

Statistické zpracování dat 6.přednáška

Mgr. Radmila Krkošková, Ph.D.



**SLEZSKÁ
UNIVERZITA**

OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ



EVROPSKÁ UNIE
Evropské strukturální a investiční fondy
Operační program Výzkum, vývoj a vzdělávání



MINISTERSTVO ŠKOLSTVÍ,
MLÁDEŽE A TĚLOVÝCHOVY



Vícenásobná lineární regresní analýza (1)



Obsah přednášky



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

- **Vícerozměrná (vícenásobná, mnohorozměrná, mnohonásobná, n-rozměrná, n-násobná) lineární regresní analýza**
- **Populační a výběrová regresní funkce**
- **Přiléhavost regresní nadroviny k datům**
- **Koeficient determinace R^2**
- **Klasický vícerozměrný lineární regresní model**



Obsah přednášky



**SLEZSKÁ
UNIVERZITA**
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

- **Multikolinearita**
- **Heteroskedasticita (H-S)**
- **Testy H-S a její odstraňování**
- **Autokorelace**
- **Nominální proměnné**



Příklad 1.



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Zajímají nás:

Y - tržby prodejny (např. spotřební elektroniky OK)

v závislosti na:

X_1 - výdaje na reklamu

X_2 - počet kolemjdoucích

X_3 - průměrný plat prodavačů

X_4 – počet konkurenčních prodejen v místě



Příklad 2.



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Zajímá nás:

Y – dětská úmrtnost /v promile/

v závislosti na:

X_1 – gramotnost žen /v procentech/

X_2 – HDP na hlavu /v USD/

X_3 – porodnost /v procentech/



Vícenásobná regresní analýzy



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

- Grafické znázornění v dimenzích $m > 2$ - obtížné(?)
- Jediné kritérium = závislá proměnná: Y
- Více prediktorů = nez. prom.: X_1, X_2, \dots, X_m ($m = 2, 3, \dots$)

regresní *nadrovina*: $Y = B_0 + B_1X_1 + B_2X_2 + \dots + B_mX_m$

regresní model: $Y_i = B_0 + B_1X_{i1} + B_2X_{i2} + \dots + B_mX_{im} + u_i$

Cíl: nalezení **nejlepších** odhadů regresních koeficientů (Excel)



Aplikace regresní analýzy



**SLEZSKÁ
UNIVERZITA**
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

- Prognózování (tržeb, nákladů, poptávky aj.)
- Rozhodování o umístění provozoven
- Analýza marketingového mixu
(vztahy mezi prvky 5P)
- Stanovení vztahů mezi kritérii a prediktory v případě konstantních efektů jiných prediktorů
- Odhady chybějících dat



Populační regresní funkce



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Týká se regresní závislosti v **celé populaci**.

Příklad:

Stanovte závislost **Tržeb prodejny** na **Počtu kolemjdoucích**, **Velikosti prodejny**, **Průměrného platu prodavačů**, **Přítomnost konkurence** v jistém prodejním řetězci v ČR – data za **všech** 25 prodejen ($n = 25$, $m = 2$ event. 4)



Příklad – Data – všechny prodejny řetězce Ř



**SLEZSKÁ
UNIVERZITA**
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Prodejna	Roční tržby tis. Kč	Poč. kolemjoucích	Velikost prodejny m2	Prům. plat prodavačů tis.Kč/měs.	Přítomnost konkurence v místě
1	7800	12	90	10,0	1
2	10500	20	150	17,1	0
3	5700	11	100	10,5	1
4	12000	30	180	20,8	0
5	8100	15	120	12,4	1
6	9600	17	90	15,7	1
7	12900	27	200	23,2	1
8	6600	13	100	12,1	1
9	19500	55	320	26,3	0
10	15600	45	220	24,8	0
11	11400	29	170	20,5	0
12	9000	15	145	13,8	1
13	10800	24	170	16,2	0
14	9900	22	130	15,4	1
15	7200	11	120	13,1	1
16	10560	16	140	14,6	1
17	11280	18	150	15,9	0
18	11700	20	190	20,5	0
19	12300	23	190	21,3	1
20	10320	31	170	14,3	1
21	8040	16	130	12,6	1
22	8760	19	140	14,2	1
23	10920	21	170	17,4	1
24	11940	24	160	21,1	1
25	12360	29	170	22,1	1



