



**SLEZSKÁ
UNIVERZITA**

OBCHODNĚ PODNIKATELSKÁ
FAKULTA V KARVINĚ

Statistické zpracování dat

Distanční studijní text

Jaroslav Ramík, Radmila Krkošková

Karviná 2023

Obor: Statistika.

Klíčová slova: Analýza rozptylu, jednoduchá regresní analýza, vícerozměrná regresní analýza, analýza časových řad, ARIMA modely.

Anotace: Publikace představuje studijní oporu předmětu Statistické zpracování dat pro navazující studium na vysoké škole ekonomického zaměření. Obsahově pokrývá základní témata: analýza rozptylu – 1 faktor, analýza rozptylu – 2 faktory, jednoduchá a vícerozměrná regresní analýza, analýza časových řad.

Autor: **Prof. RNDr. Jaroslav Ramík, CSc.**
Mgr. Radmila Krkošková, Ph.D.

Obsah

ÚVODEM.....	6
RYCHLÝ NÁHLED STUDIJNÍ OPORY.....	7
1 ANALÝZA ROZPTYLU (ANOVA) – JEDEN FAKTOR.....	9
1.1 Nezávislý a závislý faktor	11
1.2 Předpoklady analýzy rozptylu s jedním faktorem.....	12
1.3 Předpoklady analýzy rozptylu s jedním faktorem.....	13
1.4 Míra těsnosti závislosti.....	15
1.5 Analýza rozptylu v programu GRETl.....	16
2 ANALÝZA ROZPTYLU (ANOVA) – DVA A VÍCE FAKTORŮ	27
2.1 Analýza rozptylu se dvěma faktory.....	28
2.2 Předpoklady analýzy rozptylu se dvěma faktory	30
2.3 Kruskal – Wallisova analýza rozptylu	42
3 REGRESNÍ ANALÝZA – JEDNOROZMĚRNÁ LINEÁRNÍ REGRESE	46
3.1 Regresní analýza	47
3.2 Jednoduchá regresní analýza.....	48
3.3 Metoda nejmenších čtverců.....	48
3.4 Míra variability, koeficient determinace	50
3.5 Klasický lineární model	51
3.6 Diagnostická kontrola modelu	52
3.6.1 Heteroskedasticita.....	52
3.6.2 Autokorelace	53
3.6.3 Normalita	53
4 REGRESNÍ ANALÝZA – JEDNOROZMĚRNÁ: INTERVALY SPOLEHLIVOSTI, TESTY HYPOTÉZ, NELINEÁRNÍ REGRESE	65
4.1 Intervaly spolehlivosti	66
4.2 Testy hypotéz	67
4.3 Nelineární regresní analýza.....	68
4.4 Parabolická regrese	69
4.5 Törnqvistovy funkce	70
4.6 Metoda vybraných bodů.....	72
5 REGRESNÍ ANALÝZA – VÍCEROZMĚRNÁ.....	85

5.1	Vícerozměrná regresní analýza	86
5.2	Metoda nejmenších čtverců.....	86
5.3	Náhodný vektor a jeho charakteristiky.....	88
5.4	Klasický lineární model	88
5.5	Míry variability a koeficient determinace	89
5.6	Intervaly spolehlivosti a testy hypotéz	90
5.7	Individuální <i>T</i> -testy o hodnotách regresních koeficientů.....	91
5.8	<i>F</i> -test hypotézy o hodnotách regresních koeficientů.....	92
6	REGRESNÍ ANALÝZA – VÍCEROZMĚRNÁ: MULTIKOLINEARITA, HETEROSKEDASTICITA, AUTOKORELACE.....	103
6.1	Co je multikolinearita?	104
6.2	Co je heteroskedasticita?.....	107
6.2.1	Jak zjistit heteroskedasticitu?.....	108
6.2.2	Jak odstranit heteroskedasticitu?.....	110
6.3	Co znamená autokorelace?.....	115
7	ZÁKLADY ANALÝZY ČASOVÝCH ŘAD	124
7.1	Typy ekonomických časových řad.....	125
7.2	Elementární charakteristiky časových řad.....	127
7.3	Modely ekonomických časových řad.....	128
8	ANALÝZA TRENDU ČASOVÝCH ŘAD	132
8.1	Trendová složka časových řad	133
8.2	Trendové funkce.....	134
8.2.1	Lineární trend.....	134
8.2.2	Kvadratický trend.....	137
8.2.3	Mocninný trend.....	137
8.2.4	Exponenciální trend	138
8.2.5	Logistický trend	140
8.2.6	Gompertzův trend	142
8.3	Volba vhodného modelu trendu	143
8.4	Klouzavé průměry	144
8.5	Exponenciální vyrovnání.....	145
9	SEZÓNŇÍ SLOŽKA, NÁHODNÁ SLOŽKA.....	150
9.1	Model konstantní sezónnosti se schodovitým trendem.....	151
9.2	Model konstantní sezónnosti s lineárním trendem.....	152

9.3	Model proporcionální sezónnosti	152
9.4	Analýza náhodné složky.....	153
10	MODEL Y TYP U ARIMA A PREDIKCE ČASOVÝCH ŘAD.....	159
10.1	Program GRET L.....	160
10.2	Modelování časových řad pomocí ARIMA modelu.....	161
10.2.1	autoregresivní proces (ar)	162
10.2.2	proces klouzavých průměrů (ma).....	162
10.2.3	autoregresivní proces klouzavých průměrů (arma).....	163
10.2.4	autoregresivní a integrovaný proces klouzavých průměrů (arima).....	163
10.3	Box – Jenkinsova metodologie prognózování časových řad.....	164
LITERATURA		178
SHRNUTÍ STUDI JNÍ OPORY		179
PŘEHLED DOSTUPNÝCH IKON.....		180

ÚVODEM

Tento text představuje studijní oporu pro studium všech akreditovaných studijních programů v navazujícím magisterském studiu na Slezské univerzitě, Obchodně podnikatelské fakultě v Karviné. Předmět Statistické zpracování dat navazuje na předmět Statistika z bakalářského studia. V opoře je kladen důraz především na uplatnění statistických metod při zpracování ekonomických dat v aplikovaných ekonomických disciplínách, jako jsou zejména marketing a management.

Učební text této knihy nabízí studentům vysokých škol ekonomického zaměření strukturovaný a komplexní přehled o 10 důležitých tematických kapitolách. Každá kapitola je přibližně stejně rozsáhlá a obtížností vyvážená, což umožňuje učebním materiálům po-kryt dostatečný rozsah znalostí, aby se studenti mohli seznámit s klíčovými koncepty a metodami v každé oblasti.

Jednotlivé kapitoly jsou navrženy tak, aby odpovídaly délce běžné dvouhodinové prezenční přednášky. To umožňuje studentům přístup k obsáhlému materiálu v relativně stravitelném a dobře organizovaném formátu. Každá kapitola se soustředí na určitou tematickou oblast, a to umožňuje studentům hlouběji proniknout do konkrétních témat a získat ucelené porozumění ekonomickým analýzám.

V případě prezenčního studia na vysoké škole je každá přednáška doplněna seminářem. Semináře jsou klíčovým prvkem výuky, protože umožňují studentům aplikovat teoretické znalosti na praktické číselné příklady. Tímto způsobem studenti získávají dovednosti a praxi v řešení reálných ekonomických situací. Navíc jsou semináře vybaveny počítačovými technologiemi, což umožňuje efektivnější řešení složitějších problémů a analýz. Ve studijní opoře jsou použity programy Excel a GRETL.

Kombinace prezenčních přednášek a seminářů vytváří bohaté a interaktivní učební prostředí, které podporuje aktivní zapojení studentů a podporuje jejich schopnost kriticky myslet, analyzovat a aplikovat naučené koncepty. Díky této kombinaci jsou studenti připraveni na praktickou aplikaci svých znalostí v reálném světě ekonomické praxe.

RYCHLÝ NÁHLED STUDIJNÍ OPORY

Vysokoškolské studium v případě předmětu Statistické zpracování dat vyžaduje enormní úsilí studenta zaměřené na pravidelnost a vytrvalost ve studiu i samostudiu, schopnost koncentrace na předmět, aktivní přístup spočívající v samostatném řešení příkladů. V tom všem by tato studijní opora měla studentům kombinované formy studia pomoci nahradit kvalitní prezenční výuku i úlohu učebnic a skript. Studijní opora je k tomu účelu vybavena určitými nástroji, o jejichž funkcích byste měli být informováni a mohli je tudíž účelně využívat ve svůj prospěch. Pro lepší zvládnutí látky jsou vám v elektronické verzi kurzu Statistické zpracování dat k dispozici ještě doplňkové materiály v elektronické podobě. Dalšími podpůrnými zdroji ke studiu mohou být klasické učebnice a skripta a další doporučená literatura.

Předpokladem pro úspěšné zvládnutí tohoto předmětu Statistické zpracování dat je zvládnutí bakalářského předmětu Statistika na SU OPF nebo odpovídajícího základního bakalářského kurzu Pravděpodobnosti – Statistiky, a to podle typu bakalářského studia na některé VŠ v ČR.

Tato studijní opora se zaměřuje na důležité statistické metody v oblasti ekonomie a jejich aplikaci v různých ekonomických analýzách. Obsahem prvních dvou kapitol je analýza rozptylu, známá také jako ANOVA (Analysis of Variance). Tato metoda je nezbytná pro srovnání více skupin dat a zjišťování, zda existují signifikantní rozdíly mezi nimi. ANOVA poskytuje cenné informace o vztazích mezi proměnnými a je klíčovou technikou v ekonomickém výzkumu.

Následující tři kapitoly (kapitoly 3 až 6) se věnují regresní analýze. Regrese je další klíčovou metodou v ekonomických analýzách, která se zabývá predikcí a modelováním vztahů mezi závislými a nezávislými proměnnými. Tyto kapitoly se soustředí na jak jednoduchou regresní analýzu, která zkoumá vztah mezi jednou nezávislou a jednou závislou proměnnou, tak i na vícerozměrnou regresní analýzu, která zahrnuje více nezávislých proměnných. Regresní analýza má široké uplatnění v ekonomii, například při predikci ekonomických ukazatelů nebo studiu vlivu různých faktorů na ekonomické jevy.

V posledních čtyřech kapitolách (kapitoly 7 až 10) se studijní opora věnuje analýze ekonomických časových řad. Tato oblast je v ekonomii mimořádně významná, protože se zabývá analýzou a predikcí ekonomických dat, která jsou získávána v pravidelných časových intervalech. Ekonomické časové řady mohou poskytnout cenné informace o dlouhodobých trendech, sezónních vlivu, cyklech a jiných periodických vzorcích v ekonomických datech. Analýza ekonomických časových řad je klíčovým nástrojem pro ekonomické prognózování, plánování a strategické rozhodování.

Studijní opora poskytuje komplexní přehled statistických metod používaných v ekonomických analýzách.

1 ANALÝZA ROZPTYLU (ANOVA) – JEDEN FAKTOR

RYCHLÝ NÁHLED KAPITOLY



Jednofaktorová metoda ANOVA, kterou prokazujeme závislost hodnot znaků Y na faktoru X , pro něž jsou k dispozici příslušná data, spočívá v tom, že celkovou variabilitu měřenou součtem čtverců odchylek od celkového průměru rozdělíme na variabilitu uvnitř jednotlivých výběrů a na variabilitu mezi jednotlivými výběry. Cílem, k němuž směřujeme, je buď přijmout nulovou hypotézu o vzájemné nezávislosti Y na X , nebo ji zamítnout (na zvolené hladině významnosti). Jedná se tedy o běžný statistický postup nazývaný testování statistických hypotéz, známý ze základního kurzu statistiky. V případě přijetí nulové hypotézy vyvozujeme nezávislost hodnot Y na X , v opačném případě konstatujeme, že Y na X závisí.

