

Primer primerjave analize kontingenčnih tabel z modeli regresijske analize nominalnih spremenljivk

Cveto Trampuž

An Illustrative Comparison Logit Analysis with Dummy Variable Regression Analysis. Two different regression models in which the dependent variable and all the explanatory variables are nominal are presented. The comparison is made at an elementary technical level on the three variables (age, employment, opinion on possible shortening of weekly working hours) from the Public Opinion Survey 89⁷ in Slovenia.

1. Uvod

1.1 Predstavitev podatkov v obliki večrazsežnostnih frekvenčnih tabel in pojem modela analize.

V empiričnih raziskavah, predvsem v družboslovju, moramo velikokrat analizirati odnose med lastnostmi opazovanih enot, ki so opisane z nominalnimi spremenljivkami. Pri tem nastane vprašanje, kakšne in kako zanesljive informacije nam podatki nudijo, ter kako jih lahko dobimo, predstavimo in ocenimo.

Prvi korak je prav gotovo izračun eno in večrazsežnostnih frekvenčnih porazdelitev (v nadalnjem besedilu jih bomo imenovali kar tabele) za tiste spremenljivke, ki nas zanimajo. Če ne znamo nič drugega, so tako dobljene frekvence, skupaj z vsemi možnimi strukturnimi odstotki, edine informacije za zaključke o odnosih med spremenljivkami oziroma o odnosih med enotami. V tabelah so ohranjene vse informacije, ki jih podatki vsebujejo.

V naslednjem koraku nas zanima, ali je možno objektivneje oceniti zanesljivost rezultatov analiz in s tem naših odločitev. To vprašanje je še posebej pomembno zato, ker so v veliki večini raziskav opazovane enote le vzorec, vzet iz neke populacije. Radi bi vedeli, s kakšno zanesljivostjo lahko sklepamo iz vzorčnih podatkov na razmerek v celotni populaciji.

Problem rešujemo tako, da na osnovi poznavanja problematike, ki jo raziskujemo, predpostavimo določene odnose med spremenljivkami in preizkusimo, če dani podatki naše predpostavke dovolj dobro potrjujejo. Zamisel mora biti torej taka, da dopušča tudi možnosti ocenjevanja in pospoljevanja. Pravimo, da si zamislimo model odnosov med spremenljivkami oziroma enotami. Model mora biti tak, da dopušča ocenitev tako količine informacije, ki je prenesena v model (kadar podatke prikazujemo v bolj strnjeni obliki,

kot so tabele, lahko izgubimo del informacije, ki jo sicer vsebujejo), kot tudi parametre modela, ki nam v njegovem okviru ponazarjajo predpostavljene odnose.

1.2 Opis izbranih spremenljivk.

Splošna pričerjava analize nominalnih spremenljivk z različimi regresijskimi modeli bi bila v okviru tega prispevka preobsežna, zato prikazujemo primerjavo le pri analizi odnosov med tremi konkretnimi spremenljivkami z linearno analizo in regresijsko analizo nominalnih spremenljivk. Pri tem nam ne gre toliko za vsebinski pomen dobljenih rezultatov, kot za pot, po kateri ugotavljamo medsebojne odnose med nominalnimi spremenljivkami z izbranimi postopki. V bolj obsežnih primerih, ko analiziramo več spremenljivk z več vrednostmi, so postopki podobni, kot bodo opisani v naslednjih poglavjih.

Za prikaz in tudi za razlaganje primerjav smo izbrali naslednje tri spremenljivke:

Z zaposlenost vrednosti: (1=da, 0=ne).

S starost vrednosti: (1=do 30, 2=31-50, 3=nad 50)

Pri starosti smo tvorili še naslednje dihotomne spremenljivke:

S_1 ima vrednost 1, če je $S = 1$ in 0 v ostalih primerih

S_2 ima vrednost 1, če je $S = 2$ in 0 v ostalih primerih

S_3 ima vrednost 1, če je $S = 3$ in 0 v ostalih primerih.

T ali naj se zmanjša število delovnih ur v tednu

vrednosti: (1=ne, 0=ostalo: da+ne vem).