Populační regresní funkce

poskytuje (podmíněnou) *průměrnou hodnotu* \hat{Y} závisle proměnné Y v závislosti na hodnotě nezávisle proměnných X_1, X_2, X_3, X_4 :

$$\hat{Y} = B_0 + B_1 X_1 + B_2 X_2 + B_3 X_3 + B_4 X_4$$



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ



Populační regresní funkce + stochastický model



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

poskytuje **hodnotu** závisle proměnné **Y**
v závislosti na hodnotě nezávisle
proměnných X_1, X_2, X_3, X_4 až na náhodnou
(stochastickou) chybu (poruchu):

$$Y = B_0 + B_1 X_1 + B_2 X_2 + B_3 X_3 + B_4 X_4 + u$$

Náhodná chyba: $E(u) = 0$



Výběrová regresní funkce + stochastický model



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

V praxi *nejsou k dispozici* data z celé populace,
ale jen ze *vzorku* → výběrová regresní funkce :

$$\hat{Y} = b_0 + b_1X_1 + b_2X_2 + \dots$$

↙
Odhad $E(Y|X_1, X_2, \dots)$

$$Y = b_0 + b_1X_1 + b_2X_2 + \dots + e$$

↘
Odhad chyby - *reziduum*



Výběrová regresní funkce – otázky?



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

1. Jak získat odhady regresních koeficientů B_0, B_1 a B_2, \dots , tj. b_0, b_1 a b_2, \dots ?

Odpověď: Známa **metoda nejmenších čtverců (MNČ)**

2. Jak dobré (přesné) odhady to jsou?

Odpověď: Testy hypotéz za standardních předpokladů (5 předpokladů standardního modelu – viz dále).

$$Y = b_0 + b_1X_1 + b_2X_2 + e$$



Koeficient determinace



- Teoretický součet čtverců:
- teoretické hodnoty („na regr. nadrovině“)

$$S_T = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

Y_i – hodnoty z dat

- Reziduální součet čtverců: $S_R = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$

- Celkový součet čtverců: $S_y = \sum_{i=1}^n (Y_i - \bar{Y})^2$

- Platí vztah: $S_y = S_T + S_R$

- **Koeficient determinace** - míra variability:

$$R^2 = \frac{S_T}{S_y} = 1 - \frac{S_R}{S_y}$$

- **Pozor!** R^2 má platnost pro libovolný typ regresní funkce!



Příklad 1 – řešení v Excelu



**SLEZSKÁ
UNIVERZITA**
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Regresní statistika

Násobné R	0,969
Hodnota spolehlivosti R	0,940
Nastavená hodnota spolehlivosti	0,927
Chyba stř. hodnoty	780,552
Pozorování	25

	Koeficienty	Chyba stř. hodnoty	t stat	Hodnota P
Hranice	1642,641	932,471	1,762	0,093
Poč. kolemjdoucích/hod.	81,899	36,855	2,222	0,038
Velikost prodejny m ²	19,893	8,496	2,342	0,030
Prům. plat prodavačů/měs.	241,002	70,482	3,419	0,003
Přítomnost konkurence v místě	-171,803	399,213	-0,430	0,672

Násobné R	= R - koeficient korelace
Hodnota spolehlivosti R	= R ² - koeficient determinace
Nastavená hodnota spolehlivosti R	= R ² _{adj} - upravený koeficient determinace
Chyba stř. hodnoty	= s ² - směrodatná chyba (odhad směrodatné odchylky náhod. složky)



