



4EK211 Základy ekonometrie

Umělé proměnné

Predikce

Multikolinearita

Cvičení 5

Umělé proměnné

Umělé proměnné

- **Dummy** / umělé / booleovské **proměnné**
- nabývají hodnotu: 0 nebo 1
- Kvalitativní proměnné: vyjadřují přítomnost či nepřítomnost určité vlastnosti
 - přítomnost popisované vlastnosti... obvykle 1
 - absence vlastnosti ... obvykle 0
- Příklad: proměnná **žena**: $žena_i = 1$ (i -tá osoba je žena)
 $= 0$ (i -tá osoba je muž)
- Příklad: **vzdělání** zkoumaných osob

Dummy proměnná / pozorování	$D1_i$	$D2_i$	$D3_i$
i -tá osoba: základní vzdělání	1	0	0
i -tá osoba: ukončená maturita	0	1	0
i -tá osoba má VŠ titul	0	0	1

Umělé proměnné

- Příklad: **sezónnost**

Pro jednotlivá čtvrtletí vytvoříme umělé proměnné $D1$ až $D4$ (analogicky jako vzdělání), tj.

$D1_i = 1$ pokud i -té pozorování odpovídá 1. čtvrtletí

$D1_i = 0$ pokud i -té pozorování odpovídá jinému než 1. čtvrtletí

PAST UMĚLÝCH PROMĚNNÝCH

Zahrneme-li do LRM umělé proměnné odpovídající všem variantám sledované kvalitativní proměnné, budou umělé proměnné perfektně lineárně závislé (porušení 4. G.M. předpokladu) a MNČ selže, protože matice $(\mathbf{X}^T \mathbf{X})^{-1}$ neexistuje.

Příklad perfektní lineární závislosti (sezónnost):

$$D1_i + D2_i + D3_i + D4_i = 1 \quad \text{pro } \forall i$$

Proč: při práci se čtvrtletními daty každému pozorování i vždy odpovídá právě jedno čtvrtletí (jedna umělá proměnná rovna 1, ostatní rovny nule).

- **PAST UMĚLÝCH PROMĚNNÝCH - ŘEŠENÍ**
- Cíl: **vyvarovat se perfektní lineární závislosti mezi dummies (též: perfektní multikolinearita umělých proměnných)**
- do modelu zahrneme o jednu dummy proměnnou méně než je počet sledovaných vlastností
- Nezahrnutá dummy proměnná tvoří základ (bázi), ke kterému ostatní vlastnosti porovnáваме:
 - dvě pohlaví – jedna dummy
 - tři stupně vzdělání – dvě dummies
 - čtyři čtvrtletí – použijeme maximálně 3 dummies
- **Interpretace** parametrů β u dummies závisí na tom, kterou proměnnou jsem vynechal, bázi vůči které porovnáвам.

Umělé proměnné – příklady

Soubor: CV5_PR1.xls

Data: y = plat učitelů (tis. USD)
 x = roky praxe
 m = pohlaví (1 = muž, 0 = žena)

Zadání: Odhadněte model závislosti y na x a m a interpretujte získané výsledky.

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 m_i + u_i, \quad i = 1, 2, \dots, 15$$

Umělé proměnné – příklady

Soubor: CV5_PR2.xls

Data: y = výdaje na cestování (tis. USD)
 x = výše příjmu (tis. USD)
 D_1 = dosažené vzdělání (1 = základní, 0 = jiné)
 D_2 = dosažené vzdělání (1 = středoškolské, 0 = jiné)
 D_3 = dosažené vzdělání (1 = vysokoškolské, 0 = jiné)

Zadání: Odhadněte model závislosti y na x , D_2 a D_3 a interpretujte získané výsledky.

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 D_{2i} + \beta_3 D_{3i} + u_i, \quad i = 1, 2, \dots, 15$$

Umělé proměnné – příklady

Soubor: CV5_PR3.xls

Data: t = čas

R = příjmy státního rozpočtu (v mld. Kč)

Zadání: Odhadněte model závislosti R na t .

Pokuste se zachytit v modelu vliv posledního čtvrtletí v daném roce (tj. explicitně zapojit čtvrtý kvartál do modelu).

$$R_t = \beta_0 + \beta_1 t + u_t, \quad t = 1, 2, \dots, 15$$

Umělé proměnné – příklady

Soubor: CV5_PR4.xls

Data: *pocet_domu* = počet nově započatých staveb domů v USA
(v tis.)

urok_mira = úroková míra (v %)

Zadání: Odhadněte model závislosti *pocet_domu* na *urok_mira* + zohledněte sezónní vliv v modelu.