Naš problem je, da ugotovimo ali je mišljenje, da delovnega tedna ni potrebno skrajšati, različno glede na starost in zaposlitev in glede na interakcijo med zaposlitvijo in starostjo (vpliv interakcije spremenljivk S in Z na spremenljivko T si predstavljamo tako, da $S(Z)$ spremeni svoj vpliv na T pri različnih vrednostih spremenljivke $Z(S)$).

2. Nekatere metode za analizo odnosov med nominalnimi spremenljivkami.

Za analizo nominalnih spremenljivk imamo na razpolago veliko metod. Med najpogosteje uporabljenimi so tabele in statistika χ^2 . Če je spremenljivk več in če imajo več vrednosti, postane tak način analize nepregleden in naporen. Kadar odnosi niso enostavni, se namreč težko odločamo, kako potekajo spremembe frekvenc 'po eni spremenljivki vzdolž vrednosti ostalih spremenljivk' in kakšne zaključke lahko naredimo o povezanosti med temi spremenljivkami. Zato je smiseln izbrati še postopke, ki omogočajo splošnejši in celovitejši pregled informacij, ki so v podatkih na razpolago. Splošne ugotovitve lahko dopolnimo še s podrobnostmi, ki jih ugotavljamo z bolj enostavnimi postopki. Pri tem pa naše ugotovitve ne smejo biti različne na osnovi različnih metod analize. Če bi se to zgodilo, bi pomenilo, da metode niso enakopravne za našo raziskavo, da niso enako občutljive glede na dane podatke ali pa, da so naše ugotovitve na meji zanesljivosti, kakršnokoli smo si že zastavili.

V našem zgledu smo preizkusili tri možnosti.

2.1 Večrazsežnostna tabela, prikazana v dveh razsežnostih.

V celu tabeli je spremenljivka T , v glavi pa spremenljivka $ZS1S2$, ki predstavlja vse možne kombinacije vrednosti spremenljivk Z , $S1$ in $S2$.

		ZS1S2						
T		000	001	010	100	101	110	
0	f_{00}	f_{12}	f_{13}	f_{11}	f_{15}	f_{16}	f_{17}	
1	f_{21}	f_{22}	f_{23}	f_{21}	f_{25}	f_{26}	f_{27}	
	$f_{..1}$	$f_{..2}$	$f_{..3}$	$f_{..1}$	$f_{..5}$	$f_{..6}$	$f_{..7}$	

Tabela 1:

V zgornji tabeli so podani osnovni podatki za raziskovanje medsebojnih odnosov med navedenimi spremenljivkami s pomočjo navedenih regresijskih modelov. Vendar bi radi še preizkusili, če je možno s standardno analizo take tabele dobiti smiselno podobne ugotovitve, kot pri regresijski analizi.

Zanima nas, če sta spremenljivki T in $ZS1S2$ povezani. Nalogo lahko formalno rešimo takole:

Postavimo ničelno hipotezo H_0 : *spremenljivki nista povezani* in alternativno hipotezo H_1 : *spremenljivki sta povezani*. Hipotezo H_0 bomo testirali s pomočjo testa χ^2 . V primeru, da bomo ničelno hipotezo zavrnili in sprejeli z določenim tveganjem alternativno hipotezo, bomo morali ugotovljeno povezanost še razložiti. Pri tem si bomo pomagali s takoimenovanimi standardiziranimi (r'_{ij}) in popravljenimi reziduali (r''_{ij}), ki so porazdeljeni pri velikem številu enot približno po standardizirani normalni porazdelitvi [Hab78]. Da pa lahko test χ^2 uporabimo, morajo biti teoretične frekvence (F_{ij}) v vseh celicah tabele (1) vsaj 5. Poleg tega je še odprto vprašanje, če je uporaba statistike χ^2 v vseh primerih takih tabel ustrezna.

2.2 Regresijski model, ko so odvisna in vse neodvisne spremenljivke dihotomne.

Model ima naslednjo obliko:

$$T = b_0 + b_1 Z + b_2 S1 + b_3 S2 + b_4 ZS1 + b_5 ZS2 + \epsilon. \quad (1)$$

Pomen oznak:

$ZS1$ pomeni interakcijo med Z in $S1$ (produkt $Z \times S1$).