V této kapitole se naučíte, jak tento test statistické hypotézy konkrétně provést: jak vypočítat hodnotu testového kritéria a příslušnou kritickou hodnotu a jak vyvodit z těchto hodnot příslušný závěr týkající se eventuální závislosti nebo nezávislosti hodnot znaku Y na faktoru X .

CÍLE KAPITOLY



Po prostudování této kapitoly budete umět:

- vypočítat hodnotu testového kritéria,
- najít příslušnou kritickou hodnotu z tabulek Fisherova rozdělení,
- zkonstruovat tabulku ANOVA,
- přijmout nebo zamítnout nulovou hypotézu o nezávislosti hodnot znaku Y na faktoru X .

ČAS POTŘEBNÝ KE STUDIU



K prostudování této kapitoly budete potřebovat asi 60 minut.



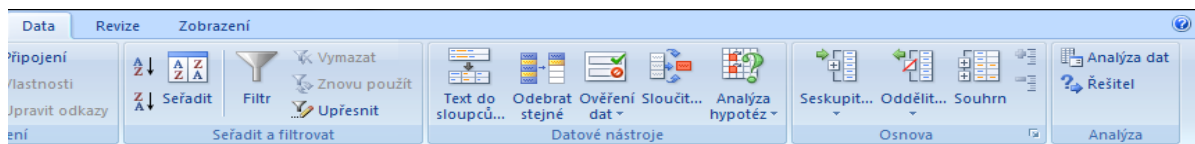
KLÍČOVÁ SLOVA KAPITOLY

Analýza rozptylu, testové kritérium, kritická hodnota, ANOVA tabulka.

Analýza rozptylu umožňuje ověřit významnost rozdílu mezi výběrovými průměry většího počtu náhodných výběrů, umožňuje posoudit vliv různých faktorů na hospodářský proces charakterizovaný kvantitativním statistickým znakem. Základní myšlenka analýzy rozptylu spočívá v rozkladu celkového rozptylu na dílčí rozptyly příslušející jednotlivým vlivům, podle nichž jsou data roztržena. Kromě dílčích rozptylů je jednou složkou celkového rozptylu tzv. reziduální rozptyl, způsobený nepostiženými vlivy. Podle počtu analyzovaných faktorů rozlišujeme *jednofaktorovou*, *dvoufaktorovou* a *vícefaktorovou* analýzu rozptylu. Všeobecně používané označení ANOVA je akronymem anglických slov „ANalysis Of VAriance“ (doslovný překlad: analýza rozptylu).

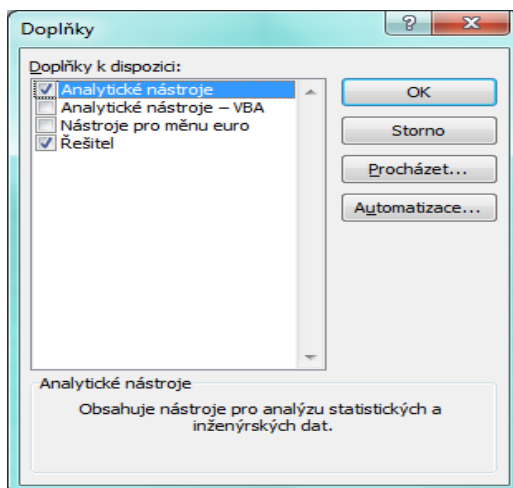
Klasická ANOVA vychází, jak uvidíte, z předpokladu normality rozdělení hodnot daného faktoru. Pokud je takový předpoklad neudržitelný, lze použít analýzu rozptylu jiného typu, konkrétně Kruskal-Wallisovu verzi ANOVA. Jednofaktorovou ANOVA se zabývá tato kapitola, vícefaktorová a Kruskal-Wallisova ANOVA je obsahem kapitoly následující.

V tomto studijním textu předpokládáme, že čtenář má k dispozici verzi Excel 2010, eventuálně vyšší. Pro zjednodušení práce je vhodné mít aktivovaný doplněk „Analýza dat“ a „Řešitel“ ve složce „Data“ (viz Obrázek 1).



Obrázek 1: Doplněk Analýza dat

V případě, že tyto doplňky nejsou ve složce „Data“, lehce je nainstalujete tímto postupem: „Tlačítko Soubor“ → „Možnosti“ → „Doplňky“ → „Přejít...“ a v dialogovém okně zaškrtnout položky „Analytické nástroje“ a „Řešitel“ (viz Obrázek 2).



Obrázek 2: Doplňky

1.1 Nezávislý a závislý faktor

Často se vyskytuje situace, kdy máme k nezávislých náhodných výběrů, které obecně nemusí pocházet z jednoho základního souboru, nebo jinak řečeno, nemusí být stejného typu, s rozsahy, tj. počty prvků n_1, n_2, \dots, n_k . Číslo k může být libovolné podle konkrétní situace, např. 2, 3, 4, ... Tyto rozsahy výběrů rovněž nemusí být stejné, v každém z nich budiž znám průměr \bar{x}_i , a také rozptyl s_i^2 , $i = 1, 2, \dots, k$. V praktických situacích obvykle tyto výběry vzniknou tak, že základní soubor rozdělíme podle určitého statistického znaku X do k skupin, např. věkových, v každé z nich pak máme n_i prvků, $i = 1, 2, \dots, k$. Znak X pak označujeme jako *nezávislý faktor*, jehož hodnoty předem stanovíme, stanovíme např. věkové skupiny takto: do 18 let, 19 až 29 let, 30 až 59 let, 60 a více let, v tomto příkladu je $k = 4$. Hovoříme proto často o *faktoru kontrolovaném*. Další příklady faktorů: velikost rodiny, měsíční příjem rodiny, velikost podniku, typ ekonomické činnosti apod. Hodnotami faktoru X jsou obvykle *kvalitativní* (nečíselné) *veličiny*, označujeme je symbolicky x_1, x_2, \dots, x_k . Tyto hodnoty mohou, ale nemusí být nutně vzájemně uspořádány.

Faktor X , jež nabývá k kvalitativních hodnot, může, ale nemusí ovlivňovat hodnoty statistického znaku Y , o kterém předpokládáme, že má na rozdíl od X *kvantitativní* (tedy číselnou) povahu.

Cílem ANOVA je právě prokázat, že hodnoty kvalitativního znaku X ovlivňují hodnoty kvantitativního znaku Y (závislého faktoru). Hodnoty znaku Y , které přísluší hodnotě x_i faktoru X , označujeme $y_{i1}, y_{i2}, \dots, y_{in_i}$. Pro analýzu rozptylu je výhodné uspořádat výchozí údaje do přehledné tabulky, viz Tabulka 1.

Tabulka 1: Schéma výchozí tabulky analýzy rozptylu pro jeden faktor

Číslo výběru	Zjištěné hodnoty sledovaného znaku	Počet prvků	Průměr	Rozptyl
1	$y_{11}, y_{12}, \dots, y_{1j}, \dots, y_{1n_1}$	n_1	\bar{y}_1	s_1^2
2	$y_{21}, y_{22}, \dots, y_{2j}, \dots, y_{2n_2}$	n_1	\bar{y}_2	s_2^2
⋮	⋮	⋮	⋮	⋮
i	$y_{i1}, y_{i2}, \dots, y_{ij}, \dots, y_{in_i}$	n_i	\bar{y}_i	s_i^2
⋮	⋮	⋮	⋮	⋮
k	$y_{k1}, y_{k2}, \dots, y_{kj}, \dots, y_{kn_k}$	n_k	\bar{y}_k	s_k^2
Celkem		n	\bar{y}	s^2

Princip metody ANOVA, kterou prokazujeme závislost Y na X , spočívá v tom, že celkovou variabilitu měřenou součtem čtverců odchylek od celkového průměru rozdělíme na variabilitu uvnitř jednotlivých výběrů a na variabilitu mezi jednotlivými výběry.

1.2 Předpoklady analýzy rozptylu s jedním faktorem

Předpokládáme, že faktor X má k úrovní (hodnot x_i), s účinkem na znak Y , který lze vyjádřit vztahem: $\mu_i = \mu + \alpha_i, i = 1, 2, \dots, k$,

kde μ_i je průměr znaku Y v i -té skupině (příslušné k hodnotě faktoru x_i),

μ je celkový průměr znaku Y ,

α_i je efekt hodnoty faktoru x_i na znak Y .

Formulujeme nyní nulovou hypotézu H_0 , že všechny výběry pocházejí ze stejné základní populace (základního souboru), jinak řečeno, že hodnoty faktoru X nemají na hodnoty znaku Y žádný efekt (vliv).

Budeme dále předpokládat, že hodnoty α_i pocházejí z normálně rozdělené populace s nulovou střední hodnotou a konstantním rozptylem σ^2 .

Formulujeme nulovou hypotézu: $H_0: E(\alpha_1) = E(\alpha_2) = \dots = E(\alpha_k) = 0$, proti alternativní hypotéze, že H_0 neplatí, že alespoň pro dvě položky, např. i a j , platí: $H_1: E(\alpha_i) \neq E(\alpha_j)$.

Symbolem $E(\alpha_i)$ označujeme střední hodnotu náhodné veličiny α_i . Předpoklad konstantního rozptylu pro všechny veličiny α_i je podstatný, je ho možno ověřit statistickým testem, a to buď tzv. Bartlettovým testem, s nímž se seznámíte později. Normalitu rozdělení veličin α_i lze taktéž ověřit příslušným testem, např. Chi-kvadrát testem dobré shody, známým ze základního kurzu statistiky, viz Ramík (2003). V praxi obvykle předpokládáme (na podkladě věcné znalosti problému), že zmíněné dva předpoklady jsou automaticky splněny a při aplikaci ANOVA je již obvykle neověřujeme.

Cílem, k němuž směřujeme, je buď *přijmout* nulovou hypotézu H_0 , nebo H_0 *zamítnout* (na zvolené *hladině významnosti*). Jedná se tedy o běžný statistický postup nazývaný *testování statistických hypotéz*, známý ze základního kurzu statistiky, viz Ramík (2003). V případě přijetí nulové hypotézy vyvozujeme nezávislost hodnot faktoru Y na faktoru X , jinak řečeno: faktor Y na faktoru X nezávisí. V opačném případě (při zamítnutí H_0), konstatujeme, že faktor Y na faktoru X závisí neboli faktor X ovlivňuje Y .

1.3 Předpoklady analýzy rozptylu s jedním faktorem

Celkovou variabilitu znaku Y změříme *výběrovým rozptylem*

$$s^2 = \frac{\sum_i \sum_j (y_{ij} - \bar{y})^2}{n-1}. \quad (1.1)$$

V souvislosti s analýzou rozptylu se budeme zabývat pouze čitatelem výše uvedeného zlomku, totiž součtem čtverců odchylek zjištěných hodnot y_{ij} od celkového průměru \bar{y} , přičemž průměr vypočítáme podle známého vztahu: sečteme všechny hodnoty a výsledek podělíme jejich počtem, tedy $\bar{y} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}$.