Predikujte ***pocet_domu*** v roce 1999.

$$pocet_domu_t = \beta_0 + \beta_1 urok_mira_t + \beta_2 Q_{2t} + \beta_3 Q_{3t} + \beta_3 Q_{4t} + u_t$$

$$t = 1, 2, \dots, 40$$

Predikce

Aplikace EM – predikce obecně

- ekonomické prognózování, předpověď, předvídání
- hlavním cílem je odhad hodnot vysvětlované proměnné mimo interval pozorování s užitím minulé i současné informace
 - extrapolace modelu do budoucna
 - extrapolace modelu do minulosti – tj. před interval pozorování (tzv. retrospektiva)
- predikcí získáváme vyrovnané hodnoty (tj. hodnoty „fitted“)

Predikce ex-ante (resp. dopředu)

- tzv. podmíněná
- podmíněná volbou vysvětlujících proměnných
- na napozorované hodnoty musíme „čekat“

Predikce ex-post (resp. dozadu)

- tzv. pseudopředpověď
- slouží k testování kvality modelu
- napozorované hodnoty jsou již k dispozici

Aplikace EM – predikce ex-ante

- volba podmíněných exogenních proměnných, možné způsoby:
 - zadáno z jiné analýzy
 - zadáno pomocí procentuální změny oproti minulému období (např. o 10 %)
 - zadáno pomocí diferencí
- predikce **bodová**
- predikce **intervalová**
 - se směrodatnou odchylkou **sigma**: $\hat{y}_p \pm \mathbf{sigma}$
 - se směrodatnou odchylkou $\tilde{\mathbf{s}}_p$: $\hat{y}_p \pm t_{1-\alpha/2(n-(k+1))}^* \tilde{\mathbf{s}}_p$
$$\tilde{\mathbf{s}}_p = \mathbf{s} \sqrt{1 + \mathbf{x}_p (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_p^T}$$
 - „s“ ve vzorci je sigma, výpočet viz předchozí cvičení nebo to je sigma ve výstupu PcGive
 - vždy platí: $\tilde{\mathbf{s}}_p > \mathbf{sigma}$
 - intervalový odhad se sigma bývá podhodnocený

Aplikace EM – predikce ex-ante + příklad

Soubor: CV5_PR5.xls

Data: ***Mira_nezam_obec*** = obecná míra nezaměstnanosti (%) → x
 Inflace = míra inflace (%) → y

Zadání: Odhadněte závislost míry inflace (y) na obecné míře nezaměstnanosti (x).

Proved'te **bodovou předpověď ex-ante** ručně i pomocí PcGivu, víte-li, že hodnota ***Mira_nezam_obec*** v roce 2008 je 4,4.

Proved'te **intervalovou předpověď ex-ante** ručně i pomocí PcGivu, víte-li, že hodnota ***Mira_nezam_obec*** v roce 2008 je 4,4 a $\alpha = 1 \%$.

$$y_t = \beta_0 + \beta_1 x_t + u_t \quad t = 1, 2, \dots, 15$$

Aplikace EM – predikce ex-post + příklad

- testuje se kvalita modelu
- vyřadíme určitý počet pozorování z modelu
- odhadneme model
- provedeme predikci vynechaných hodnot
- porovnáme získané předpovědi se skutečnými hodnotami
- obecně platí, že predikce je dobrá, pokud je absolutní hodnota chyby predikce menší než 5 % ze skutečné hodnoty pro dané období

Příklad: Proveďte predikci/předpověď ex-post pro roky 2006 a 2007 na datech CV5_PR5.xls.

Multikolinearita

Gauss-Markovy předpoklady

Opakování: Gaussovy-Markovovy předpoklady

1. $E(\mathbf{u}) = 0$
 - průměrná hodnota náhodné chyby je nula
2. $E(\mathbf{u} \mathbf{u}^T) = \sigma^2 I_n$
 - diagonála: konečný a konstantní rozptyl = homoskedasticita
 - mimo diagonálu: náhodné složky jsou sériově nezávislé
3. \mathbf{X} je nestochastická matice – $E(\mathbf{X}^T \mathbf{u}) = 0$
 - veškerá náhodnost \mathbf{y} je obsažena v náhodné složce
4. \mathbf{X} má plnou hodnost $k+1$ (hodnost matice \mathbf{X} = počet sloupců \mathbf{X})
 - matice \mathbf{X} neobsahuje žádné perfektně lineárně závislé sloupce
pozorování vysvětlujících proměnných
→ porušení: multikolinearita

Multikolinearita - definice

- Podmínky aplikace MNČ:
 - lineární nezávislost sloupců matice \mathbf{X}
 - matice $\mathbf{E}(\mathbf{X}^T \mathbf{X})$ není singulární,
 - existuje její nenulový determinant,
 - lze spočítat odhadovou funkci $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Porušení vede ke kolinearitě, resp. **multikolinearitě**:
 - **Multikolinearita**: existence více vztahů lineární závislosti mezi pozorováními vysvětlujících proměnných
 - multikolinearita silná (statisticky významná)
 - multikolinearita perfektní (přítomnost umělých proměnných)
 - **Kolinearita**: existence pouze jednoho lineárního vztahu (obvykle hovoříme o multikolinearitě i když jde pouze o kolinearitu).
- Zkoumání multikolinearity: intenzita závislosti mezi dvěma nebo více vysvětlujícími proměnnými (je či není multikolinearita únosná?)