$ZS2$ pomeni interakcijo med Z in $S2$ (produkt $Z \times S2$).

b_i je odvisna spremenljivka;

Z , $S1$, $S2$, $ZS1$, $ZS2$ so neodvisne spremenljivke;

b_0 je regresijska konstanta;

b_1, b_2, b_3, b_4, b_5 so parametri ozitoma koeficienti regresijskega modela

ϵ je člen napake.

Pomen koeficientov b_p ($p = 0, \dots, 5$) si lahko razlagamo tudi takole: Vrednost spremenljivke T , pri danih vrednostih neodvisnih spremenljivk, pomeni delež enot, ki pripadajo razredu, ki ga določajo vrednosti neodvisnih spremenljivk in vrednost 1 pri T . Člen napake mora zaradi statističnega ocenjevanja parametrov modela zadoščati posebnim pogojem [Gur86]. Zanima nas, od katerih neodvisnih spremenljivk je T odvisna in kako lahko z danimi vrednostmi neodvisnih spremenljivk ocenimo vrednost spremenljivke T .

Formalno bomo postopali takole:

Najprej postavimo ničelno hipotezo H_0 za celoten model. $H_0: b_j = 0$ za vsak j ($j = 1, \dots, 5$), ali še enostavnejše $H_0: R = 0$, kjer je R multipli korelacijski koeficient. Če bomo lahko na osnovi testa F ničelno hipotezo zavrnili, bomo lahko z določenim tveganjem sprejeli alternativno hipotezo H_1 , namreč, da je odvisna spremenljivka odvisna vsaj od ene izmed neodvisnih spremenljivk, oziroma, da je vsaj en b_j različen od nič ali pa, da je $R > 0$.

Poleg tega nas seveda zanima, kako je odvisna spremenljivka odvisna od posameznih neodvisnih spremenljivk. Zato postavimo še nadaljnje hipoteze za vsako posamezno neodvisno spremenljivko: $H_{0j}: b_j = 0$ in $H_1: b_j \neq 0$ ali $b_j > 0$ ali $b_j < 0$. H_0 testiramo s pomočjo testa F, H_{0j} pa s pomočjo testov T ali F.

V modelu (1) lahko obravnavamo interakcije tako, da imajo spremenljivke $ZS1$ in $ZS2$ vrednost 1, če sta vrednosti Z in $S1(S2)$ enaki, in vrednost 0 v ostalih primerih:

$$\begin{aligned} ZS1 &= (Z)(S1) + (1 - Z)(1 - S1) \\ ZS2 &= (Z)(S2) + (1 - Z)(1 - S2) \end{aligned} \quad (2)$$

Koeficienti b_p ($p = 0, \dots, 5$) imajo lahko tudi v tem primeru že opisani pomen.

Enake rezultate (razen b_0), kot v zgornjem primeru dobimo, če sprememimo vrednost 0 pri vseh spremenljivkah v vrednost -1. Ocene vrednosti spremenljivke T iz modela (1) pa v tem primeru pomenijo razlike deležev enot med razredoma, ki ju določajo iste vrednosti neodvisnih spremenljivk in vrednosti 1 oziroma -1 spremenljivke T .

2.3 Loglinearni modeli.

Z loglinearnimi modeli analiziramo pri večrazsežnostnih tabelah frekvence in njihova razmerja kot funkcije parametrov, ki na osnovi modela pojasnjujejo karakteristike spremenljivk in njihove medsebojne odnose. Pa tudi obratno je res. Parametri so funkcije teoretičnih frekvenc oziroma njihovih razmerij. Ocenimo jih iz dejanskih (empiričnih) frekvenc. Predpostavimo, da v tabeli nobena empirična frekvanca ni enaka 0, čeprav obstojajo posebne tehnike v okviru loglinearnih modelov, ki bolj ali manj dobro rešujejo tudi ta problem. Model izberemo podobno, kot pri regresijski analizi. Tudi loglinearni model predstavlja naše teoretične predpostavke o odnosih med spremenljivkami. Na osnovi modela izračunamo poleg parametrov lambda (λ) še teoretične frekvence. Če se teoretične frekvence dovolj dobro ujemajo z empiričnimi, lahko sklepamo, da so bile naše predpostavke o povezanosti spremenljivk z določenim tveganjem pravilne. Oceniti skušamo tudi zanesljivost informacij, ki nam jih posredujejo parametri. Opišimo celoten postopek nekoliko bolj podrobno.