Příklad 1 – řešení – interpretace výsledků



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

- Kritérium: Y - tržby z prodeje (v tis.Kč/rok)
- Prediktory: X_1 - poč. kolemjdoucích
 X_2 - velikost prodejny v m²
 X_3 - průměrný plat prodavačů v tis.Kč/měs.
 X_4 - přítomnost konkurence (binární)

Regresní rovnice: $y = 1642,6 + 81,9x_1 + 19,9x_2 + 241,0x_3 - 171,8x_4$



Příklad 1 – řešení – interpretace výsledků



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Hypotézy o statistické významnosti regres. koeficientů a R^2 :

H_0 : koeficient = 0

$b_0 = 1642,6$ (p -hodnota = 0,093 $\Rightarrow H_0$ zamítáme)

$b_1 = 81,9$ (p -hodnota = 0,038 $\Rightarrow H_0$ zamítáme)

$b_2 = 19,9$ (p -hodnota = 0,030 $\Rightarrow H_0$ zamítáme)

$b_3 = 241,0$ (p -hodnota = 0,003 $\Rightarrow H_0$ zamítáme)

$b_4 = -171,8$ (p -hodnota = 0,672 $\Rightarrow H_0$ nezamítáme)

Koeficient determinace (přiléhavost): $R^2 = 0,940$

(p -hodnota = 0,005 $\Rightarrow H_0$ zamítáme)

Závěr: Přítomnost konkurence nemá na tržby prodejny vliv. Tržby nové prodejny jsou na základě modelu prognózovány ve výši 10700 tis. Kč.



Předpoklady lineárního regresního modelu



1. Střední hodnota náhodné poruchy u je 0, tj. $E(u) = 0$
2. Náhodná chyba má **normální rozdělení**, tj. $u \sim N(0, \sigma^2)$
3. Vysvětlující proměnné X_1, X_2, \dots, X_m **nejsou kolineární**
(JINAK: **MULTIKOLINEARITA**)
4. Rozptyl náhodné chyby u je konstantní - **homoskedasticita** tj.
$$\text{Var}(u) = \sigma^2$$

(JINAK: **HETEROSKEDASTICITA**)
5. Náhodné chyby u **jsou nekorelované**, tj.
$$\text{Cov}(u_i, u_j) = 0 \text{ pro } i \neq j$$

(JINAK: **AUTOKORELACE**)

Co se může stát, když některý z předpokladů není splněn?



Poznámky:



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

1. Předpoklady kromě 3. jsou stejné jako v jednoduchém lineárním regresním modelu.
2. Kolinearita znamená, že žádná vysvětlující proměnná není přesnou lineární kombinací některých ostatních vysvětlujících proměnných.

Příklad: $X_{1i} = 2X_{2i} + X_{3i}$ pro všechna $i=1,2,\dots,n$

3. Problém tzv. **multikolinearity** spočívá v tom, že některé vysvětlující proměnné jsou **téměř** kolineární (lin. kombinacemi jiných proměnných).



Multikolinearita (MK)



**SLEZSKÁ
UNIVERZITA**
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Co je to multikolinearita?

Mezi vysvětlujícími proměnnými existuje (téměř) dokonalý lineární vztah (potvrzený daty), tzv. vysoká multikolinearita (high multicollinearity).

Otázky:

Jaké jsou příčiny MK?

Je MK skutečný problém?

Jaké jsou teoretické důsledky MK?

Jaké jsou praktické důsledky MK?

Jak MK v praxi zjišťovat (měřit)?

Pokud je zjištěna MK, je ji nezbytné odstranit a když, tak jak?



Jaké jsou příčiny multikolinearity?



