



SEMINÁŘ VÝPOČETNÍ STATISTIKY

P9
2008-04-14

VÍCENÁSOBNÁ KORELACE A REGRESE:

Regrese a korelace:

- ✓ **Jednoduchá:** Y, X
- ✓ **Vícenásobná:** Y, X_1, X_2, \dots, X_k

Příklad:

Student	Hodiny	IQ	Body
1	9	99	56
2	6	100	45
3	12	119	80
4	14	95	73
5	11	110	71
6	6	117	55
7	19	98	95
8	16	101	86
9	3	100	34
10	9	115	66

Zkoumáme:

- ✓ Hledáme **regresní funkci** ve tvaru $y' = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$. Koeficientu b_0 říkáme absolutní člen (intercept), koeficienty b_1 až b_k jsou parciální (dílní) regresní koeficienty.
- ✓ **Koeficient mnohonásobné korelace**, značíme ho velkým R , interval $0 \leq R \leq 1$, stupnice stejná jako u jednoduché. Vynásobíme-li ho 100, říká, z kolika procent je veličina Y vysvětlována (ovlivňována) nezávislými veličinami X_1 až X_k
- ✓ **Koeficient mnohonásobné determinace** – R^2 .

Předpoklady použitelnosti mnohonásobné regrese a korelace:

- ✓ **Normalita** rozdělení analyzovaných proměnných.
- ✓ Vysvětlující proměnné by měly být navzájem **nezávislé**, pokud by byly závislé, opakovaly by hodnoty jedné veličiny hodnoty druhé veličiny a nedozvěděli bychom se nic nového. V reálném světě bude závislost vždy, ale chceme, aby byla co nejmenší.

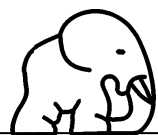
Ověření nezávislosti:

• Výpočet korelační matice:

- Obecný tvar korelační matice:

	X_1	X_2	...	X_k
X_1	1	$r_{x_1x_2}$...	$r_{x_1x_k}$
X_2		1		
...			...	
X_k				1

- Na diagonále jsou jedničky, neboť korelace mezi jednou a tou samou je maximálně silná
- Na ostatních pozicích jsou korelační koeficienty
- Pokud jsou korelační koeficienty v absolutní hodnotě menší než 0,75 ($|r_{x_jx_k}| < 0,75$), považujeme podmínku nezávislosti za splněnou.
- Je-li hodnota překročena, vzniká zde multikolinearita, předpoklad nezávislosti byl tedy porušen
- **VIF – Variance Inflation Factor**
 - Faktor zvětšení rozptylu
 - Pokud $VIF > 10$, znamená to multikolinearitu, tedy nežádoucí stav.



- ✓ **Rezidua**, tzn. rozdíly $y_i - y_i'$, $i = 1, 2, \dots, n$, by měla mít normální rozdělení s nulovou střední hodnotou a konstantním rozptylem. Nulová střední hodnota vyjadřuje, že model je umístěn natolik ideálně, že představuje osu korelačního pole.

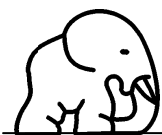
SAS:

- ✓ **Posouzení normality rozdělení analyzovaných proměnných:**

```
ods exclude Moments BasicMesures TestsForLocation Quantiles ExtremeObs;  
proc univariate data=svs normal plot;  
run;  
quit;
```
- ✓ Máme-li pouze 10 pozorování, je testování normality riskantní, testy fungují přesně až u větších datových souborů (od 30 pozorování výše). Testy většinou řeknou, že data mají normální rozdělení.
- ✓ **Výpočet korelační matice všech sledovaných proměnných:**
 - **Procedura:**

```
proc corr data=svs pearson spearman;  
run;  
quit;
```
 - **Příkazy:**
 - Pearson - spočítá tradiční korelační koeficienty, funguje, pokud je normální rozdělení a neobjevují se problematické údaje, pokud by normalita nebyla splněna, je lepší použít koeficienty pořadové korelace
 - Spearman – koeficienty pořadové korelace, Spearmanovy koeficienty, je neparametrický, takže je univerzálnější, což je zapláceno menší kvalitou (přesností).
 - Příkazy psát nemusíme, pokud je nenapišeme, vypočítají se pouze pearsonovy koeficienty.
 - Před příkaz run lze vsunout příkaz var (variables) a jména proměnných, výsledek by však byl stejný, pokud příkaz neuvedeme, jelikož automaticky se počítá pro všechny proměnné.
- ✓ **Model mnohonásobné regrese, doplněný o korelační matici všech proměnných a o regresní diagnostiku:**
 - **Procedura:**

```
proc reg data=svs corr;  
model body=hodiny IQ/r influence vif spec;  
plot r.*p;  
plot cookd.*p;  
symbol v=dot c=green;  
output out=diag r=rezid;  
run;  
quit;
```
 - **Příkazy:**
 - Procedura corr nabízí detailnější výstupy než pouze reg.
 - Model – píše se jméno vysvětlované proměnné a za rovná se se píše vysvětlující proměnné. Lomítkem naznačujeme, že budeme vypisovat další doplňující příkazy, tedy regresní diagnostiku. Ta slouží ke kontrole jednotlivých předpokladů
 - r influence – počítá rezidua a modifikace reziduí, spočítají se také studentizovaná rezidua, která využijeme k tomu, abychom posoudili, zda ve vstupních datech nejsou nějaké problémy. Bereme absolutní hodnotu studentizovaných reziduí a porovnáváme ji s hodnotou 2: $|SR| > 2$, pokud je hodnota větší, značí to odlehlé pozorování (outlier).
 - Odlehlost znamená, že údaj se vychýlil z množiny Y podstatným způsobem, ale odlehlost nemusí znamenat negativní dopad na výsledný model, je nutné zjistit, zda pozorování je vlivné, tedy zda má vliv na výsledný model \Rightarrow Cookova vzdálenost (D), $D_i > \frac{4}{n}$ - pozorování je vlivné.
 - Hledáme také odlehlou hodnotu mezi proměnnými X, charakteristika leverage umožňuje zhodnotit, zda v množině vysvětlujících proměnných není odlehlý údaj, $h_{ii} = 2 \frac{p}{n}$, p je počet parametrů regresního modelu, tedy b_0, b_1 atd., u příkladu výše jsou 3. Pokud je hranice překročena, je obsažena problematická hodnota.