Najprej si ogrejmo tako imenovani 'hierarhični saturirani' loglinearni model:

$$\log(F_{ijk}) = \lambda_0 + \lambda_i^T + \lambda_j^Z + \lambda_k^S + \lambda_{ij}^{TZ} + \lambda_{ik}^{TS} + \lambda_{jk}^{ZS} + \lambda_{ijk}^{TGS} \quad (3)$$

in še tako imenovani 'saturirani' logit model:

$$\frac{1}{2} \log\left(\frac{F_{ijk}}{F_{ijk}}\right) = \lambda + \lambda_j^Z + \lambda_k^S + \lambda_{jk}^{ZS} \quad (4)$$

Pri obeh modelih imajo indeksi naslednje vrednosti: $i = 1, 2; j = 0, 1; k = 1, 2, 3$.

Pomen označ:

F_{ijk} je teoretična frekvence, izračunana pri pogoju, da model velja.

λ_0 je logaritem geometričnega povprečja dejanskih frekvenc (f_{ijk}) vseh celic tabele.

Razlagamo si ga podobno, kot konstanto pri regresijski analizi. Parametri λ predstavljajo vplive posameznih spremenljivk ali njihovih medsebojnih interakcij na teoretične frekvence v posameznih celicah tabele.

Razmišljamo podobno, kot pri regresijski analizi. Če je λ enak nič, vpliva, ki ga parameter predstavlja, ni. Če je dovolj različen od 0, pa poenostavljeno pomeni, da je v celici empirična frekvencia prevelika (pri pozitivnih λ) ali premajhna (pri negativnih λ), da bi lahko trdili, da vpliva ni. Pri vplivih interakcij dveh ali večih spremenljivk moramo ugotoviti po lastni presoji, katera izmed njih je za to odgovorna. Pri tem moramo seveda upoštevati, da so vsote parametrov λ po ustreznih indeksih enake 0 [Hab78].

λ_i^T predstavlja vpliv spremenljivke T .

λ_j^Z predstavlja vpliv spremenljivke Z .

λ_k^S predstavlja vpliv spremenljivke S .

λ_{ij}^{TZ} predstavlja vpliv interakcije spremenljivk T in Z .

λ_{ik}^{TS} predstavlja vpliv interakcije spremenljivk T in S .

λ_{jk}^{ZS} predstavlja vpliv interakcije spremenljivk Z in S .

λ_{ijk}^{TGS} predstavlja vpliv interakcije spremenljivk T, Z in S .

λ v modelu logit razlagamo podobno, kot regresijsko konstanto.

Statistične ocene parametrov λ in rezidualev (razlik med teoretičnimi in empiričnimi frekvencami).

Predpostavimo, da nas zanimajo od 0 različni λ . Postavimo ničelno hipotezo $H_0: \lambda = 0$ in alternativno $H_1: \lambda \neq 0$. Kako veliko absolutno vrednost mora imeti λ , da ničelno hipotezo lahko zavrnemo? Vsaj približen odgovor na zastavljeno vprašanje dobimo, če izračunamo njegovo standardizirano vrednost. Označimo jo z λ' . λ' je pri velikih vzorcih porazdeljen približno po standardizirani normalni porazdelitvi. Če je njegova absolutna vrednost večja od 2.0, lahko s tveganjem, manjšim od 0.05, H_0 zavrnemo. Seveda pa ta test ni popolnoma zanesljiv, saj v splošnem ne vemo, kdaj je vzorec dovolj velik.

Druga možnost je, da izračunamo reziduale $r_{i,k} = f_{i,k} - F_{ijk}$. Za kritično oceno rezidualov izračunamo njihove standardizirane ($r'_{i,k}$) in popravljene vrednosti ($r''_{i,k}$). Oboji so porazdeljeni približno po standardizirani normalni porazdelitvi, zato jih znamo verjetnostno oceniti. Statistično ocenjevanje parametrov modela (3) s pomočjo testov Pearsonov χ^2 in L.R. (likelihood ratio) χ^2 je opisano v [Hab79].

