 1 - zamężny / zamężna, 2 - wdowiec / wdowa, 3 - rozwiedziony / rozwiedziona), interakcja między płcią a wymiarem zatrudnienia. Oszacowania parametrów znajdują się poniżej. Hipotezy testować na poziomie istotności 0,05. Odpowiedzi uzasadnić podając  $p$ -value.

Source	SS	df	MS	Number of obs = 1967		
Model	211.507908	12	17.625659	F( 12, 1954) = 118.05		
Residual	291.740412	1954	.149304203	Prob > F = 0.0000		
-----				R-squared = 0.4203		
Total	503.24832	1966	.255975748	Adj R-squared = 0.4167		
-----				Root MSE = .3864		
$\ln\_placa$	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
wyksz_1	.2021485	.0303484	6.66	0.000	.1426299	.2616671
wyksz_2	.2945919	.0362894	8.12	0.000	.2234219	.3657619
wyksz_3	.4560558	.0378275	12.06	0.000	.3818694	.5302422
wies	-.0433189	.0180795	-2.40	0.017	-.078776	-.0078618
wiek	.0388162	.0056042	6.93	0.000	.0278253	.0498072
wiek_2	-.0004376	.0000676	-6.47	0.000	-.0005702	-.0003049
plec	-.1072408	.0505662	-2.12	0.034	-.2064101	-.0080716
etat	.780341	.0444233	17.57	0.000	.693219	.8674629
plecXetat	-.1257776	.0534618	-2.35	0.019	-.2306257	-.0209296
stan_1	.0951863	.0240566	3.96	0.000	.0480071	.1423656
stan_2	-.1416434	.0605747	-2.34	0.019	-.2604411	-.0228456
stan_3	-.0243112	.0445971	-0.55	0.586	-.111774	.0631517
_cons	5.470661	.106983	51.14	0.000	5.260849	5.680474
-----				Breusch-Pagan LM statistic: Chi2(1) = 284.26 p-value = 0.000		
-----				White's general test statistic: Chi2(70) = 508.28 p-value = 0.000		
-----				Jarque-Berra test statistic: Chi2(2) = 3185.88 p-value = 0.000		
-----				Ramsey RESET test statistic: F(3,1951) = 44.72 p-value = 0.000		

1. Czy zmienne objaśniające są łącznie istotne?
2. Zinterpretować wartość współczynnika determinacji.
3. Ocenic, które zmienne są istotne.
4. Wykonać następujące polecenia:
  - (a) Policzyc i zinterpretować semielastycznosc cząstkową dla *wieku*. Średni wiek respondentów w próbie wynosi 42 lata.
  - (b) Zintpretować wielkość współczynnika przy interakcji pomiędzy *płcią a etatem*.
  - (c) Zintpretować wielkość współczynnika przy zmiennej *etat*.
5. Zbadać, czy w modelu występuje heteroskedastyczność.

6. Zbadać, czy błąd losowy ma rozkład normalny.
7. Sprawdzić, czy forma funkcyjna modelu jest poprawna.
8. Jeżeli model nie spełnia założeń KMRL określić:
  - (a) Które założenia nie są spełnione?
  - (b) Jakie to ma konsekwencje dla interpretacji modelu i wnioskowania statystycznego?
  - (c) W jaki sposób można rozwiązać problemy zasygnalizowane przez wyniki testów?

