

# MSP projekt     Timotej Ponek, xponek00

## Úkol 1

Tabuľka zo zadania:

	Praha	Brno	Znojmo	Tišňov	Rokytnice nad Jizerou	Jablunkov	Dolní Věstonice	okolí studenta
Zimní čas	510	324	302	257	147	66	87	14
Letní čas	352	284	185	178	87	58	65	13
střídání časů	257	178	124	78	44	33	31	1
nemá názor	208	129	70	74	6	19	32	2

Nasledující hypotézy sú testované na hladine významnosti  $\alpha = 0,05$ .

a)

Pre riešenie si vytvoríme tabuľku odhadovaných četností (ktorú som vyrátal v python notebooku skrz funkciu `chi2_contingency`, ide o 4 parameter ktorý vracia táto funkcia)

	Praha	Brno	Znojmo	Tišňov	Rokytnice nad Jizerou	Jablunkov	Dolní Věstonice	okolí studenta
Zimní čas	537.411	370.559	275.793	237.725	115.015	71.277	87.071	12.149
Letní čas	384.720	265.274	197.433	170.181	82.336	51.025	62.332	8.698
střídání časů	234.862	161.943	120.528	103.891	50.264	31.150	38.052	5.310
nemá názor	170.007	117.224	87.246	75.203	36.384	22.548	27.544	3.843

Táto tabuľka bude používaná aj ďalej v a-c

Hypotéza: V městech, obcích a v okolí studenta (8. průzkumů) je stejné procentuální zastoupení obyvatel, co preferují zimní čas.

Použijem test dobré zhody. Vypočítám si p-hodnotu (skrz funkci `chisquare`), která vyjadruje, akú veľkú oporu má naša hypotéza v pozorovaných dátach (čím nižšia hodnota, tým nižšia opora). Dá sa definovať aj ako najmenšia hladina významnosti testu, pri ktorej zamietnem nulovú hypotézu.

P-hodnotu porovnám s mojou hladinou významnosti  $\alpha$

$$p = 0.0067$$

$$p < \alpha$$

Hypotézu zamietam.

b)

Hypotéza: V městech, obcích a v okolí studenta (8. průzkumů) je stejné procentuální zastoupení obyvatel, co preferují letní čas.

Použijem test dobré zhody. Vypočítám si p-hodnotu skrz funkci `chisquare` a tú porovnám s mojou hladinou významnosti  $\alpha$

$$p = 0.3258$$

$$p > \alpha$$

Hypotézu nezamietam.

c)

Hypotéza: V městech, obcích a v okolí studenta (8. průzkumů) je stejné procentuální zastoupení obyvatel, co preferují střídání času.

Použijem test dobré zhody. Vypočítám si p-hodnotu skrz funkci `chisquare` a tú porovnám s mojou hladinou významnosti  $\alpha$

$$p = 2.4094 \cdot 10^{-5}$$

$$p < \alpha$$

Hypotézu zamietam.

d)

Pre riešenie si vytvoríme tabuľku odhadovaných četností (ktorú som vyrátal v python notebooku skrz funkci `chi2_contingency`, ide o 4 parameter ktorý vracia táto funkcia)

	větší města	menší města	obce
Zimní čas	906.979	512.957	273.065
Letní čas	647.689	366.311	195.000
střídání časů	399.114	225.725	120.161
nemá názor	288.219	163.007	86.774

Táto tabuľka bude používaná aj ďalej v d-e

Hypotéza: U větších měst, menších měst a obcí (3. průzkumy) je stejné procentuální zastoupení obyvatel, co preferují zimní čas.

Použijem test dobrej zhody. Vypočítam si p-hodnotu skrz funkciu `chisquare` a tú porovná s mojou hladinou významnosti  $\alpha$

$p = 0.0018$

$$p < \alpha$$

Hypotézu zamietam.

d)

Hypotéza: U větších měst, menších měst a obcí (3. průzkumy) je stejné procentuální zastoupení nerozhodnutelných obyvatel.

Použijem test dobrej zhody. Vypočítam si p-hodnotu skrz funkciu `chisquare` a tú porovná s mojou hladinou významnosti  $\alpha$

$$p = 3.5051 \cdot 10^{-5}$$

$$p < \alpha$$

Hypotézu zamietam.

f)

Použil som Pearsonov korelačný koeficient

hodnota pre koeficientu pre okolie študenta a veľké mestá: 0.9204

hodnota pre koeficientu pre okolie študenta a menšie mestá: 0.9085

hodnota pre koeficientu pre okolie študenta a dediny: 0.9261

Najväčšia korelácia je medzi okolím študenta a dedinou.