3. Rezultati primerjav analiz.

Analiza odvisnosti spremenljivke T od S in Z s pomočjo tabele (1) je v našem primeru e'-vivalentna analizi rezultatov z regresijsko analizo (1).

To je ravno iz odnosov med frekvencami v tabeli (1) in koeficienti modela (1). Vpeljimo naslednje oznake:

$$T' = B_0 + B_1 Z + B_2 S1 + B_3 S2 + B_4 ZS1 + B_5 ZS2,$$

kjer so B_p ($p = 0, \dots, 5$) ocene za koeficiente b_p in $p_{i/j} = f_{ij}/f_j$; ($i = 1, 2$; $j = 1, \dots, 6$).

Med koeficienti B in količinami p veljajo naslednje zvezce:

$$\begin{aligned} B_0 &= p_{2/1} & p_{2/1} &= B_0 \\ B_1 &= p_{2/1} - p_{2/1} & p_{2/1} &= B_0 + B_1 \\ B_2 &= p_{2/3} - p_{2/1} & p_{2/3} &= B_0 + B_2 \\ B_3 &= p_{2/2} - p_{2/1} & p_{2/2} &= B_0 + B_3 \\ B_4 &= p_{2/6} - p_{2/3} - p_{2/1} + p_{2/1} & p_{2/6} &= B_0 + B_1 + B_2 + B_4 \\ B_5 &= p_{2/5} - p_{2/2} - p_{2/1} + p_{2/1} & p_{2/5} &= B_0 + B_1 + B_3 + B_5 \end{aligned}$$

Gornje zvezce je seveda možno posložiti za poljubno število spremenljivk za tabele (1) in modelje (1).

Zvezce med ocenami regresijskih koeficientov B_p in $p_{i/j}$ je možno dobiti na preprost način tudi v primeru, ko so produkti $ZS1$ in $ZS2$ računani po formulah (2), oziroma, ko so vrednosti 0 vseh spremenljivk spremenjene v -1.

Tudi med statistiko χ^2 , izračunano iz tabele (1), in multiplim korelačijskim koeficientom R iz (1) velja zveza $R^2 = \frac{\chi^2}{n}$ (n je število enot), kar ni težko dokazati.

S tem je dana zveza tudi med statistiko F pri regresijskih modelih (1) in χ^2 pri tabelah.

Potrebno pa bi bilo še ugotoviti, kakšni so odnosi med pogoji, ki jim morajo zadoščati podatki za uporabo testa χ^2 pri tabelah (1) in testov F pri regresijskih modelih (1). Vprašanje, ali so tabele in ustrezni regresijski modeli v splošnem ekvivalentni, ostaja še odprto.

Rezultati analiz odvisnosti spremenljivke T od spremenljivk Z in S s pomočjo večrazsežnostnih frevenčnih porazdelitev in z regresijsko analizo nominalnih spremenljivk¹.

Opis spremenljivk za večrazsežnostno tabelo (2) je podan na strani 2. Spremenljivko $ZS1S2$ smo izračunali po formuli $ZS1S2 = 100Z + 10S1 + S2$. Vrednosti spremenljivke $ZS1S2$ bi lahko določili (šifrirali) tudi drugače, na primer s števili 1,2,...,6. Predlagano formulo smo izbrali le zato, da so vrednosti že kar kombinacije vrednosti spremenljivk Z in S .

Test χ^2 nam s tveganjem $sig\chi^2 = 0.0003$ omogoča sprejetje hipoteze, da sta spremenljivki T in $ZS1S2$ povezani. Reziduali (Std Res in Adj Res) nam kažejo, kako lahko

¹Vsi rezultati na tej in naslednjih straneh so bili dobljeni s programi CROSSTABS, REGRESSION in HILOGLINEAR v okviru SPSSPC+ in LOGLINEAR v okviru SPSS-X.