Multikolinearita – příčiny a důsledky

Příčiny

- tendence časových řad ekonomických ukazatelů (makroúdatů) vyvíjet se stejným směrem (např. HDP, C, I, S, Ex, Im)
- použití průřezových dat
- zahrnutí zpožděné endogenní nebo exogenní proměnné
- chybné použití dummies (past umělých proměnných)

Důsledky

- snížená přesnost odhadů regresních koeficientů
- velké standardní chyby odhadové funkce MNČ
 - pochybnosti či nejistota pokud jde o správnost specifikace modelu
- odhady zůstávají nestranné, vydatné
- velká citlivost odhadové funkce MNČ na velmi malé změny v matici X
- obtížné vyjádření odděleného působení silně kolineárních proměnných

Měření multikolinearity – metoda I

- použití párových korelačních koeficientů
- pro pouze 2 vysvětlující proměnné:

$$r_{x_1, x_2} = \frac{\text{cov}(x_1, x_2)}{s_{x_1} s_{x_2}} \in \langle -1, 1 \rangle$$

- multikolinearita je únosná, pokud:

$$r_{x_1, x_2} \leq 0,9 \quad \text{a současně}$$

$$r_{x_1, x_2}^2 < R_{(y, x_1, x_2)}^2 \rightarrow \text{koeficient vícenásobné determinace modelu}$$

- modul PcGive \rightarrow *Package* \rightarrow Descriptive Statistics
 \rightarrow *Model* \rightarrow Formulate $\rightarrow x_1, x_2$,
- zvolit nabídku korelační matice

Měření multikolinearity – příklad na metodu I

Soubor: CV3_PR1.xls

Data: y = maloobchodní obrat potřeb pro domácnost v mld. CZK
 x_1 = disponibilní příjem v mld. CZK
 x_2 = cenový index

Zadání: Odhadněte závislost maloobchodního obratu (y) na disponibilním příjmu (x_1) a cenovém indexu (x_2).
Vyhodnoťte multikolinearitu.

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + u_i, \quad i = 1, 2, \dots, 8$$

Měření multikolinearity – metoda II

- více vysvětlujících proměnných (tj. nestačí zkoumat párové korelační koeficienty)
- využívá se pomocné regrese a koeficientů R_i^2
- $y = f(x_1, x_2, x_3) \rightarrow \text{z modelu} \rightarrow R^2$
 - $x_1 = f(x_2, x_3) \rightarrow R_1^2$
 - $x_2 = f(x_1, x_3) \rightarrow R_2^2$
 - $x_3 = f(x_1, x_2) \rightarrow R_3^2$
- jsou-li všechna $R_i^2 < R^2$, pak je multikolinearita únosná

Měření multikolinearity – příklad na metodu II

Soubor: CV5_PR6.xls

Data: y = počet prodaných kuřat (v desítkách milionů kusů)
 x_1 = výše dotace do zemědělství (v mld. Kč)
 x_2 = cena za kuře (Kč/kg)
 x_3 = cena vepřového (Kč/kg)

Zadání: Odhadněte závislost počtu prodaných kuřat (y) na proměnných x_1 , x_2 a x_3 .

Vyhodnoťte multikolinearitu.

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \beta_3 x_{3t} + u_t, \quad t = 1, 2, \dots, 23$$

Měření multikolinearity – příklad na metodu I a II

Soubor: CV5_PR7.xls

Data: *HDP* = hrubý domácí produkt (v mld. Kč)
C = spotřeba (v mld. Kč)
I = investice (v mld. Kč)
G = vládní výdaje (v mld. Kč)

Zadání – metoda I:

Odhadněte závislost *HDP* na proměnných *C* a *I*. Vyhodnoťte multikolinearitě.

$$HDP_t = \beta_0 + \beta_1 C_t + \beta_2 I_t + u_t \quad t = 1, 2, \dots, 19$$

Zadání – metoda II:

Odhadněte závislost *HDP* na proměnných *C*, *I* a *G*. Vyhodnoťte multikolinearitě.

$$HDP_t = \beta_0 + \beta_1 C_t + \beta_2 I_t + \beta_3 G_t + u_t \quad t = 1, 2, \dots, 19$$

Možnosti řešení neúnosné multikolinearity v LRM:

- získání dalších pozorování
- snížení počtu exogenních proměnných
- použití jiného modelu
- použití jiné odhadové techniky
- transformace pozorování
 - první difference - pozor na autokorelaci
 - poměrové veličiny - pozor na heteroskedasticitu