Môjho prieskumu sa zúčastnili moji spolužiaci zo strednej školy, ktorí sú poväčšinou z dedín a jedného menšieho mesta, takže výsledok by sedel.

## Úkol 2

a)

Zostrojím prvý model pre pôvodnú regresnú funkciu  $Z = \beta_1 + \beta_2 X + \beta_3 Y + \beta_4 X^2 + \beta_5 Y^2 + \beta_6 X*Y$ .

Použijem na to funkciu `OLS`, ktorá mi vytvorí regresný model s nasledujúcimi parametrami:

				Interval spořádlivosti	
coef	std err	t	P> t	[0.025	0.975]

$\beta_1$	-124.7148	24.310	-5.130	0.000	-173.280	-76.149
$\beta_2$	2.8798	3.793	0.759	0.450	-4.697	10.456
$\beta_3$	2.3539	7.161	0.329	0.743	-11.952	16.659
$\beta_4$	9.9022	0.169	58.542	0.000	9.564	10.240
$\beta_5$	-1.9379	0.631	-3.071	0.003	-3.198	-0.677
$\beta_6$	3.1224	0.285	10.944	0.000	2.552	3.692

hodnota  $R^2 = 0.99 \rightarrow$  model je validný

Na základe intervalov spoľahlivosti z regresného modelu odstránime parametre, ktorých interval spoľahlivosti zahŕňa hodnotu 0. Odstránime teda parametre  $\beta_2$  a  $\beta_3$ .

Vytvorím nový model pre upravenú regresnú funkciu  $Z = \beta_1 + \beta_4 X^2 + \beta_5 Y^2 + \beta_6 X*Y$ .

Funkcia `OLS`, mi vytvorí regresný model s nasledujúcimi parametrami:

	coef	std err	t	P> t	Interval spoľahlivosti	
					[0.025	0.975]
$\beta_1$	-109.1347	11.436	-9.543	0.000	-131.967	-86.302
$\beta_4$	10.0118	0.072	138.671	0.000	9.868	10.156
$\beta_5$	-1.8156	0.283	-6.413	0.000	-2.381	-1.250
$\beta_6$	3.2253	0.243	13.291	0.000	2.741	3.710

hodnota  $R^2 = 0.99 \rightarrow$  model je validný

Odstraňovanie ďalších regresných parametrov  $\beta$  by nevedlo k nájdeniu nového lepšieho modelu, keďže interval spoľahlivosti žiadneho z regresných parametrov neobsahuje hodnotu 0. Preto ukončujem hľadanie nového modelu, a prehlasujem model  $Z = \beta_1 + \beta_4 X^2 + \beta_5 Y^2 + \beta_6 X*Y$  za najlepší

b)

Funkcia `OLS` vytvára regresný model a vypočítava jeho parametre metódou najmenších štvorcov, takže pre tento pod stačí do tabuľky skopírovať hodnoty ktoré sú výstupom tejto funkcie, a sú zapísané aj vyššie.

		Interval spoľahlivosti 95%	
Parameter	Odhad hodnoty parametru	0.025	0.975
$\beta_1$	-109.1347	-131.967	-86.302
$\beta_4$	10.0118	9.868	10.156
$\beta_5$	-1.8156	-2.381	-1.250
$\beta_6$	3.2253	2.741	3.710

c)

neustranný odhad rozptylu získam z druhého modelu skrz premennú `mse_resid`

neustranný odhad rozptylu = 2525.6887

d)

zvolím si 2 ľubovoľné regresné parametre z najlepšieho modelu, napr.  $\beta_4$  a  $\beta_5$  a tie budem ďalej uvažovať v c-d.

Hypotéza: Zvolené dva regresní parametry jsou současně nulové.

Použijem f-test,  $\alpha = 0.05$

p-hodnota =  $9.5423 \cdot 10^{-90}$

$p < \alpha$

Hypotézu zamietam.

e)

Hypotéza: Zvolené dva regresní parametry jsou stejné.

Použijem t-test,  $\alpha = 0.05$

p-hodnota =  $3.5725 \cdot 10^{-53}$

$p < \alpha$

Hypotézu zamietam.