Crosstabulation: T Zmanjsati st. del. ur na teden

By ZS1S2 Zaposlenost-Starost

Count I

Exp Val I

Row Pct I

Col Pct I

Tot Pct I

Residuals

Std Res I S2=0 I S2=1 I S2=0 I S2=0 I S2=1 I S2=0 I

Adj Res I 000 I 001 I 010 I 100 I 101 I 110 I Total

-----+-----+-----+-----+-----+-----+-----+-----+

0	I	222	I	67	I	42	I	48	I	317	I	162	I	858
	I	189.7	I	53.4	I	39.5	I	58.0	I	342.8	I	174.5	I	42.1%
ostalo	I	25.9%	I	7.8%	I	4.9%	I	5.6%	I	36.9%	I	18.9%	I	
	I	49.2%	I	52.8%	I	44.7%	I	34.8%	I	38.9%	I	39.0%	I	
	I	10.9%	I	3.3%	I	2.1%	I	2.4%	I	15.5%	I	7.9%	I	
	I	32.3	I	13.6	I	2.5	I	-10.0	I	-25.8	I	-12.5	I	
	I	2.3	I	1.9	I	.4	I	-1.3	I	-1.4	I	.3	I	
	I	3.5	I	2.5	I	.5	I	-1.6	I	-1.4	I	-1.4	I	
	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	
1	I	229	I	60	I	52	I	60	I	498	I	253	I	1182
	I	261.3	I	73.6	I	54.5	I	80.0	I	472.2	I	240.6	I	57.9%
ce	I	19.4%	I	5.1%	I	4.4%	I	7.6%	I	42.1%	I	21.4%	I	
	I	50.8%	I	47.2%	I	55.3%	I	65.2%	I	61.1%	I	51.0%	I	
	I	11.2%	I	2.9%	I	2.5%	I	4.4%	I	24.4%	I	12.4%	I	
	I	-32.3	I	-13.6	I	-2.5	I	10.0	I	25.8	I	-2.5	I	
	I	-2.0	I	-1.6	I	-.3	I	1.1	I	1.2	I	.3	I	
	I	-3.5	I	-2.5	I	-.5	I	1.8	I	2.9	I	1.4	I	
	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	
Column		451		127		94		138		815		415		2640
Total		22.1%		6.2%		4.6%		6.8%		40.0%		20.3%		100.0%

Chi-Square	D.F.	Significance	Min E.P.	Cells with E.P. < 5
23.62992	5	.0003	39.535	None

Tabela 2: Večrazsežnostna frekvenčna tabela.

Dependent Variable.. T Zmanjšati st. del. ur na teden

Multiple R	.10763	Analysis of Variance			
R Square	.01158	DF	Sum of Squares	Mean Square	
Adjusted R Square	.00915	Regression	5	5.75846	1.15169
Standard Error	.49151	Residual	2034	491.37683	.24158

F = 4.76730 Signif F = .0002

----- Variables in the Equation -----

Variable	B	SE B	Beta	T	Sig T
Z	.14441	.04781	.13749	3.020	.0026
S2	-.03532	.04937	-.03567	-.715	.4745
ZS1	-.08797	.07375	-.07173	-1.193	.2331
S1	.04543	.05573	.03982	.815	.4150
ZS2	-.00581	.06697	-.00577	-.087	.9309
(Constant)	.50776	.02314		21.939	.0000

Tabela 3: Regresijski model

povezanost na hitro razložimo. Največje vrednosti so v stolpcih 000, 001 in 101. Tisti, ki niso zaposleni ($Z = 0$) in so v starostni skupini $S3$ ($S1 = 0$ in $S2 = 0$) in delno v $S2$ se manj odločajo, da delovni teden ne bi bil skrajšan. Obratno pa velja za zaposlene ($Z = 1$) predvsem v starostnem razredu $S2 = 1$. Tabelo lahko seveda še natančneje raziščemo po standardnih postopkih.

V tabeli (3) je prikazan le zadnji korak ‘forward’ analize, ko so v model vključene vse neodvisne spremenljivke.

Iz rezultatov regresijske analize lahko opazimo, da je T odvisna le od Z ($sigT = 0.0026$) in da konstanta, ki pomeni skupino $Z = 0$, $S = 3$ v veliki meri prispeva k večji vrednosti T , kar je v skladu z ugotovitvami pri analizi tabele.

3.1 Na osnovi rezultatov, dobljenih iz regresijske analize (1) in loglinearnih modelov (3) in (4), ugotavljamo v našem konkretnem primeru enake odnose med spremenljivkami T , Z in S .