### Rozwiązanie Zadanie 2

1. Test na łączną istotność regresji:  $F = 118.05$ ,  $p - value = 0.000 < 0.05$  odrzucamy hipotezę zerową o łącznej nieistotności regresji.
2. 42.03% zmienności zmiennej zależnej zostało wyjaśnione za pomocą zmienności zmiennych niezależnych.
3. Istotne zmienne, to te dla których  $p - value$  jest mniejsze od przyjętego poziomu istotności wynoszącego 0.05. Czyli istotne zmienne to:
  - (a)  $wyksz\_1$  ( $t = 6.66$ ,  $p - value = 0.000$ )
  - (b)  $wyksz\_2$  ( $t = 8.12$ ,  $p - value = 0.000$ )
  - (c)  $wyksz\_3$  ( $t = 12.06$ ,  $p - value = 0.000$ )
  - (d)  $wies$  ( $t = -2.4$ ,  $p - value = 0.017$ )
  - (e)  $wiek$  ( $t = 6.93$ ,  $p - value = 0.000$ )
  - (f)  $wiek\_2$  ( $t = -6.47$ ,  $p - value = 0.000$ )
  - (g)  $plec$  ( $t = -2.12$ ,  $p - value = 0.034$ )
  - (h)  $etat$  ( $t = 17.57$ ,  $p - value = 0.000$ )
  - (i)  $plec \times etat$  ( $t = -2.35$ ,  $p - value = 0.019$ )
  - (j)  $stan\_1$  ( $t = 3.96$ ,  $p - value = 0.000$ )
  - (k)  $stan\_2$  ( $t = -2.34$ ,  $p - value = 0.019$ )
  - (l)  $const$  ( $t = 51.14$ ,  $p - value = 0.000$ )

4. Wykonać następujące polecenia:

- (a) semielastyczność cząstkowa dla wieku:

$$\frac{\partial E(\ln\_placa)}{\partial wiek} = \beta_{wiek} + 2 \cdot \beta_{wiek\_2} \cdot \bar{wiek} = 0,002058$$

wzrost wieku osób 42 letnich o 1 rok powoduje wzrost wynagrodzenia średnio o 0.2% *ceteris paribus*.

- (b) kobiety pracujące na cały etat mają średnio o 92.3% ( $100\% \cdot [\exp(0,78) \cdot \exp(-0,126) - 1] = 100\% \cdot [\exp(0,654) - 1]$ ) wyższe wynagrodzenie niż kobiety pracujące na niepełny etat *ceteris paribus*.
- (c) mężczyźni pracujący na cały etat mają średnio o 118% ( $100\% \cdot [\exp(0,78) - 1]$ ) wyższe wynagrodzenie niż mężczyźni pracujący na niepełny etat *ceteris paribus*.

5. Występowanie heteroskedastyczności testujemy za pomocą:

- (a) testu White'a:
- hipoteza zerowa: homoskedastyczność składnika losowego.
  - wartość statystyki testowej wynosi:  $\chi^2(70)=508.28$  oraz  $p\text{-value} = 0.000 < 0.05$ , więc odrzucamy hipotezę zerową o homoskedastyczności.
- (b) testu Breuscha-Pagana:
- hipoteza zerowa: homoskedastyczność składnika losowego.
  - wartość statystyki testowej wynosi  $\chi^2(1) = 284.26$  oraz  $p\text{-value} = 0.000 < 0.05$ , więc odrzucamy hipotezę zerową o homoskedastyczności.
6. Normalność zaburzenia losowego testujemy za pomocą:
- (a) testu Jarque-Bera:
- hipoteza zerowa: zaburzenie losowe ma rozkład normalny.
  - wartość statystyki testowej wynosi  $\chi^2(2) = 3185.88$  oraz  $p\text{-value} = 0.000 < 0.05$ , czyli odrzucamy hipotezę zerową o normalności zaburzenia losowego.
7. Poprawność przyjętej formy funkcyjnej modelu testujemy za pomocą:
- (a) test RESET:
- hipoteza zerowa: przyjęta postać funkcyjna modelu jest prawidłowa.
  - wartość statystyki testowej  $F(3, 1951) = 44.72$  i  $p\text{-value} = 0.000 < 0.05$ , więc odrzucamy hipotezę zerową o poprawności przyjętej formy funkcyjnej.
8. Odpowiedzi są następujące:
- (a) Nie są spełnione następujące założenie:
- o homoskedastyczności zaburzenia losowego,
  - o sposobie „generowania danych”:  $y = \beta x + \varepsilon$  (czyli założenie o liniowej zależności między zmienną zależną i zmiennymi niezależnymi),
  - nie jest także spełnione dodatkowe założenie o normalności składnika losowego.
- (b) Konsekwencje dla interpretacji modelu i wnioskowania statystycznego są następujące:
- W przypadku nie spełnienia założenia o homoskedastyczności zaburzenia losowego, estymator  $b$  jest co prawda nieobciążony i zgodny, ale nieefektywny. Estymator macierzy wariancji-kowariancji  $b$  jest już obciążony i niezgodny. Macierz wariancji-kowariancji jest wykorzystywana do testowania hipotez na temat istotności zmiennych, więc poprawność wnioskowania statystycznego jest podważona.
  - Odrzucenie hipotezy o poprawności przyjętej formy funkcyjnej podważa interpretację ekonomiczną modelu (interpretacja oszacowanych parametrów). Takie własności jak nieobciążoność czy efektywność estymatora MNK są wyprowadzane przy założeniu prawdziwości przyjętej formy funkcyjnej modelu.
  - Próba zawiera 1967 obserwacji (można przyjąć, iż jest to duża próba). Dla dużych prób rozkłady statystyk są bliskie standardowym rozkładom.
- (c) Rozwiązanie problemów zasygnalizowanych przez wyniki testów:
- Niepoprawna forma funkcyjna:** możemy próbować poprawić formę funkcyjną modelu wprowadzając do modelu interakcje między zmiennymi, dokonać przekształceń zmiennych (np. przekształcenie Boxa-Coxa), zastosować model wielomianowy, schodkowy lub krzywej łamanej.
  - Problem heteroskedastyczności** można rozwiązać za pomocą Stosowalnej UMNK lub odpornego estymatora White'a macierzy wariancji kowariancji.