- DFFITS – charakteristika nazývaná Velschova-Kuhova vzdálenost, slouží k posouzení vlivnosti. Může být kladná i záporná, bereme jí tedy v absolutní hodnotě: $|DFFITS| = 2\sqrt{\frac{p}{n}}$, pokud je hodnota překročena, je pozorování vlivné.
- Rozdíl mezi Cookovou a Velschovou-Kuhovou vzdáleností:
 - Cookova je obecnější, říká, jak je ovlivněn celý model.
 - Velschova-Kuhova vzdálenost – do jaké míry ovlivňuje pozorování jednu hodnotu veličiny Y, hodnotu, u které byla spočtena.
- Spec – počítá tzv. Whiteův test, který hodnotí jednu z důležitých vlastností reziduí, umožňuje posoudit, zda rezidua mají konstantní rozptyl.
- plot r.*p – konstrukce reziduálního grafu, lze posoudit, zda jsou vlastnosti splněny a model je kvalitní.
- plot cookd.*p – grafické zobrazení Cookových hodnot.
- symbol v=dot c=green – provedení grafů, v=dot znamená, že body budou jako tečky, barva bude zelená.
- output out=diag r=rezid – vytvoříme pomocný soubor nazvaný podle slova diagnostika „diag“, bude obsahovat jedinou proměnnou, tedy „rezid“, která je tvořena reziduálními proměnnými.

✓ **Testování normality reziduí a ověření, zda mají nulovou střední hodnotu:**

- **Procedura:**

```
ods exclude Moments Quantiles ExtremeObs;  
proc univariate data=diag normal plot;  
var rezid;  
run;  
quit;
```

- **Popis:**

- Využíváme soubor diag a zaměřujeme se na rezidua

✓ **Výstupy:**

- **Výstup:**

yscup.

Procedura CORR						
3 Proměnné: hodiny IQ body						
Jednoduché statistiky						
Proměnná	N	Průměr	Std odch	Medián	Minimum	Maximum
hodiny	10	10.50000	4.92725	10.00000	3.00000	19.00000
IQ	10	105.40000	8.90942	100.50000	95.00000	119.00000
body	10	66.10000	18.84705	68.50000	34.00000	95.00000
Pearsonovy korelační koeficienty, N = 10						
Prob > r pro H0: Rho=0						
		hodiny	IQ	body		
hodiny		1.00000	-0.23539 0.5127	0.96138 <.0001		
IQ		-0.23539 0.5127	1.00000	0.04010 0.9124		
body		0.96138 <.0001	0.04010 0.9124	1.00000		
Spearmanovy korelační koeficienty, N = 10						
Prob > r pro H0: Rho=0						
		hodiny	IQ	body		
hodiny		1.00000	-0.25382 0.4792	0.98173 <.0001		
IQ		-0.25382 0.4792	1.00000	-0.08511 0.8152		
body		0.98173 <.0001	-0.08511 0.8152	1.00000		



- **Procedura REG:**

- R-kvadrát je vysoký, bodový výsledek je popsán časem a IQ
- Individuální p-hodnoty ($Pr > t$) hodnotí zobecnitelnost každého koeficientu, všechny jsou statisticky významné
- Inlace proměnné = VIF – jsou malé, není přítomna multikolinearita
- Test první a druhé specifikace momentu – Whiteův test, p-hodnota menší než 5% => konstantnost rozptylů není

- **Další výstupy:**