Tento *celkový součet čtverců* budeme označovat symbolem S_y , tj.

$$S_y = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2. \quad (1.2)$$

Celkovému součtu čtverců přísluší počet stupňů volnosti $df_y = n-1$.

Variabilitu mezi skupinami budeme měřit *meziskupinovým součtem čtverců* $S_{y,m}$, který definujeme následovně

$$S_{y,m} = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2. \quad (1.3)$$

Meziskupinovému součtu čtverců přísluší počet stupňů volnosti $df_m = k-1$.

Variabilitu uvnitř skupin označujeme jako *vnitroskupinovou*, nebo také *reziduální* a používáme přitom označení $S_{y,v}$, přičemž definujeme *vnitroskupinový (reziduální) součet čtverců* takto

$$S_{y,v} = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2. \quad (1.4)$$

Vnitroskupinovému součtu čtverců přísluší počet stupňů volnosti $df_v = n-k$.

Aritmetickými úpravami výše uvedených vzorců lze snadno dokázat základní vztah analýzy rozptylu, totiž, že *celkový součet čtverců je roven sumě meziskupinového a vnitroskupinového součtu čtverců*, symbolicky:

$$S_y = S_{y,m} + S_{y,v}. \quad (1.5)$$

Pro ověření nulové hypotézy H_0 použijeme statistiku:

$$F = \frac{\frac{S_{y,m}}{k-1}}{\frac{S_{y,v}}{n-k}} = \frac{\frac{S_{y,m}}{df_m}}{\frac{S_{y,v}}{df_v}} \quad (1.6)$$

kteřá má při platnosti nulové hypotézy *Fisherovo rozdělení* $F(k-1, n-k)$. Kritické hodnoty Fisherova rozdělení $F_\alpha(df_1, df_2)$ jsou tabelovány pro různé hodnoty hladiny významnosti α a různé hodnoty parametrů (stupňů volnosti: *degree of freedom*) df_1 a df_2 . Někdy se namísto kritických hodnot tabelují kvantily Fisherova rozdělení $F_{1-\alpha}^k(df_1, df_2)$. Vztah mezi kritickými hodnotami a kvantily je jednoduchý:

$$F_\alpha(df_1, df_2) = F_{1-\alpha}^k(df_1, df_2).$$

Např. 5 % kritická hodnota je rovna 95 % kvantilu při stejných hodnotách parametrů df_1 a df_2 .

Pro výpočet kritických hodnot lze využít Excelu. Postupuje se přitom takto: v hlavním menu postupně vybíráte: Vložit → Funkce → Statistické → FINV ($\alpha, df_1; df_2$).

Postup testování hypotézy H_0 charakterizujeme následujícími 3 kroky:

Krok 1. Zvolte hladinu významnosti α , která představuje chybu 1. druhu, tj. pravděpodobnost zamítnuti správné hypotézy. Praktické hodnoty hladiny významnosti α jsou: 0,1; 0,05; 0,01 neboli v procentech: 10%, 5%, 1%.

Krok 2. Vypočtete hodnotu statistiky F podle vzorce (1.6), přičemž pro hodnoty meziskupinového součtu čtverců $S_{y,m}$ a pro výpočet vnitroskupinového součtu čtverců $S_{y,v}$ použijte vzorce (1.3) a (1.4). Výpočetně výhodnější, např. pro výpočet na kalkulačce, jsou následující vzorce:

$$S_y = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}^2 - \frac{1}{n} \left(\sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij} \right)^2, \quad (1.7)$$

$$S_{y,m} = \sum_{i=1}^k n_i \bar{y}_i^2 - \frac{1}{n} \left(\sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij} \right)^2. \quad (1.8)$$

K výpočtu $S_{y,v}$ lze využít základního vztahu (1.5) a právě uvedených vztahů (1.7) a (1.8): $S_{y,v} = S_y - S_{y,m}$.

Krok 3. Porovnejte hodnotu statistiky F vypočtené v Kroku 2 s kritickou hodnotou $F_{\alpha}(k - 1, n - k)$. Výsledek tohoto porovnání může být dvojit:

I. Platí $F \leq F_{\alpha}(k - 1, n - k)$.

Potom se nulová hypotéza H_0 *přijímá* (nezamítá) a tudíž se konstatuje, že hodnoty faktoru X *nemají* na hodnoty znaku Y *statisticky významný vliv* (na zvolené hladině významnosti). Jinak řečeno, faktor X je *neúčinný*.

II. Platí $F > F_{\alpha}(k - 1, n - k)$.

Potom se nulová hypotéza H_0 *zamítá*, přijímá se hypotézu alternativní H_1 , a tudíž se konstatuje, že hodnoty faktoru X *mají* na hodnoty znaku Y *statisticky významný vliv*. Jinak řečeno, faktor X je *účinný*.

Podaří-li se výše uvedeným testem prokázat, že hodnoty faktoru X mají na hodnoty znaku Y *statisticky významný vliv*, mohou nás zajímat další informace o tom, které skupiny se významně odlišují od průměru, eventuálně jak skupinové průměry seřadit, případně zařadit do společných celků. V krajním případě by se totiž mohlo stát, že významnost rozdílnosti k skupin způsobuje jediná skupina a ostatní skupiny se navzájem neliší. Touto problematikou se zabývají metody tzv. *simultánního testování*, z nichž nejznámější je metoda Shaffého. Vy se touto problematikou zde nezabývat nebudete, zájemce odkazujeme na literaturu, viz např. Anděl (2007).

Metoda analýzy rozptylu je založena na předpokladech shody rozptylů v jednotlivých k skupinách. Pokud jsou předpoklady splněny, pak popsání metoda ANOVA poskytuje nejlepší výsledky – je nejúčinnější. Není-li tento předpoklad splněn, pak použití výše uvedeného testu může poskytnout nesprávný výsledek. V takovém případě lze použít jiné metody, např. Kruskal-Wallisova ANOVA, která používá Chi-kvadrát test, s níž se seznámíte v příští kapitole.

V Excelu jsou k dispozici funkce, které umožňují řešit jednofaktorové i vícefaktorové úlohy ANOVA. Naleznete je v hlavním menu: *Nástroje* → *Analýza dat* → *ANOVA: jeden faktor*.

1.4 Míra těsnosti závislosti

Variabilita podmíněných (skupinových) průměrů \bar{y}_i kolem celkového průměru \bar{y} je způsobena závislostí znaku Y na znaku X . Tuto variabilitu jsme vyjádřili meziskupinovým součtem čtverců $S_{y,m}$. Variabilita znaku Y uvnitř jednotlivých skupin – vyjádřena vnitroskupinovým (reziduálním) součtem čtverců $S_{y,v}$, je způsobena jinými (neuvažovanými) činiteli. Čím větší je $S_{y,m}$, tím větší je těsnost závislosti znaků X a Y . Protože však jsou jednotlivé součty čtverců vzájemně vázány vztahem (1.5), lze míru těsnosti závislosti vyjádřit

jako podíl meziskupinového a celkového součtu čtverců. Zavádíme proto jako míru těsnosti závislosti znaku Y na znaku X poměr determinace P^2 takto:

$$P^2 = \frac{S_{y,m}}{S_y}. \quad (1.9)$$

Odmocninu z poměru determinace P nazýváme *poměr korelace*.

Poměr determinace nabývá hodnot z intervalu $[0,1]$. Čím těsnější je závislost Y na X , tím více se hodnota poměru determinace blíží k 1, tím více se také vnitroskupinový součet čtverců blíží k celkovému součtu čtverců, přičemž meziskupinový součet čtverců se blíží k nule. Naopak, čím více se poměr determinace blíží k 0, tím menší část z celkového součtu čtverců tvoří meziskupinový součet čtverců (na úkor vnitroskupinového), a tím menší je těsnost závislosti znaku Y na X . Způsob výpočtu determinačního a korelačního poměru si procvičíte na numerických příkladech. V Excelu bohužel funkce pro výpočet poměru determinace nebo korelace chybí, musí se proto k výpočtu použít vzorce (1.9).

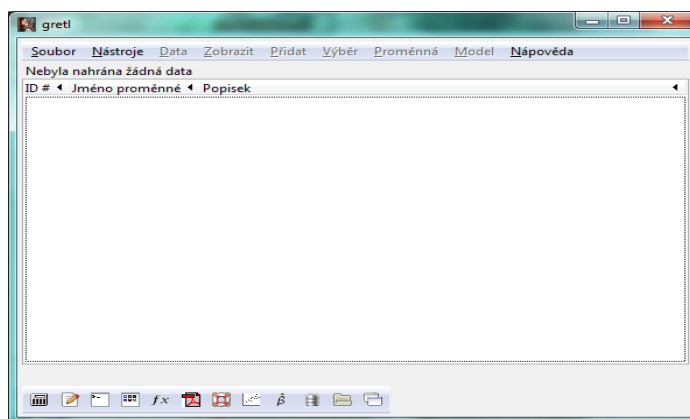
Uvědomte si však, že poměr determinace P^2 je náhodná veličina (jakožto podíl dvou veličin – součtu čtverců, které jsou samy náhodnými veličinami), proto může být výsledkem kladné číslo i v případě, že výsledkem ANOVA je fakt, že zkoumaný faktor není statisticky významný neboli sledovaná veličina na faktoru nezávisí. V takovém případě by logicky mělo platit, že poměr determinace P^2 je nulový, tj. $P^2 = 0$. Tento zdánlivý rozpor vysvětlujeme statistickým přístupem: testem statistické hypotézy. Nulová hypotéza $H_0: P^2 = 0$. Jako testové kritérium se použije statistika F ze vzorce (1.6).

Pokud platí $F \leq F_\alpha(k-1, n-k)$, potom nulovou hypotézu H_0 nelze zamítnout a hodnoty faktoru X nemají na hodnoty znaku Y statisticky významný vliv na zvolené hladině významnosti a poměr determinace (samozřejmě i poměr korelace) je roven nule, jinak řečeno, je statisticky nevýznamný.

V opačném případě se nulová hypotéza zamítá a poměr determinace je statisticky významný. Hodnota poměru determinace i poměru korelace je nenulová. V tom případě má smysl hovořit o síle závislosti veličiny Y na faktoru X .

1.5 Analýza rozptylu v programu GRETL

GRETL je volně dostupný produkt se zaměřením na statistické metody, které podporují ekonometrické analýzy. Název je akronymem pro *GNU Regression, Econometric and Time-series Library*. Systém GRETL se dá používat dvěma způsoby. Snaha tvůrců systému od začátku směřovala k přiblížení ekonometrie široké veřejnosti a bylo vytvořeno grafické uživatelské rozhraní (GUI – Graphical User Interface), které je pro většinu běžných uživatelů přijatelnější. Po spuštění programu se objeví hlavní okno (Obrázek 3). V horní části je *hlavní menu* a v dolní části se nachází *panel nástrojů*.



Obrázek 3: Hlavní okno programu GRETL

Po instalování program obsahuje velký počet datových souborů, které se dají otevřít z hlavního menu – Soubor – Otevřít data – Vzorový soubor. Je zde možno vybírat z databáze Ramanathan, Greene, Stock and Watson. Záložka Data poskytuje velký prostor na přizpůsobení databáze podmínkám modelování.