**Zadanie 3**

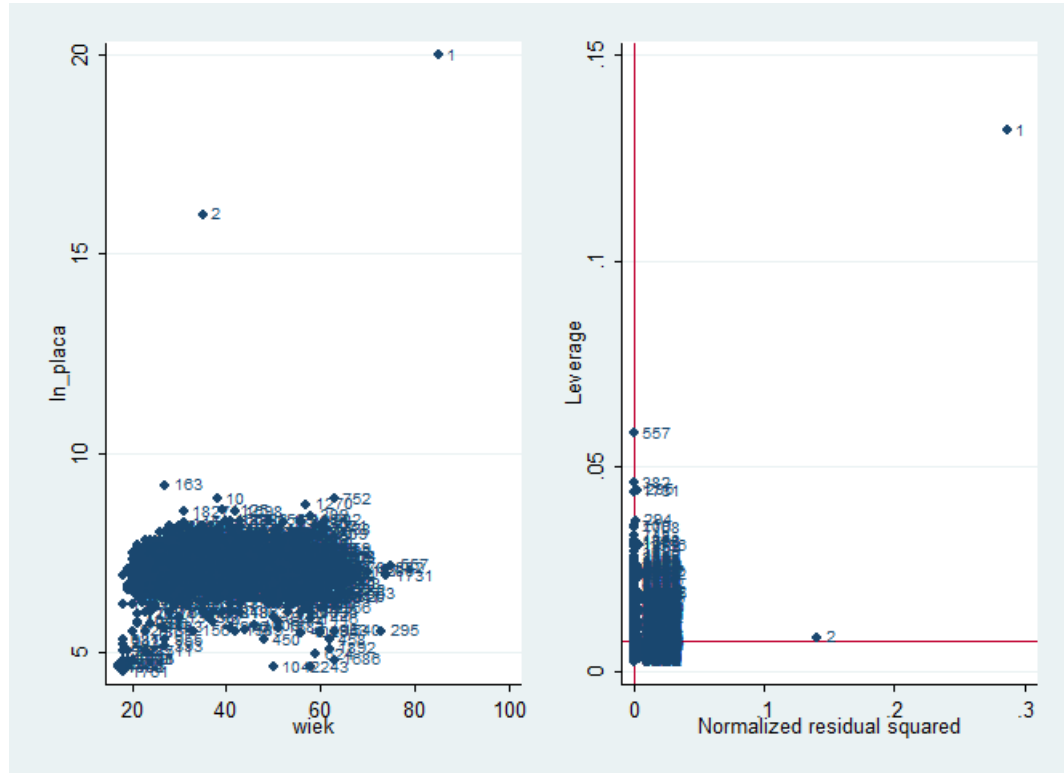
Na podstawie danych BAEL z 2010 roku oszacowano wysokość wynagrodzenia dla pracowników w mikroprzedsiębiorstwach ( $\ln\_placa$  - logarytm wynagrodzenia). Zmiennymi objaśniającymi są wiek ( $wiek$ ), wiek do kwadratu ( $wiek\_2$ ), miejsce zamieszkania ( $wies$ : 0 - miasto, 1 - wieś), płeć ( $plec$ : 0 - mężczyzna, 1 - kobieta), wykształcenie ( $wyksz$ : 0 - podstawowe, 1 - zawodowe, 2 - średnie, 3 - wyższe), wymiar zatrudnienia ( $etat$ : 0 - niepełny etat, 1 - cały etat), stan cywilny ( $stan$ : 0 - kawaler / panna, 1 - zamężny / zamężna, 2 - wdowiec / wdowa, 3 - rozwiedziony / rozwiedziona), interakcja między płcią a wymiarem zatrudnienia. Oszacowania parametrów znajdują się poniżej. Hipotezy testować na poziomie istotności 0,05. Odpowiedzi uzasadnić podając  $p$ -value.