Žvezne med parametri λ pri loglinearnih modelih in parametri B pri regresijskem modelu niso linearne. Zato je njihov prikaz težji in ga tu le nakazujemo brez poglobljene analize.

Za vse ocene parametrov λ v loglinearnih modelih (3) in (4) veljata formuli (5) in (6)

{vrednosti spremenljivk T in Z smo zato, da bi se ujemali z indeksi v formulah, povečali za eno);

$$\lambda_{ijk}^{TZS} = \frac{1}{tzs} \sum_{i=1}^t \sum_{j=1}^z \sum_{k=1}^s (t\delta_{ij} - 1)(z\delta_{pj} - 1)(s\delta_{jk} - 1)\mu_{ijk} \quad (5)$$

$$\lambda_{ijk'}^{TZS} = \frac{1}{tzs} \sum_{i=1}^t \sum_{j=1}^z \sum_{k=1}^s (z\delta_{ij} - 1)(s\delta_{jk} - 1) \frac{\mu_{ijk}}{\mu_{ijk'}} \quad (6)$$

Oznake imajo naslednji pomen:

λ_{ijk}^{TZS} je oznaka za parametre λ v modelih (3) in (4), indeks pri λ dobimo tako, da upoštevamo vse smiselne kombinacije (glede na model) indekov i' , j' in k' prvega, drugega in tretjega reda. Ko računamo na primer, λ_{ijk}^{TZS} si mislimo, da indeksov j' in k' in konstant z in s v vsoti formule (5) nij.

$$\lambda_{ijk}^{TZS} = \frac{1}{tzs} \sum_{i=1}^t \sum_{j=1}^z \sum_k (t\delta_{ij} - 1)\mu_{ijk}$$

Podobno reduciramo formulo (5) za računanje na primer

$$\lambda_{ijk'}^{TZS} = \frac{1}{tzs} \sum_{i=1}^t \sum_j \sum_k (z\delta_{ij} - 1)(s\delta_{jk} - 1)\mu_{ijk}$$

α_{ijk} so naravni logaritmi empiričnih frekvenc, $\mu_{ijk} = \ln(f_{ijk})$.

Indeks i imajo naslednje vrednosti:

$i, i' = 1, \dots, t$ (v našem primeru je $t = 2$);

$j, j' = 1, \dots, z$ (v našem primeru je $z = 2$);

$k, k' = 1, \dots, s$ (v našem primeru je $s = 3$).

δ_{ij} so simboli, ki imajo vrednost 1, če sta indeksi i in i' enaka ($i = i'$) in vrednost 0, če velja $i \neq i'$

Rezultati analize odvisnosti spremenljivke T od spremenljivk Z in S s pomočjo loglinearnih modelov.

Hierarhični popolni (saturirani) loglinearni model.

Rezultati, ki jih dobimo s tem modelom, nam omogočajo ugotoviti, kateri λ v modelu (3) so dovolj različni od 0. Kaj pomenijo vsi rezultati, ki jih izpisuje računalniški program (ilogilnear – backward elimination), je možno prebrati v SPSS-X Advanced Statistics Guide, itd. Izmed njih naj omenimo le naslednje ugotovitve:

Končni hierarhični model je na osnovi testa χ^2 izločen s $T * Z$ in $Z * S$ in ne s $T * Z * S$ ter $T * S$. Če upoštevamo rezultate logit modela (4), kjer kombinacija $Z * S$ ne nastopa, $T * Z * S$ pa zaradi gornje ugotovitve ne pride v pesteve. Lahko igotovimo, da je mnenje, da delovnega tedna ni potrebno skrajšati, odvisno predvsem od zaposlitve. Tisti, ki so zaposleni, razmišljajo, da delovnega tedna ni potrebno skrajšati, kar se ujema z ugotovitvami, ki smo jih dobili na osnovi rezultatov regresijske analize. Seveda pa obstaja velika povezanost med starostjo in zaposlitvijo, kar je trivialna ugotovitev.

Na tem mestu ponovno poudarjammo, da nam v tem članku ne gre za vsebinsko raznščlanjanje o odnosih med spremenljivkami T , Z in S (zvezo med njimi bi lahko raziskali še z drugimi metodami (ne da bi spremenljivke razvrstili v razrede), temveč le za preizkus možnosti primerjav analiz s pomočjo tabel (1), regresijske analize (1) in loglinearnih modelov (3) in (4).