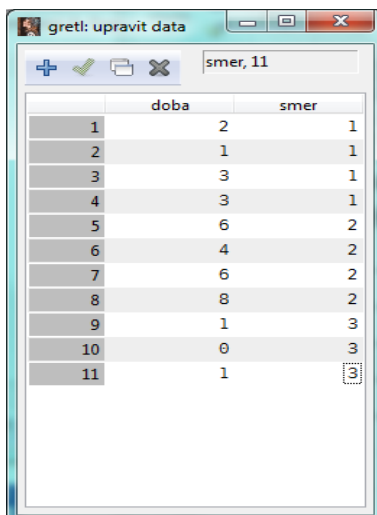
Na následujícím příkladu si ukážeme, jak se zadávají data do programu GRETL. Tabulka 2 uvádí, kolik dnů po přiletu trvá adaptace na časový posun (JETLAG). Na hladině významnosti 5 % ověříme, má-li směr letu vliv na délku adaptace (zotavení).

Tabulka 2: Doba adaptace ve dnech

Směr	Doba adaptace ve dnech			
Západ	2	1	3	3
Východ	6	4	6	8
Stejný	1	0	1	

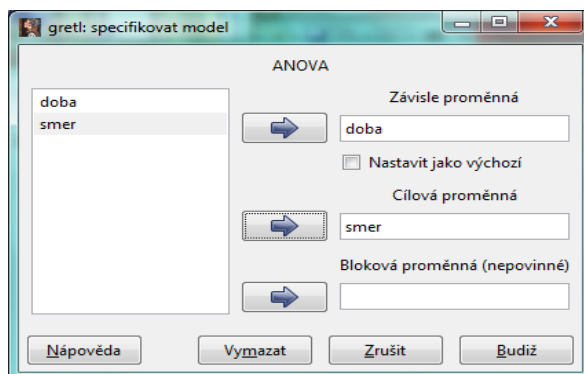
Nulová hypotéza tvrdí, že doba adaptace nezávisí na časovém posunu. Alternativní hypotéza tvrdí, že doba adaptace závisí na časovém posunu.

V hlavním menu vybereme nový soubor dat – počet pozorování=11. Struktura souboru dat = průřezová. Kvantitativní proměnnou jsme pojmenovali doba a jednotlivé varianty kvalitativního znaku (směr) musí být přiřozena čísla (1-západ, 2-východ, 3-stejný).



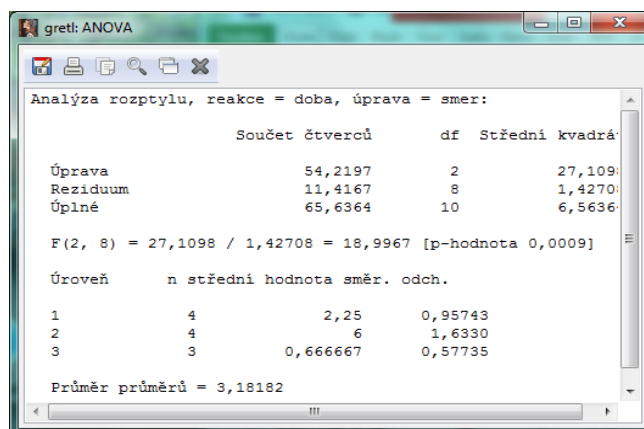
Obrázek 4: Zadávání hodnot do GRETLU

Dále vybereme posloupnost příkazů: *Model – Další lineární modely – ANOVA*.



Obrázek 5: Analýza rozptylu v GRETLU – zadání

Zadání hodnot potvrdíme tlačítkem *Budíž* a dostáváme následující výsledek.



Obrázek 6: Analýza rozptylu v GRETLU - výstup

Protože p -hodnota = 0,0009 je menší než hladina významnosti (0,05), nulovou hypotézu zamítáme. Můžeme tedy z 95 % tvrdit, že doba adaptace závisí na časovém posunu. Tato skutečnost byla dokonce prokázána i na hladině významnosti 0,01.

ŘEŠENÁ ÚLOHA 1.1



Na testovacím okruhu byla testována průměrná spotřeba tří automobilů téže třídy různých výrobců Škoda, Renault a Fiat. Řidič absolvoval s každým automobilem 5 testovacích jízd. Tabulka ukazuje spotřebu benzínu na 100 kilometrů v jednotlivých jízdách.

Automobil	Spotřeba				
Škoda	7,4	7,8	6,8	7,6	8,1
Renault	6,7	7,2	8,3	7,1	7,5
Fiat	6,8	6,9	7,3	7,9	7,6

Na hladině významnosti $\alpha = 0,05$ zjistěte, zda má typ automobilu vliv na spotřebu benzínu. V kladném případě vypočítejte determinační a korelační poměr.

Řešení:

Chceme zjistit závislost znaku Y (průměrná spotřeba) na jediném znaku X (výrobce automobilu). Provedeme proto jednofaktorovou analýzu rozptylu.

Faktor X má tři hodnoty: $x_1 = \text{Škoda}$, $x_2 = \text{Renault}$, $x_3 = \text{Fiat}$, tzn. $k = 3$, s počty hodnot $n_1 = n_2 = n_3 = 5$ v každé z nich budeme testovat nulovou hypotézu

$H_0: E(\alpha_1) = E(\alpha_2) = E(\alpha_3) = 0$, tj. průměrná spotřeba je u všech vozidel stejná.

Alternativní hypotéza H_1 je negací nulové hypotézy.

Nejprve vypočítáme podmíněné průměry $\bar{y}_1, \bar{y}_2, \bar{y}_3$

$$\bar{y}_1 = \frac{\sum_{j=1}^5 y_{1j}}{5} = \frac{7,4 + 7,8 + \dots + 8,1}{5} = 7,54$$

$$\bar{y}_2 = \frac{\sum_{j=1}^5 y_{2j}}{5} = \frac{6,7 + 7,2 + \dots + 7,5}{5} = 7,36$$

$$\bar{y}_3 = \frac{\sum_{j=1}^5 y_{3j}}{5} = \frac{6,8 + 6,9 + \dots + 7,6}{5} = 7,3$$

a celkový průměr znaku Y

$$\bar{y} = \frac{\sum y_{ij}}{n} = \frac{7,4 + 7,8 + \dots + 7,6}{15} = 7,4.$$

Dále vypočítáme pomocí vztahů (1.2), (1.3), popř. (1.7), (1.8) součty S_y a S_{ym} .

$$\begin{aligned} S_y &= \sum_{i=1}^3 \sum_{j=1}^5 (y_{ij} - \bar{y})^2 = (7,4 - 7,4)^2 + (7,8 - 7,4)^2 + \dots + (8,1 - 7,4)^2 + \\ &\quad + (6,7 - 7,4)^2 + (7,2 - 7,4)^2 + \dots + (7,5 - 7,4)^2 + \\ &\quad + (6,8 - 7,4)^2 + \dots + (7,6 - 7,4)^2 = 3,4 \end{aligned}$$

$$S_{ym} = \sum_{i=1}^3 n_i (y_{ij} - \bar{y})^2 = 5(\bar{y}_1 - \bar{y})^2 + 5(\bar{y}_2 - \bar{y})^2 + 5(\bar{y}_3 - \bar{y})^2 =$$

$$= 5(7,54 - 7,4)^2 + 5(7,36 - 7,4)^2 + 5(7,3 - 7,4)^2 = 0,16.$$

Součet S_{ym} má $k-1$ stupňů volnosti, v našem případě $df_m = 3-1 = 2$.

Pomocí součtů S_y a S_{ym} dopočítáme součet S_{yv} , neboť $S_y = S_{yv} + S_{ym}$.
Proto $S_{yv} = S_y - S_{ym} = 3,4 - 0,16 = 3,24$.

Součet S_{yv} má $n-k$ stupňů volnosti, proto $df_v = 15-3 = 12$.

Testové kritérium F vypočítáme podle vztahu (1.6):

$$F = \frac{\frac{S_{ym}}{k-1}}{\frac{S_{yv}}{n-k}} = \frac{\frac{0,16}{2}}{\frac{3,24}{12}} = 0,296.$$

Pro stanovení kritického oboru C najdeme v tabulkách kritických hodnot $F_{\alpha}(k-1, n-k)$ kritickou hodnotu $F_{0,05}(2, 12) = 3,89$ (ověřte v Excelu pomocí funkce FINV). Kritický obor je proto interval od 3,89 do nekonečna, tj. $C = (3,89, +\infty)$. Zřejmě platí $0,296 < 3,89$, tzn. $F \notin C$, proto nulovou hypotézu H_0 přijímáme.

Znamená to, že faktor X -výrobce automobilu je neúčinný neboli, že průměrná spotřeba benzínu není statisticky významně ovlivněna výrobcem automobilu. Poměr determinace i korelace je tedy 0.



ŘEŠENÁ ÚLOHA 1.2

Rozhodněte, zda velikost výnosů petržele (faktor Y) závisí na použitém druhu hnojiva (faktor X). Pokud závisí, pak pomocí determinačního poměru zjistíte těsnost této závislosti. Data jsou uvedena v následující tabulce, použijte hladinu významnosti 0,05.

Hnojivo	Výnosy (1kg/10 m ²)					
A	40	42	45	40	44	47
B	76	75	82	68		
C	60	58	62	64	70	

Řešení:

U tohoto příkladu si ukážeme řešení s pomocí Excelu. Nejprve však příklad vyřešíme klasickým postupem.

K výpočtu hodnot součtů čtverců S_{ym} a S_y , potřebujeme znát celkový průměr \bar{y} a podmíněné průměry $\bar{y}_1, \bar{y}_2, \bar{y}_3$.

$$\bar{y}_1 = \frac{\sum_{j=1}^6 y_{1j}}{n_1} = \frac{40 + 42 + \dots + 47}{6} = 43,$$

$$\bar{y}_2 = 75,25; \bar{y}_3 = 62,8,$$

$$\bar{y} = \frac{\sum_{i=1}^3 n_i \bar{y}_i}{n} = \frac{43 \cdot 6 + 75,25 \cdot 4 + 62,8 \cdot 5}{15} = 58,2.$$

Nyní již můžeme vypočítat součty S_{ym} a S_y , podle vztahů (1.2), (1.3)

$$S_y = \sum_{ij} (y_{ij} - \bar{y})^2 = (40 - 58,2)^2 + \dots + (47 - 58,2)^2 + \\ + (76 - 58,2)^2 + \dots + (68 - 58,2)^2 + \\ + (60 - 58,2)^2 + \dots + (70 - 58,2)^2 = 2878,4.$$

$$S_{ym} = \sum_{i=1}^3 n_i (\bar{y}_i - \bar{y})^2 = 6(43 - 58,2)^2 + 4(75,25 - 58,2)^2 + 5(62,8 - 58,2)^2 = 2654,85.$$

$$\text{Hodnota testového kritéria je } F = \frac{\frac{S_{ym}}{k-1}}{\frac{S_{yv}}{n-k}} = \frac{\frac{2654,85}{2}}{\frac{2878,4 - 2654,85}{12}} = 71,26.$$

Kritická hodnota je $F_{0,05}(2, 12) = 3,89$ a je mnohem menší než hodnota testového kritéria F . Proto nulovou hypotézu zamítáme a konstatujeme, faktor hnojiva významně ovlivňuje hodnoty výnosů petržele.

Hodnotu determinačního poměru P^2 zjistíme dosazením hodnot S_{ym} a S_y do vztahu (1.9).

$$P^2 = \frac{2654,85}{2878,4} = 0,92.$$

Hodnoty determinačního poměru blízké 1 svědčí o vysoké závislosti faktoru Y na faktoru X . Hodnota 0,92 proto znamená, že závislost výnosů petržele na použitém druhu hnojiva je vysoká.