1. Przed oszacowaniem regresji uzyskano statystyki opisowe dla zmiennej wysokość wynagrodzenia dla pracowników w mikroprzedsiębiorstwach ( $placa$ ). Wyjaśnić, które obserwacje budzą podejrzenia i dlaczego?

Variable	Obs	Mean	Std. Dev.	Min	Max
placa	4534	57188.11	48913.66	90	99999

```
count if placa==99999
2567
```

2. Po oszacowanej regresji badacz uzyskał wykres rozrzutu pomiędzy wiekiem a logarytmem płacy oraz wykres przedstawiający zależność pomiędzy dźwignią a wystandaryzowanymi resztami podniesionymi do kwadratu. Wyjaśnić, które obserwacje budzą podejrzenia i dlaczego?



3. W modelu uwzględniono dodatkowo dwie zmienne: staż w obecnym miejscu pracy (*staz*), ogólny staż pracy (*staz\_ogolny*). Następnie obliczono współczynniki korelacji *Spearmana*. Jaki problem występuje w tym modelu?

Source	SS	df	MS	Number of obs = 1967		
Model	216.967865	14	15.4977047	F( 14, 1952)	=	105.67
Residual	286.280455	1952	.146660069	Prob > F	=	0.0000
-----				R-squared	=	0.4311
-----				Adj R-squared	=	0.4271
Total	503.24832	1966	.255975748	Root MSE	=	.38296
-----						
ln_placa	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
wyksz_1	.1889953	.0302267	6.25	0.000	.1297153	.2482753
wyksz_2	.2786084	.0360692	7.72	0.000	.2078703	.3493465
wyksz_3	.4458425	.0375599	11.87	0.000	.3721808	.5195043
wies	-.0512427	.0179803	-2.85	0.004	-.0865054	-.01598
wiek	.0358145	.0056632	6.32	0.000	.0247081	.046921
wiek_2	-.0004565	.0000674	-6.77	0.000	-.0005887	-.0003243
plec	-.1157975	.0501914	-2.31	0.021	-.2142319	-.0173632
etat	.7589422	.0441721	17.18	0.000	.6723127	.8455717
plecXetat	-.1165905	.0530097	-2.20	0.028	-.220552	-.012629
stan_1	.0928414	.0238605	3.89	0.000	.0460467	.1396362
stan_2	-.1109334	.0602626	-1.84	0.066	-.2291193	.0072524
stan_3	-.0083049	.0442809	-0.19	0.851	-.0951477	.0785379
staz	.0076781	.0014517	5.29	0.000	.0048311	.0105252
staz_ogolny	.0026003	.001987	1.31	0.191	-.0012967	.0064972
_cons	5.566119	.1095404	50.81	0.000	5.351291	5.780947

Macierz korelacji Spearmana:

	wiek	staz	staz_ogolny	wyksz
wiek	1.0000			
staz	0.4875	1.0000		
staz_ogolny	0.9201	0.5743	1.0000	
wyksz	-0.1575	0.0251	-0.1777	1.0000



4. Przeprowadzono regresję wysokości wynagrodzenia dla pracowników w mikroprzedsiębiorstwach ( $\ln\_placa$  - logarytm wynagrodzenia) na latach poświęconych na naukę ( $nauka$ ).