V nadaljnjem bi bilo potrebno raziskati (a to presega okvir tega članka), kdaj se lahko ugotovitve o odnosih med spremenljivkami razlikujejo z uporabo loglinearnih in regresijskih modelov. [Goo76] navaja, da nastajajo razlike takrat, kadar so pogojne verjetnosti v (1) izven intervala 0.25 – 0.75. V našem primeru ta pogoj ni bil izpolnjen.

		STAROST											
		do 30 let			od 31 do 50 let			nad 50 let			SKUPAJ		
		ZAPOSLENI			ZAPOSLENI			ZAPOSLENI			ZAPOSLENI		
		da	ne	\sum	da	ne	\sum	da	ne	\sum	da	ne	\sum
da	I	253	52	305	198	60	258	90	229	319	811	311	1122
ne	I	-0,057	0,057	0,026	0,025	-0,025	-0,052	0,032	-0,032	0,026	0,116	-0,166	0,31
	F	1,280	1,280	0,588	0,603	-0,603	-1,283				3,850	-3,850	4,23
da+	I	612	12	204	317	67	381	48	222	270	527	311	858
ne	I	0,057	-0,057	-0,026	-0,025	0,025	0,052	-0,032	0,032	-0,026	-0,116	0,116	-0,11
vem	I	1,280	-1,280	-0,588	-0,603	0,603	1,283				-3,850	3,850	-1,59
Σ	I	115	94	509	815	127	942	138	451	589	1363	672	2040
	F	0,386	-0,386	-0,237	0,571	-0,571	0,251	-0,957	0,957	-0,014	0,343	-0,343	
	F	8,628	-8,628	-5,298	11,01	-11,01	6,153	-22,61	22,61		11,11	-11,11	

Tabela 4: Rezultati analize modela (3).

Pomen oznak v tabeli (4):

f empirična frekvence,

l parameter lambda pri modelu (3),

l' standardizirani lambda.

Ustrezeni parametri l pri modelu (4) so enaki parametrom v modelu (3).

Literatura

- [Agr84] Alan Agresti. *Analysis of Ordinal Categorical Data*. John Wiley & Sons, New York, 1984.
- [BFH75] Yvonne M. M. Bishop, Stephen E. Fienberg, and Paul W. Holland. *Discrete Multivariate Analysis: Theory and Practice*. MIT Press, Massachusetts, London, 1975.
- [Fie77] Stephen E. Fienberg. *The Analysis of Cross-Classified Categorical Data*. MIT Press, Massachusetts, London, 1977.
- [Gil81] G. Nigel Gilbert. *Modelling Society. An introduction to Loglinear Analysis for Social Researches*. George Allen & Unwin, London, 1981.
- [Goo76] Leo A. Goodman. The relationship between modified and usual multiple-regression approaches to the analysis of dichotomous variables. *Sociological Methodology*, 83–110, 1976.
- [Goo78] Leo A. Goodman. *Analyzing Qualitative/Categorical Data. Log-linear Models and Latent Structure Analysis*. Abt Books, 1978.

- [Gur86] Damodar Gurajati. *Basic Econometrics*. International student edition. McGraw-Hill Co, Singapore, 1986.
- [Hab78] Shelby J. Haberman. *Analysis of Qualitative Data*. Volume 1, Academic press, New York, 1978.
- [Hab79] Shelby J. Haberman. *Analysis of Qualitative Data*. Volume 2, Academic press, New York, 1979.
- [Kno80] David Knoke and J. Peter Burke. *Log-linear Models*. Sage University Paper 20, Sage Publications, London, 1980.
- [Mom88] Konstantin Momirović. *Uvod u analizu nominalnih varijabli*. Volume 2 of *Metodološke sveske*, FSPN, Ljubljana, 1988.
- [Nor88] J. Marija Norusis. *SPSS-X Advanced Statistics Guide*. SPSS Inc.Chicago, Cambridge,Massachusetts, 1988.
- [Upt70] Graham J. G. Upton. *The Analysis of Cress-tabulated Data*. John Wiley & Sons, New York, 1970.