Řešení pomocí Excelu:

Nejprve je zapotřebí připravit v Excelu data. Jednotlivé hodnoty y_{ij} pro faktoru Y pro hodnotu x_i faktoru X uspořádáme do řádků, podobně jako v tabulce v zadání. V prvním sloupci umístíme kvůli lepší orientaci název hodnoty faktoru (popisky) x_i , v tomto případě název hnojiva: A, B, C. Data ve worksheetu Excelu vypadají tedy například takto:

	A	B	C	D	E	F	G	H
1	A	40	42	45	40	44	47	
2	B	76	75	82	68			
3	C	60	58	62	64	70		
4								

Data je možné uspořádat také do sloupců, přitom do prvního řádku umístíme názvy hodnot faktoru X (popisky). To je výhodné zejména u velkého množství dat, tj. pro velkou hodnotu počtu dat n .

Dále otevřeme v hlavním menu postupně položky:

Data → Analýza dat... → ANOVA: jeden faktor

Analýza rozptylu (ANOVA) – JEDEN FAKTOR

Pokud se tam položka Analýza dat nevyskytuje je ji zapotřebí doinstalovat (viz začátek této kapitoly).

Zvolíte-li pak první položku ANOVA: jeden faktor, otevře se zadávací okno, kde postupně zadáte:

Vstupní oblast: \$A\$1:\$G\$3

Sdružit: zakliknete tlačítko *Řádky* (je možné uspořádat data do sloupců, pak ovšem zakliknete tlačítko *Sloupce*

Popisky v prvním sloupci – zakliknete

Alfa: 0,05 (hladina významnosti je předvolena, lze ji však změnit)

Výstupní oblast: \$A\$5 (levý horní roh výstupní oblasti). Potvrdíte OK

Výběr	Počet	Součet	Průměr	Rozptyl
A	6	258	43	8
B	4	301	75,25	32,91667
C	5	314	62,8	21,2

roj variabil	SS	Rozdíl	MS	F	Hodnota P	F krit
Mezi výbě	2654,85	2	1327,425	71,2552	2,19E-07	3,885294
Všechny v	223,55	12	18,62917			
Celkem	2878,4	14				

V první tabulce s názvem Faktor jsou uvedeny základní statistické údaje o datech: Počet, Součet, Průměr a Rozptyl.

Ve druhé tabulce nazvané ANOVA jsou uvedeny výpočty metodou ANOVA, jednotlivé položky mají následující význam:

Mezi výběry = meziskupinový

Všechny výběry = vnitroskupinový

Celkem = celkový

SS = Součet čtverců (Sum of Squares)

Rozdíl = stupeň volnosti (DF – Degree of Freedom)

MS = Průměr čtverců (Mean Square)

F = testové kritérium = 71,25

Hodnota P = Signifikance (p-hodnota) = 0,000000219 < 0,05 = α

F krit = kritická hodnota rozdělení F = 3,89

Hodnoty získané řešením v Excelu jsou stejné jako při použití „ručního“ výpočtu, proto i závěry jsou stejné. V Excelu máme navíc vypočtenou p -hodnotu testu (tzv. signifikanci), která, pokud je menší než zvolená hladina významnosti α , znamená, že nulovou hypotézu zamítáme. V opačném případě nulovou hypotézu nezamítáme (přijímáme).

ŘEŠENÁ ÚLOHA 1.3

Firma Dekorace domu má své prodejny ve čtyřech městech (Ostrava, Karviná, Bohumín, Český Těšín). Tabulka zobrazuje tržby firmy v posledních pěti měsících. Testujte na hladině významnosti 5 %, zda výše tržeb závisí na lokalitě, ve které se prodejna nachází. Pokud bude prokázána závislost, pak pomocí determinačního poměru zjistíte sílu této závislosti. Jak se změní výsledek v případě testování na hladině významnosti 0,01?

Město	Tržby (v tis.Kč)				
Ostrava	71	83	65	77	84
Karviná	60	51	54	80	55
Bohumín	55	55	62	65	63
Český Těšín	68	73	67	59	53

Řešení:

Tento příklad vyřešíme s pomocí Excelu. Nejprve si napíšeme hypotézu, kterou budeme testovat:

H_0 : výše tržeb nezávisí na lokalitě, ve které se prodejna nachází,

H_1 : výše tržeb závisí na lokalitě, ve které se prodejna nachází.

Zadáme posloupnost příkazů: Data → Analýza dat... → ANOVA: jeden faktor

A dostaneme následující výstup, ve kterém můžeme vidět hodnoty podmíněných průměrů, hodnotu meziskupinového součtu = 860, hodnotu vnitroskupinového součtu = 1142, hodnotu testového kritéria $F = 4,016$, kritickou hodnotu Fisherova rozdělení = 3,23, a konečně hodnotu $P = 0,026$.

Anova: jeden faktor						
Faktor						
Výběr	Počet	Součet	Průměr	Rozptyl		
Ostrava	5	380	76	65		
Karviná	5	300	60	135,5		
Bohumín	5	300	60	22		
Český Těšín	5	320	64	63		
ANOVA						
Zdroj variability	SS	Rozdíl	MS	F	Hodnota P	F krit
Mezi výběry	860	3	286,6667	4,016346	0,026236	3,238872
Všechny výběry	1142	16	71,375			
Celkem	2002	19				

Hodnota P představuje minimální hodnotu, od které lze nulovou hypotézu zamítnout. Proto v případě, že testujeme na hladině významnosti 0,05; tak nulovou hypotézu zamítáme, protože $0,026 < 0,05$. Tzn., že z 95 % můžeme tvrdit, že výše tržeb závisí na lokalitě, ve které se prodejna nachází. Kdežto v případě, že testujeme hypotézu na hladině významnosti 0,01; tak nulovou hypotézu nelze zamítnout, protože $0,026 > 0,01$; takže z 99 % nebyla závislost mezi výší tržeb a lokalitou prokázána.

Sílu závislosti posoudíme pomocí poměru determinace. Jde o poměr meziskupinové variability na celkové variabilitě. Výsledek je možné vyjádřit v procentech.

$$p^2 = \frac{S_{ym}}{S_y} = \frac{860}{2002} = 42,96 \%$$



SAMOSTATNÉ ÚKOLY

1.1 Pan Novák může jet do zaměstnání čtyřmi různými trasami. Čtyřikrát projel jednotlivé trasy a zaznamenal si dobu, po kterou jel do zaměstnání. Na hladině významnosti $\alpha = 0,01$ zjistěte, zda záleží na tom, kterou trasou pojede.

Cesta 1	Cesta 2	Cesta 3	Cesta 4
22	27	26	28
26	29	33	30
25	26	25	32
30	28	30	26

1.2 Učitel fyziky zkoumal, jaký vliv má druh zkušebního testu na jeho úspěšnost. Vytvořil tři typy stejně obtížných testů a náhodně je rozdělil mezi studenty ve třídě. Tabulka uvádí bodové zisky studentů v jednotlivých testech. Na hladině významnosti $\alpha = 0,05$ zjistěte, zda má typ testu vliv na úspěšnost studentů.

Typ testu		
T1	T2	T3
75	72	64
90	78	78
70	94	70
90	78	90
85		50

1.3 Ve vepřině zjišťovali, jestli váhové přírůstky vepřů závisí na použitém druhu krmiva, či nikoli. Na hladině významnosti $\alpha = 0,05$ rozhodněte, zda jsou váhové přírůstky pro různá krmiva různé, eventuálně zjistěte, který druh krmiva dává nejmenší váhové přírůstky.

Krmivo		
A	B	C
21,5	19,9	23,7
22,8	24,3	22,5
26,3	20,1	20,6
24,2	20,9	21,4
25,6	21,1	
28,1		

1.4 Výroba součástek může v podniku probíhat na jednom ze čtyř rozdílných strojů. I když každý stroj provádí stejné operace, má každý svá specifika. Na hladině významnosti $\alpha = 0,01$ testujte hypotézu o tom, že počet vyrobených součástek není ovlivněn volbou stroje.

Stroj			
A	B	C	D
93	108	123	133
98	153	143	163
80	123	150	168
88	158	165	145
60	143	140	130

1.5 Školský úřad Karviná chtěl srovnat úroveň znalostí maturantů gymnázií okresu Karviná. Za tímto účelem byl vytvořen test zahrnující otázky ze všech oblastí učiva a zadán náhodně vybraným studentům jednotlivých škol. Bodové výsledky studentů jsou uvedeny v následující tabulce.

Gymnázium Karviná	Gymnázium Český Těšín	Gymnázium Bohumín	Gymnázium Orlová	Gymnázium Havířov
79	62	74	73	86
86	54	81	67	52
49	88	64	59	61
72			76	

- a. Na hladině významnosti $\alpha = 0,05$ zjistěte, je-li průměrná úroveň maturantů jednotlivých škol stejná.
- b. Jak ovlivní výsledek průzkumu změna hladiny významnosti na 0,01?



ODPOVĚDI

- 1.1 $F = 1,0$ $F_{\text{krit}} = 5,95$ p -hodnota = 0,43 – H_0 přijímáme (doba nezávisí na trase).
- 1.2 $F = 1,43$ $F_{\text{krit}} = 3,98$ p -hodnota = 0,28 – H_0 přijímáme (typ testu nemá vliv na úspěch).
- 1.3 $F = 4,7$ $F_{\text{krit}} = 3,89$ p -hodnota = 0,03 – H_0 zamítáme (krmivo má vliv, nejvíce A).
- 1.4 $F = 15,02$ $F_{\text{krit}} = 5,29$ p -hodnota = 0,000 – H_0 zamítáme (typ stroje má vliv).
- 1.5 a) $F = 0,12$ $F_{\text{krit}} = 3,26$ p -hodnota = 0,97 – H_0 přijímáme (škola nemá vliv).
 b) $F = 0,12$ $F_{\text{krit}} = 5,41$ p -hodnota = 0,97 – H_0 přijímáme (škola nemá vliv).



SHRNUTÍ KAPITOLY

Formálně vzato je ANOVA, ať jednofaktorová nebo vícefaktorová, testem statistické hypotézy, s níž jste se seznámili v základním kurzu statistiky. Cílem ANOVA je prokázat, že hodnoty kvalitativního znaku X ovlivňují hodnoty kvantitativního znaku Y (závislého faktoru). Princip metody ANOVA, kterou prokazujeme závislost Y na X , spočívá v tom, že celkovou variabilitu měřenou součtem čtverců odchylek od celkového průměru rozdělíme na variabilitu uvnitř jednotlivých výběrů a na variabilitu mezi jednotlivými výběry.

2 ANALÝZA ROZPTYLU (ANOVA) – DVA A VÍCE FAKTORŮ

RYCHLÝ NÁHLED KAPITOLY



Jednofaktorová metoda ANOVA, kterou prokazujeme závislost znaků (faktorů) Y na X , pro něž jsou k dispozici příslušná data, spočívá v tom, že celkovou variabilitu měřenou součtem čtverců odchylek od celkového průměru rozdělíme na variabilitu uvnitř jednotlivých výběrů a na variabilitu mezi jednotlivými výběry. Cílem, k němuž směřujeme nyní, je situace, kdy budeme uvažovat, že se kromě třídění do skupin vyskytují další faktory, říkáme jim bloky, podle nichž výsledky (tj. hodnoty znaku Y) rovněž třídíme.