Source	SS	df	MS			
Model	13.7328051	1	13.7328051	Number of obs =	1968	
Residual	735.877367	1966	.374301814	F( 1, 1966) =	36.69	
Total	749.610172	1967	.381093122	Prob > F =	0.0000	
				R-squared =	0.0183	
				Adj R-squared =	0.0178	
				Root MSE =	.6118	

$\ln\_placa$	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
$nauki$	.1009973	.0166741	6.06	0.000	.0682967	.133698
$\_cons$	6.948158	.0260974	266.24	0.000	6.896976	6.999339

(a) Jaki będzie prawdopodobny kierunek obciążenia oszacowania parametru przy zmiennej  $nauka$  w wyniku pominięcia zmiennej związanej z wielkością miejscowości.

**Rozwiązanie Zadanie 3**

1. Największe podejrzenia budzą obserwacje, dla których  $placa=99999$ . Zwraca uwagę fakt, iż dla 2567 obserwacji spośród 4534 wysokość wynagrodzenia jest zaskakująco wysoka i dokładnie taka sama. Można podejrzewać, że nie są to rzeczywiste obserwacje dla wynagrodzenia, ale kody braków (np. odpowiedź: nie wiem). Jeśli rzeczywiście obserwacje, dla których  $placa=99999$  są błędne, to powinniśmy je usunąć z modelu.
2. Na podstawie wykresu rozrzutu pomiędzy wiekiem a logarytmem wynagrodzenia podejrzenia budzą obserwacje **nr 1** i **nr 2**. Obserwacja **nr 1** jest nietypowa na tle pozostałych obserwacji ze względu na wysoką wartość zmiennej *wiek* i wysoką wartość *logarytmu wynagrodzenia*. Obserwacja **nr 2** jest nietypowa na tle pozostałych obserwacji ze względu na wysoką wartość *logarytmu wynagrodzenia*. Na podstawie wykresu przedstawiającego zależność pomiędzy *dźwignią* a *wystandaryzowanymi resztami podniesionymi do kwadratu* możemy potwierdzić, iż obserwacja **nr 1** charakteryzuje się zarówno wysoką wartością *dźwigni* jak i wysokimi wartościami *wystandaryzowanych reszt podniesionymi do kwadratu*. Znaczący wpływ na oszacowania ma dodanie obserwacji o **nr 1**, która jest nietypowa ze względu na zmienną objaśniającą oraz nie pasuje do linii regresji (otrzymujemy duże reszty). Jest to obserwacja nietypowa, która nie pasuje do prostej regresji. Natomiast obserwacja **nr 2** cechuje się wysoką wartością *wystandaryzowanych reszt podniesionymi do kwadratu*. Dodanie obserwacji **nr 2** do próby nie powoduje znacznych zmian w oszacowanych parametrach (jest to obserwacja o stosunkowo dużej reszcie, ale nie jest nietypowa ze względu na zmienną objaśniającą).
3. Na podstawie macierzy korelacji Spearmana można zauważyć silną zależność pomiędzy zmienną *wiek* a zmienną *staz\_ogolny*. Jednocześnie zmienna *staz\_ogolny* jest nieistotna w modelu ( $p\text{-value}=0.191$ ). W związku z powyższym w modelu wystąpi problem **niedokładnej współliniowości**.
4. W dużych miastach wynagrodzenia są średnio wyższe niż w małych miejscowościach. Jednocześnie w dużych miastach jest też więcej osób dobrze wykształconych. Pominięcie tej zmiennej doprowadzi więc do dodatniego obciążenia oszacowania parametru.