Y = Roční tržby tis. Kč	X1 = Poč. kolemjdoucích/hod.	X2 = Velikost prodejny m ²
7800	12	111
10500	20	145
5700	11	107
12000	30	187
8100	15	124
9600	17	132
12900	27	174
6600	13	115
19500	55	292
15600	45	250
11400	29	182
9000	15	124
10800	24	161
9900	22	153
7200	11	107
10560	16	128
11280	18	136
11700	20	145
12300	23	157
10320	31	191
8040	16	128
8760	19	140
10920	21	149
11940	24	161
12360	29	182

Příklad 1*: Roční tržby závisí na velikosti prodejny a počtu kolemjdoucích: $R^2 = 0,84$

Významnost F
1,42421E-09

	Koeficienty	Hodnota P
Hranice	9053,10	0,856
Poč. kolemjdoucích/hod.	555,85	0,872
Velikost prodejny m ²	-70,73	0,932

Statistická významnost regresních koeficientů: katastrofa!!!

Důvod: téměř perfektní kolinearita X1 a X2
 $X2 = 4 \cdot X1 + 60$



Je multikolinearita skutečný problém?



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

- Příklad perfektní MK je **patologický extrém!**
- MK může být v praxi **vysoká**, nikoliv však perfektní!
- V Příkladu 1* však z ANOVA vyplývá, že
Počet kolemjdoucích a Velikost prodejny mají společný vliv na Tržby! (Celý model je statisticky významný – *F*-test v Regrese)
- Jak měřit vysokou MK? – v případě 2 korelovaných proměnných je mírou **korelační koeficient**, v případě MK více proměnných to však neplatí!!! (viz dále)



Jaké jsou teoretické a praktické důsledky MK?



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

- MK není problémem *populace*, nýbrž je problémem *vzorku* (data ve vzorku jsou „špatně“ vybrána)
- **Jinak řečeno:** vzorek nepotvrzuje teorii závislosti vysvětlované proměnné na vysvětlujících proměnných
- Hypotéza o nulovosti regresních koeficientů se přijímá, i když ve skutečnosti (tj. v populaci) neplatí
- Intervaly spolehlivosti regres. koeficientů jsou velmi široké
- Veškeré odhady regresních koeficientů jsou citlivé na jakékoliv změny dat
- Regresní koeficienty mohou mít špatná znaménka
- Regresní funkce je nevhodná pro predikci



Jak MK v praxi zjišťovat (měřit)?



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

- Na detekci MK se **nepoužívají** statistické testy!
- MK je **problém „stupně“**, nikoliv „existence“ jako takové
- **K určování stupně MK se používají (heuristická!) pravidla:**
 1. **Vysoký koeficient determinace R^2** , přitom vysoká p -hodnota regresních koeficientů (tj. Sig.- blízka k 1)
 2. **Vysoké hodnoty párových korelací** mezi vysvětlujícími proměnnými (např. $> 0,8$)
 3. **Významné regrese některých vysvětlujících proměnných na jiných vysvětlujících proměnných**
(viz Příklad 1*: závislost X_2 na X_1)



Jak odstranit multikolinearitu?



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Neexistuje zaručená metoda, protože MK je problémem vzorku, nikoliv nutně populace, z níž vzorek pochází

Možné (doporučené) metody:

1. **Vypustit některou vysvětlující proměnnou** – pozor: nevylít s vodou z vaničky i dítě! (Ekonomický model)
2. **Pořídit nový vzorek**, eventuálně doplnit starý
3. **Promyslet znovu ekonomický a matematický model** (nebylo něco opomenuto?, zjednodušeno?,...)
4. **Transformace proměnných**, např. namísto celkové spotřeby použít spotřebu na hlavu apod.