CÍLE KAPITOLY



Po prostudování této kapitoly budete umět:

- pomocí Excelu vypočítat analýzu rozptylu se dvěma faktory,
 - pomocí GRETLU vypočítat analýzu rozptylu se dvěma faktory,
 - použít Kruskal-Wallisovu verzi analýzy rozptylu.
-

ČAS POTŘEBNÝ KE STUDIU



K prostudování této kapitoly budete potřebovat asi 90 minut.

KLÍČOVÁ SLOVA KAPITOLY



Analýza rozptylu se dvěma faktory, Kruskal-Wallisova ANOVA.

2.1 Analýza rozptylu se dvěma faktory

ANOVA vychází z předpokladu normality rozdělení hodnot uvažovaných faktorů. Pokud U analýzy rozptylu s jedním faktorem jste uvažovali výsledky tříděné podle jistého kvalitativního znaku X do několika (konkrétně do k) skupin o rozsazích n_1, n_2, \dots, n_k . Proto v tomto případě hovoříme také o ANOVA při jednoduchém třídění neboli třídění podle jednoho faktoru. V této kapitole budeme uvažovat situaci, kdy se kromě třídění do skupin, vyskytují další faktory, říkáme jim bloky, podle nichž výsledky (tj. hodnoty znaku Y) rovněž třídíme. Přehledná situace vzniká, když kromě prvního faktoru uvažujeme ještě faktor druhý, říkáme pak, že je třídíme do bloků a v takovém případě se jedná o dvoufaktorovou ANOVA. Formálně vzato je ANOVA, ať jednofaktorová, dvoufaktorová nebo vícefaktorová, parametrickým testem statistické hypotézy, s nímž jste se seznámili v základním kurzu statistiky. Tato tzv. klasická je takový předpoklad neudržitelný, lze použít jiného typu ANOVA, tedy neparametrického testu statistické hypotézy (tento pojem si připomeňte ze základního kurzu statistiky). Konkrétně se v této kapitole seznámíte s Kruskal-Wallisovou verzí ANOVA, která využívá Chi-kvadrát test statistické hypotézy.

U analýzy rozptylu s jedním faktorem jsme uvažovali výsledky tříděné podle jistého kvalitativního znaku X do několika (konkrétně do k) skupin o rozsazích n_1, n_2, \dots, n_k . V tomto odstavci budeme uvažovat situaci, kdy se kromě třídění do skupin, vyskytuje další faktor, podle něhož výsledky (tj. hodnoty znaku Y) rovněž třídíme, říkáme, že je třídíme do bloků. Začneme výklad příkladem známým již z předchozí kapitoly.

Příklad 1. Testovacími jízdami na zkušebním okruhu se zjišťuje průměrná spotřeba paliva automobilu Octavia při použití benzínu od různých výrobců (např. Aral, Shell, Benzina, Slovnaft). Všechny testy provede jeden řidič, když s každým druhem benzínu uskuteční několik testovacích jízd, a to tak, že pro každou značku benzínu uskuteční jiný počet jízd. Zjištěné výsledky testů, tj. změřené průměrné spotřeby na 100 km, podrobíme jednofaktorové analýze rozptylu, která nám umožní zjistit, zda značka (tj. výrobce) použitého benzínu má vliv na průměrnou spotřebu automobilu.

Příklad 2. Nyní budeme uvažovat podobnou situaci, kdy výsledky testů byly získány různými řidiči (např. A, B, C, D, E, F), a to tak, že každý řidič uskutečnil jednu testovací jízdu s každou značkou benzínu (tím se myslí čerpací stanice, ze kterých pocházely pohonné hmoty). Výsledky testů proto budeme členit nejen podle značky benzínu (1. faktor), ale také podle testovacích řidičů (2. faktor). Podle předpokladů je nyní počet výsledků ve všech skupinách stejný a je roven počtu řidičů (každý řidič jel s jednou značkou benzínu jedenkrát). Zjištěné výsledky podrobíme dvoufaktorové analýze rozptylu, která umožní jednak zjistit, zda značka (tj. výrobce) použitého benzínu má vliv na průměrnou spotřebu automobilu, jednak zjistit, zda různí řidiči mají vliv na tuto spotřebu.

Příklad 3. Nyní budeme uvažovat stejnou situaci jako v příkladu 2, přitom výsledky testů byly získány různými řidiči (např. A, B, C, D, E, F), a to tak, že každý řidič uskutečnil tři testovací jízdy s každou značkou benzínu. Zjištěné výsledky podrobíme dvoufaktorové

analýze rozptylu s opakováním, která umožní jednak zjistit, zda značka (tj. výrobce) použitého benzínu má vliv na průměrnou spotřebu automobilu, jednak zjistit, zda různí řidiči mají vliv na tuto spotřebu.

Na konci této kapitoly budou všechny tři příklady podrobně analyzovány na konkrétních číselných datech. Nyní budeme postupovat ve výkladu s obecnými daty, nejprve pro případ popsany v příkladu 2. Taková data, podobně jako u jednofaktorové analýzy rozptylu, uspořádáme do přehledné Tabulky 3.

Tabulka 3: Schéma výchozí tabulky analýzy rozptylu pro dva faktory

	Hodnoty sledovaného znaku						Průměr skupiny
	Číslo bloku						
Číslo skupiny	1	2	...	j	...	r	
1	y_{11}	y_{12}	...	y_{1j}	...	y_{1r}	$\bar{y}_{1\cdot}$
2	y_{21}	y_{22}	...	y_{2j}	...	y_{2r}	$\bar{y}_{2\cdot}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
i	y_{i1}	y_{i2}	...	y_{ij}	...	y_{ir}	$\bar{y}_{i\cdot}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
k	y_{k1}	y_{k2}	...	y_{kj}	...	y_{kr}	$\bar{y}_{k\cdot}$
Průměr bloku	$\bar{y}_{\cdot 1}$	$\bar{y}_{\cdot 2}$...	$\bar{y}_{\cdot j}$...	$\bar{y}_{\cdot r}$	\bar{y}

V Tabulce 3 značíme symbolem $\bar{y}_{i\cdot}$ průměr v i -té skupině, symbolem $\bar{y}_{\cdot j}$ označujeme průměr hodnot v j -tém bloku, symbolem \bar{y} značíme celkový průměr.

Celkový součet čtverců (celkovou variabilitu) označujeme stejně, jako v (1.2), tedy:

$$S_y = \sum_{i=1}^k \sum_{j=1}^r (y_{ij} - \bar{y})^2 \quad (2.1)$$

Variabilitu mezi skupinami budeme měřit *meziskupinovým součtem čtverců* $S_{y,m}$, který definujeme následovně:

$$S_{y,m} = r \sum_{i=1}^k (\bar{y}_{i\cdot} - \bar{y})^2 \quad (2.2)$$

Meziskupinovému součtu čtverců přísluší počet stupňů volnosti $df_m = k-1$.

Variabilitu mezi bloky budeme měřit *mezblokovým součtem čtverců* $S_{y,b}$, který definujeme následovně:

$$S_{y,b} = k \sum_{j=1}^r (\bar{y}_{\cdot j} - \bar{y})^2 \quad (2.3)$$

Meziskupinovému součtu čtverců přísluší počet stupňů volnosti $df_b = r-1$.

Variabilitu uvnitř skupin označujeme jako vnitroskupinovou, nebo také reziduální a použijeme přitom označení $S_{y,v}$, přičemž definujeme *vnitroskupinový (reziduální) součet čtverců* takto

$$S_{y,v} = \sum_{i=1}^k \sum_{j=1}^r (y_{ij} - \bar{y}_{i\cdot} - \bar{y}_{\cdot j} + \bar{y})^2 \quad (2.4)$$

Vnitroskupinovému součtu čtverců přísluší počet stupňů volnosti $df_v = (k-1)(r-1)$.

Aritmetickými úpravami výše uvedených vzorců lze dokázat totiž, že *celkový součet čtverců je roven sumě meziskupinového, vnitroskupinového a blokového součtu čtverců*, symbolicky

$$S_y = S_{y,m} + S_{y,v} + S_{y,b} \quad (2.5)$$

Tento vztah se nazývá *základní vztah dvoufaktorové analýzy rozptylu*.

2.2 Předpoklady analýzy rozptylu se dvěma faktory

Předpokládáme, že faktor X_1 má k úrovní, faktor X_2 má r úrovní s efektem na znak Y , který lze vyjádřit vztahem

$$\mu_{ij} = \mu + \alpha_i + \beta_j, \quad i = 1, 2, \dots, k, \quad j = 1, 2, \dots, r, \quad (2.6)$$

kde μ_{ij} je průměr znaku Y v i -té skupině a j -tém bloku, μ je celkový průměr znaku Y , α_i je efekt hodnoty faktoru X_1 na znak Y , β_j je efekt hodnoty faktoru X_2 na znak Y .

V modelu (2.6) nejprve předpokládáme, že efekty obou faktorů na znak Y jsou *aditivní a vzájemně nezávislé*, tj. bez vzájemných interakcí. Tento předpoklad nám umožní oddělit od sebe hypotézy o efektech jednotlivých faktorů.

Formulujeme nejprve nulovou hypotézu, že všechny *skupiny* pocházejí ze stejné základní populace (základního souboru), jinak řečeno, že hodnoty faktoru X_1 nemají na hodnoty znaku Y žádný efekt (vliv). Budeme tedy v nulové hypotéze předpokládat, že α_i pocházejí z *normálně rozdělené* populace s nulovou střední hodnotou a konstantním rozptylem σ^2 .

Formulujeme nulovou hypotézu $H_0: E(\alpha_1) = E(\alpha_2) = \dots = E(\alpha_k) = 0$, proti *alternativní hypotéze*, že H_0 neplatí, že alespoň pro dvě hodnoty, např. i a j , platí: $H_1: E(\alpha_i) \neq E(\alpha_j)$.

Cílem, k němuž směřujeme, je *přijmout* nulovou hypotézu H_0 , eventuálně H_0 *zamítnout* (na zvolené hladině významnosti). Pro ověření nulové hypotézy H_0 použijeme statistiku:

$$F_1 = \frac{\frac{S_{y,m}}{k-1}}{\frac{S_{y,v}}{(k-1)(r-1)}}, \quad (2.7)$$

kteřá má při platnosti nulové hypotézy *Fisherovo rozdělení* $F(k-1, (k-1)(r-1))$. Kritické hodnoty lze nalézt v tabulkách, nebo lze využít funkce z Excelu: FINV ($\alpha; k-1; (k-1)(r-1)$).

Dále formulujeme nulovou hypotézu, že všechny *bloky* pocházejí ze stejné základní populace (základního souboru), jinak řečeno, že hodnoty faktoru X_2 nemají na hodnoty znaku Y žádný efekt. Budeme tedy v nulové hypotéze předpokládat, že β_j pocházejí z *normálně rozdělené* populace s nulovou střední hodnotou a konstantním rozptylem σ^2 .

Formulujeme nulovou hypotézu $H_0': E(\beta_1) = \dots = E(\beta_r) = 0$, proti *alternativní hypotéze*, že H_0' neplatí, že alespoň pro dvě hodnoty, např. $i' \neq i''$, platí $H_1': E(\beta_{i'}) \neq E(\beta_{i''})$.

Pro ověření nulové hypotézy H_0' použijeme statistiku:

$$F_2 = \frac{\frac{S_{y,b}}{r-1}}{\frac{S_{y,v}}{(k-1)(r-1)}}, \quad (2.8)$$

kteřá má při platnosti nulové hypotézy *Fisherovo rozdělení* $F(r-1, (k-1)(r-1))$.

Zásadní rozdíl mezi dvoufaktorovou a jednofaktorovou analýzou rozptylu spočívá v tom, že u jednofaktorové ANOVA *neuvažujeme* působení dalšího faktoru, zatímco u dvoufaktorové ANOVA tak činíme. Tento rozdíl je vyjádřen ve výpočtu testového kritéria (2.7) a (2.8), kde se ve jmenovateli zlomku vyskytuje člen $(k-1)(r-1)$. Kdybychom na stejnou situaci aplikovali pouze jednofaktorovou ANOVA, pak by ve výpočtu hodnoty testového kritéria podle vztahu (1.6) byl na stejném místě člen $(n-k)$ nebo člen $(n-r)$, podle toho, zda bychom brali v úvahu skupiny nebo bloky. Tento rozdíl může zapříčinit rozdílné výsledky získané jednofaktorovou nebo dvoufaktorovou ANOVA!

ŘEŠENÁ ÚLOHA 2.1



Testovacími jízdami na zkušební okruhu se zjišťuje průměrná spotřeba benzínu Natural 95 automobilu Octavia při použití benzínu od různých výrobců (Aral, Shell, Benzina, Slovnaft). Bylo vybráno 6 řidičů A, B, C, D, E, F, z nichž každý absolvoval s každým typem benzínu jednu zkušební jízdu. Na hladině významnosti 0,05 testujte, je-li průměrná spotřeba paliva závislá na typu použitého benzínu a na tom, který řidič s vozem jel.

Značka benzínu	Řidiči					
	A	B	C	D	E	F
Aral	7,5	6,9	7,9	7,3	6,9	7,8
Shell	7,6	7,2	7,5	8,0	7,3	8,2
Benzina	7,2	8,1	7,8	7,6	7,8	6,9
Slovnaft	7,0	7,3	7,2	7,5	8,2	7,7

Řešení:

Máte za úkol prozkoumat závislost průměrné spotřeby (znak Y) na typu použitého benzínu (znak X_1) a na řidiči (znak X_2), který s vozem jel. Znak X_1 má $k = 4$ skupiny, znak X_2 má $r = 6$ bloků.

Pro faktor X_1 formulujeme nulovou hypotézu:

$$H_0: E(\alpha_1) = E(\alpha_2) = E(\alpha_3) = E(\alpha_4), \quad (2.9)$$

proti H_1 : neplatí (2.9), tj. průměrná spotřeba závisí na použitém druhu benzínu.

Pro faktor X_2 formulujeme nulovou hypotézu

$$H'_0: E(\beta_1) = E(\beta_2) = \dots = E(\beta_6), \quad (2.10)$$

proti alternativní hypotéze H'_1 : neplatí (2.10), tj. průměrná spotřeba benzínu závisí na řidiči, který s vozem jel.

Pro ověření těchto hypotéz, tj. pro výpočet testových kritérií, musíme znát hodnotu součtů $S_{y,m}$, $S_{y,v}$ a S_y . Vypočítáme podmíněné průměry $\bar{y}_i, i = 1, 2, 3, 4, \bar{y}_j, j = 1, 2, \dots, 6$ (výpočty jsou v Tabulce 4) a také celkový průměr \bar{y} .

$$\bar{y}_1 = \frac{7,5 + 6,9 + \dots + 7,8}{6} = 7,38, \text{ další průměry } \bar{y}_2, \bar{y}_3, \bar{y}_4, \text{ vypočítáme analogicky.}$$

$$\bar{y}_{.1} = \frac{7,5 + 7,6 + 7,2 + 7}{4} = 7,33, \text{ další průměry } \bar{y}_{.2}, \dots, \bar{y}_{.6} \text{ vypočítáme analogicky.}$$

$$\text{Celkový průměr je } \bar{y} = \frac{7,5 + 6,9 + \dots + 7,7}{24} = 7,50.$$

Hodnoty všech průměrů jsou uvedeny v tabulce. Nyní lze přistoupit k výpočtu jednotlivých součtů:

$$S_{y,m} = r \sum_{i=1}^4 (\bar{y}_i - \bar{y})^2 = 6 \cdot [(7,38 - 7,5)^2 + \dots + (7,48 - 7,5)^2] = 0,21.$$

$$S_{y,b} = k \sum_{j=1}^6 (\bar{y}_{.j} - \bar{y})^2 = 4 \cdot [(7,33 - 7,5)^2 + \dots + (7,38 - 7,5)^2] = 0,35.$$

Potřebujeme znát i hodnotu součtu $S_{y,v}$, z praktického hlediska je však výhodnější vypočítat hodnotu součtu S_y . Součet $S_{y,v}$ pak snadno dopočítáme, neboť $S_y = S_{y,m} + S_{y,v} + S_{y,b}$.

$$S_y = \sum_{i=1}^4 \sum_{j=1}^6 (y_{i,j} - \bar{y})^2 = (7,5 - 7,5)^2 + (6,9 - 7,5)^2 + \dots + (7,8 - 7,5)^2 + \\ + (7,6 - 7,5)^2 + \dots + (8,2 - 7,5)^2 + \dots + (7,7 - 7,5)^2 = 3,79.$$

Potom vypočítáme $S_{y,v} = S_y - S_{y,m} - S_{y,b} = 3,79 - 0,21 - 0,36 = 3,22$.

Pro ověření hypotézy H_0 určíme testové kritérium F_1

$$F_1 = \frac{\frac{S_{y,m}}{k-1}}{\frac{S_{y,v}}{(k-1)(r-1)}} = \frac{0,21}{\frac{3}{3 \cdot 5}} = 0,32.$$

V tabulce kritických hodnot F -rozdělení nebo pomocí Excelu najdeme $F_{0,05}(3,15) = \text{FINV}(0,05; 3,15) = 3,29$.

Protože $0,32 < 3,29$, přijímáme H_0 , což znamená, že použitá značka benzínu nemá na průměrnou spotřebu vliv.

Pro ověření hypotézy H'_0 určíme testové kritérium F_2

$$F_2 = \frac{\frac{S_{y,b}}{r-1}}{\frac{S_{y,v}}{(k-1)(r-1)}} = \frac{0,36}{\frac{3,22}{3 \cdot 5}} = 0,33.$$

V tabulce kritických hodnot F -rozdělení nebo pomocí Excelu najdeme $F_{0,05}(5,15) = \text{FINV}(0,05; 5,15) = 2,9$. Protože $0,33 < 2,9$, přijímáme i hypotézu H'_0 , tzn., že ani volba řidiče nemá na průměrnou spotřebu statisticky významný vliv.

Na rozdíl od jednofaktorové ANOVA jsme zde v obou situacích uvažovali současné působení *dvou faktorů!*

Tabulka 4: Podmíněné průměry

Zn. benzínu	Řidiči						Průměry
	A	B	C	D	E	F	
Aral	7,5	6,9	7,9	7,3	6,9	7,8	7,38
Shell	7,6	7,2	7,5	8,0	7,3	8,2	7,63
Benzina	7,2	8,1	7,8	7,6	7,8	6,9	7,57
Slovnaft	7,0	7,3	7,2	7,5	8,2	7,7	7,48
Průměry	7,33	7,38	7,6	7,6	7,55	7,65	7,50

Nakonec si ještě ukážeme řešení pomocí Excelu. Využijeme přitom funkci menu:

Nástroje → Analýza dat... → ANOVA: dva faktory bez opakování

Nejprve je zapotřebí připravit v Excelu data. Jednotlivé hodnoty y_{ij} pro faktorů Y pro hodnoty faktorů $X_1 = \text{benzín}$ a $X_2 = \text{řidič}$ uspořádáme do řádků a sloupců, podobně jako v tabulce v zadání. Data ve worksheetu Excelu vypadají tedy například takto:

	A	B	C	D	E	F	G	I
1	benzín/řidič	A	B	C	D	E	F	
2	Aral	7,5	6,9	7,9	7,3	6,9	7,8	
3	Shell	7,6	7,2	7,5	8	7,3	8,2	
4	Benzina	7,2	8,1	7,8	7,6	7,8	6,9	
5	Slovnaft	7	7,3	7,2	7,5	8,2	7,7	
6								

Dále otevřeme v hlavním menu postupně položky:

Data → Analýza dat... → ANOVA: dva faktory bez opakování

Po volbě třetí položky ANOVA: dva faktory bez opakování, se otevře zadávací okno:

Vstupní oblast: \$A\$1:\$G\$5

Popisky – zakliknete

analýza rozptylu (ANOVA) – DVA a více faktorů

Alfa: 0,05 (hladina významnosti je předvolena, lze ji však změnit). Potvrdíte OK.

Anova: dva faktory bez opakování

Faktor	Počet	Součet	Průměr	Rozptyl
Aral	6	44,3	7,383333	0,185667
Shell	6	45,8	7,633333	0,154667
Benzina	6	45,4	7,566667	0,194667
Slovnaft	6	44,9	7,483333	0,181667
A	4	29,3	7,325	0,075833
B	4	29,5	7,375	0,2625
C	4	30,4	7,6	0,1
D	4	30,4	7,6	0,086667
E	4	30,2	7,55	0,323333
F	4	30,6	7,65	0,296667

ANOVA

Zdroj variability	SS	Rozdíl	MS	F	Hodnota P	F krit
Řádky	0,21	3	0,07	0,325581	0,806868	3,287383
Sloupce	0,358333	5	0,071667	0,333333	0,884913	2,901295
Chyba	3,225	15	0,215			
Celkem	3,793333	23				

V první tabulce jsou uvedeny základní statistické údaje o datech: Faktor, Počet, Součet, Průměr a Rozptyl.

Ve druhé tabulce nazvané ANOVA jsou uvedeny výpočty metodou ANOVA: dva faktory bez opakování, jednotlivé položky mají následující význam:

Řádky = meziskupinový

Sloupce = vnitroskupinový

Chyba = meziblokový

Celkem = celkový

SS = Součet čtverců (Sum of Squares)

Rozdíl = stupeň volnosti (DF – Degree of Freedom)

MS = Průměr čtverců (Mean Square)

F = testové kritérium

Hodnota P = Signifikance (p-hodnota)

F krit = kritická hodnota rozdělení F

Hodnoty získané řešením v Excelu jsou stejné jako při použití „ručního“ výpočtu, proto i závěry jsou stejné. V Excelu máme navíc vypočtenou p-hodnotu testu (tzv. signifikanci), která, pokud je menší než zvolená hladina významnosti α , znamená, že nulovou hypotézu zamítáme. V opačném případě nulovou hypotézu přijímáme.

V předchozích úvahách jsme měli situaci právě jednoho výskytu všech kombinací hodnot skupin a bloku obou uvažovaných faktorů. Například každý řidič absolvoval jedinou jízdu s každým typem benzínu. Dále budeme uvažovat situaci vícenásobného opakování všech kombinací hodnot skupin a bloku obou uvažovaných faktorů. Například každý řidič absolvuje několik jízd (například 3 jízdy – viz řešená úloha 2.2) s každým typem benzínu, přitom samozřejmě mohou být dosažené hodnoty průměrné spotřeby různé. Zda se tyto

výsledky odlišují výrazně či nikoliv, se opět zjišťuje statistickým testem. Podrobnou analýzu situace, která je analogická analýze případu bez opakování, již zde uvádět nebudeme. Omezíme se pouze na řešení příkladu s využitím Excelu, konkrétně položky ANOVA: dva faktory s opakováním (řešená úloha 2.2)

ŘEŠENÁ ÚLOHA 2.2



Podobně jako v příkladu 2.1 se zjišťuje průměrná spotřeba benzínu Natural 95 automobilu Octavia při použití benzínu od různých výrobců (Aral, Shell, Benzina, Slovnaft). Bylo vybráno 6 řidičů A, B, C, D, E, F, z nichž každý absolvoval s každým typem benzínu tři zkušební jízdy. Na hladině významnosti 0,05 testujte, je-li průměrná spotřeba paliva závislá na typu použitého benzínu a na řidiči. Údaje jsou uvedeny v následující Tabulce 5.