Heteroskedasticita



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Rozptyl náhodné chyby u je konstantní, tj.

$$\text{Var}(u) = \sigma^2$$

Graficky: Hodnoty jsou rozptýleny ve stejně širokém pásu kolem regresní funkce (regresní nadroviny)

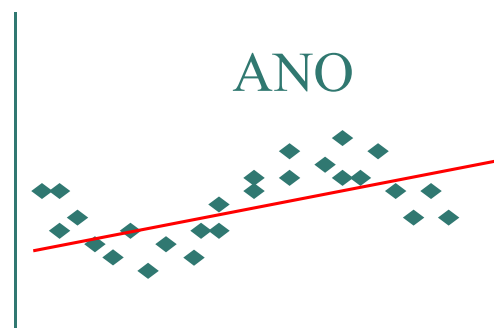
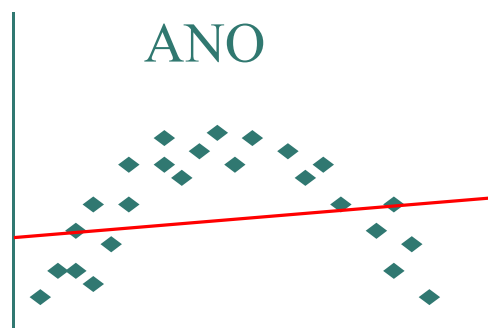
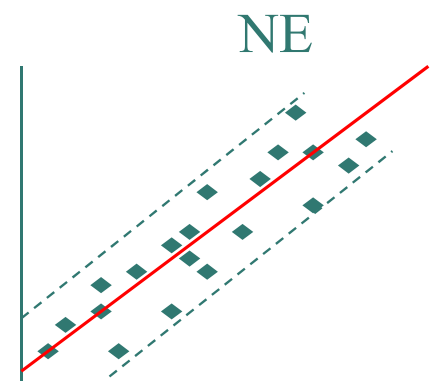
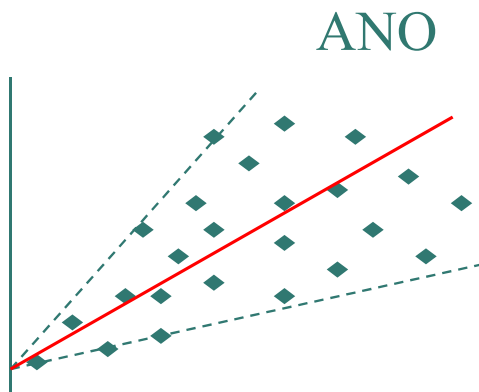
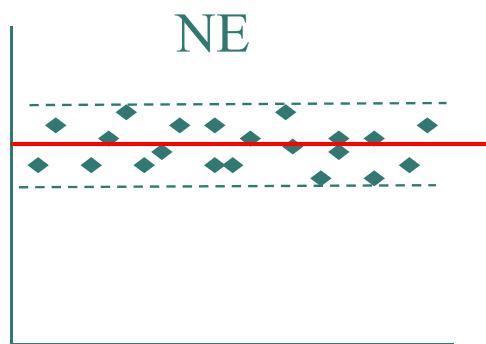
Otázky:

1. Co je podstatou heteroskedasticity (H-S)?
2. Jaké jsou důsledky H-S?
3. Jak zjišťovat H-S v dané situaci?
4. Jak odstraňovat H-S?



Jak vypadá H-S?

Grafická analýza reziduí:



Co je podstatou H-S?



**SLEZSKÁ
UNIVERZITA**
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Jedná se o rozptyl náhodné chyby u_i v regresním (populačním)

modelu, např. $Y_i = B_0 + B_1 X_{i1} + B_2 X_{i2} + u_i$



Co je podstatou H-S?



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Některé důvody nekonstantnosti rozptylu:

1. Učení se z chyb: rozptyl počtu chyb se s rostoucím časem zmenšuje
2. S rostoucím věkem roste rozptyl příjmů zaměstnanců
3. S lepšími technikami sběru dat klesá rozptyl chyb v datech



Co je podstatou H-S?



SLEZSKÁ
UNIVERZITA
OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

4. S přítomností odlehlých hodnot roste rozptyl
5. U špatně specifikovaného modelu dochází k proměnlivosti rozptylu
6. Šikmost rozdělení vysvětlujících proměnných zvětšuje rozptyl
7. Panelová (průřezová) data mívají proměnlivý rozptyl





Děkuji Vám za pozornost!!!