Tabulka 5: Analýza rozptylu se dvěma faktory s opakováním

benzin/řidič	Aral	Shell	Benzina	Slovnaft
A	7,5	7,6	7,2	7
	7,7	7,4	7,6	7,4
	8	7,3	8,1	7,7
B	6,9	7,2	8,1	7,3
	6,7	7,4	8,5	7,6
	6,6	7,6	8,8	7,8
C	7,9	7,5	7,8	7,2
	8	7,8	7,7	7,1
	8,3	8,1	7,6	7
D	7,3	8	7,6	7,5
	7,2	8	7,8	7,7
	7,1	7,9	8	7,8
E	6,9	7,3	7,8	8,2
	6,8	7,2	8	8,1
	6,7	7	8,1	8
F	7,8	8,2	6,9	7,7
	7,7	8,4	7,5	7,7
	7,5	8,5	7,9	7,7

Řešení:

Data ve worksheetu Excelu vypadají přesně tak jako v Tabulce 5, jsou umístěny např. v poli A1 až E19. Dále otevřeme v hlavním menu postupně položky:

Data → Analýza dat... → ANOVA: dva faktory s opakováním

Po volbě druhé položky ANOVA: dva faktory s opakováním, se otevře zadávací okno, kde postupně zadáte:

Vstupní oblast: \$A\$1:\$E\$19

Řádků na výběr: 3 (tj. počet opakování)

Alfa: 0,05 (hladina významnosti je předvolena, lze ji však změnit)

Výstupní oblast: např. \$L\$1 (levý horní roh výstupní oblasti)

Potvrdíte OK.

analýza rozptylu (ANOVA) – DVA a více faktorů

Obdržíte následující výstup, kterého “levý horní roh” začíná v buňce L1 nadpisem ANOVA: dva faktory s opakováním. V první tabulce jsou uvedeny základní statistické údaje o datech: Faktor, Počet, Součet, Průměr a Rozptyl.

Anova: dva faktory s opakováním

Faktor	Aral	Shell	Benzina	Slovnaft	Celkem
A					
Počet	3	3	3	3	12
Součet	23,2	22,3	22,9	22,1	90,5
Průměr	7,73	7,43	7,63	7,37	7,54
Rozptyl	0,06	0,02	0,20	0,12	0,10
B					
Počet	3	3	3	3	12
Součet	20,2	22,2	25,4	22,7	90,5
Průměr	6,73	7,40	8,47	7,57	7,54
Rozptyl	0,02	0,04	0,12	0,06	0,46
C					
Počet	3	3	3	3	12
Součet	24,2	23,4	23,1	21,3	92
Průměr	8,07	7,80	7,70	7,10	7,67
Rozptyl	0,04	0,09	0,01	0,01	0,16
D					
Počet	3	3	3	3	12
Součet	21,6	23,9	23,4	23	91,9
Průměr	7,200	7,967	7,800	7,667	7,658
Rozptyl	0,010	0,003	0,040	0,023	0,103
E					
Počet	3	3	3	3	12
Součet	20,4	21,5	23,9	24,3	90,1
Průměr	6,80	7,17	7,97	8,10	7,51
Rozptyl	0,01	0,02	0,02	0,01	0,33
F					
Počet	3	3	3	3	12
Součet	23	25,1	22,3	23,1	93,5
Průměr	7,67	8,37	7,43	7,70	7,79
Rozptyl	0,02	0,02	0,25	0,00	0,19
Celkem					
Počet	18	18	18	18	
Součet	132,6	138,4	141	136,5	
Průměr	7,37	7,69	7,83	7,58	
Rozptyl	0,28	0,20	0,19	0,13	

Ve druhé tabulce nazvané ANOVA jsou uvedeny výpočty metodou ANOVA: dva faktory s opakováním.

ANOVA						
Zdroj variability	SS	Rozdíl	MS	F	Hodnota P	F krit
Výběr	0,69	5	0,14	2,64	0,03	2,41
Sloupce	2,08	3	0,69	13,23	0,00	2,80
Interakce	10,23	15	0,68	12,99	0,00	1,88
Dohromady	2,52	48	0,05			
Celkem	15,53	71				

Jednotlivé položky mají následující význam:

Výběr = meziskupinový

Sloupce = vnitroskupinový
Interakce = meziblokový
Celkem = celkový
SS = Součet čtverců (Sum of Squares)
Rozdíl = stupeň volnosti (DF – Degree of Freedom)
MS = Průměr čtverců (Mean Square)
F = testové kritérium
Hodnota P = Signifikance (p-hodnota)
F krit = kritická hodnota rozdělení F

Hodnoty získané řešením v Excelu jsou analogické jako v příkladu 2.1, tedy v případě ANOVA bez opakování. Navíc je tu p-hodnota uvedena v řádce Interakce, která se týká testu vzájemné závislosti faktorů. Nulová hypotéza předpokládá, že faktorů jsou vzájemně nezávislé. Pokud je tato hodnota menší než zvolená hladina významnosti α , znamená to, že nulovou hypotézu zamítáme. V opačném případě nulovou hypotézu přijímáme.

V této kapitole jsme uvažovali situaci, kdy se kromě třídění do skupin vyskytují další faktory, říkáme jim bloky, podle nichž výsledky (tj. hodnoty znaku Y) rovněž třídíme. Přehledná situace vzniká, když kromě prvního faktoru uvažujeme ještě faktor druhý, říkáme pak, že je třídíme do bloků a v takovém případě se jedná o dvoufaktorovou ANOVA. Formálně vzato je ANOVA, ať jednofaktorová, dvoufaktorová nebo vícefaktorová, parametrickým testem statistické hypotézy, s nímž jste se seznámili v základním kurzu statistiky. Nejprve jsme měli situaci právě jednoho výskytu všech kombinací hodnot skupin a bloku obou uvažovaných faktorů. Například každý řidič absolvoval jedinou jízdu s každým typem benzínu. Poté jsme uvažovali situaci vícenásobného opakování všech kombinací hodnot skupin a bloku obou uvažovaných faktorů. Například každý řidič absolvuje několik jízd s každým typem benzínu, přitom samozřejmě mohou být dosažené hodnoty průměrné spotřeby různé. Zda se tyto výsledky odlišují výrazně či nikoliv, se opět zjistilo statistickým testem. K řešení příkladů jsme použili Excel, konkrétně položku Analýza dat, podobně budeme postupovat v řešené úloze 2.3. A v řešené úloze 2.4 si ukážeme řešení dvoufaktorové analýzy rozptylu v programu GRETl.

ŘEŠENÁ ÚLOHA 2.3



Po půlročním zkušebním období firmy „Dům a zahrada“ bylo vybráno 12 obchodů (6 internetových obchodů, 6 kamenných prodejen) se sortimentem zahrada, dům a byt, dílna a náradí, a byly zaznamenány tržby z prodeje v jednotlivých sortimentech. Testujte na hladině významnosti 0,05; zda je výše tržeb ovlivněna typem obchodu nebo sortimentem.

analýza rozptylu (ANOVA) – DVA a více faktorů

prodej\kategorie zboží	ZAHRAIDA	DŮM A BYT	DÍLNA A NÁŘADÍ
INTERNETOVÝ OBCHOD	44	48	30
	42	34	18
	52	35	75
	70	2	70
	35	41	62
	20	33	68
KAMENNÝ OBCHOD	33	38	30
	12	1	42
	13	50	18
	22	5	27
	64	44	34
	35	47	30

Řešení:

Testujeme tedy hypotézy: H_0 : výše tržeb není ovlivněna nabízeným sortimentem,
 H_1 : výše tržeb je ovlivněna nabízeným sortimentem.

A dále hypotézy: H'_0 : výše tržeb není ovlivněna typem obchodu,
 H'_1 : výše tržeb je ovlivněna typem obchodu.

V Excelu zvolíme následující posloupnost příkazů: Data → Analýza dat... → ANOVA: dva faktory s opakováním. Vstupní oblast musí obsahovat i záhlaví tabulky a v každé cele musí být stejný počet hodnot.

Po volbě druhé položky ANOVA: dva faktory s opakováním, se otevře zadávací okno, kde postupně zadáte:

Obdržíte následující zkrácený výstup (pouze tabulka ANOVA):

ANOVA						
Zdroj variability	SS	Rozdíl	MS	F	Hodnota P	F krit
Výběr	1521	1	1521	4,530831	0,041616	4,170877
Sloupce	661,5556	2	330,7778	0,985337	0,385072	3,31583
Interakce	752,6667	2	376,3333	1,121041	0,339213	3,31583
Dohromady	10071	30	335,7			

Výběr – meziskupinový SS (1.faktor), Sloupce – meziblokový SS (2.faktor), Interakce – SS pro interakce mezi faktory 1 a 2, Dohromady – vnitroskupinový SS. Z výše uvedeného výstupu tedy vidíme, že není rozdíl mezi nabízeným sortimentem (sloupce), ale je rozdíl mezi typem obchodu (řádky – výběr).

Tedy H_0 : výše tržeb není ovlivněna nabízeným sortimentem, *nelze zamítnout*, protože *Hodnota P = 0,385 což je větší než hladina významnosti 0,05*, na které testujeme. Proto nemůžeme tvrdit, že by mezi výši tržeb a nabízeným sortimentem byla závislost.

V případě nulové hypotézy: H'_0 : výše tržeb není ovlivněna typem obchodu, vidíme, že *hodnota P = 0,041 což je menší než hladina významnosti 0,05*, na které testujeme, proto *nulovou hypotézu zamítáme*. A tedy můžeme tvrdit, že výše tržeb je z 95 % ovlivněna typem obchodu.

V případě Interakce nulová hypotéza předpokládá, že faktory jsou vzájemně nezávislé. Protože *Hodnota P = 0,339 je větší než hladina významnosti 0,05*; nulovou hypotézu nelze zamítnout, a tedy nelze tvrdit, že by faktory (nabízený sortiment a typ obchodu) byly závislé.

ŘEŠENÁ ÚLOHA 2.4



Ve třech městech okresu Karviná jsme v jednotlivých dnech sledovali průměrnou spotřebu pitné vody (v m^3) na jednoho obyvatele. Zjistěte, zda je průměrná spotřeba vody závislá na dni v týdnu, a je-li spotřeba v různých městech různá. Uvažujte hladinu významnosti 0,05. Zjištěné údaje jsou uvedeny v Tabulce 6.

Tabulka 6: Spotřeba pitné vody (m^3)

	Karviná	Petřvald	Bohumín
Po	0,6	0,7	0,5
Út	0,7	0,6	0,6
St	0,9	0,8	0,7
Čt	0,6	0,6	0,5
Pá	1	1,3	0,8
So	1,2	1,6	1,3
Ne	1	1,2	1,3

Řešení:

Formulace první dvojice hypotéz:

H_0 : spotřeba pitné vody nezávisí na dnu v týdnu,

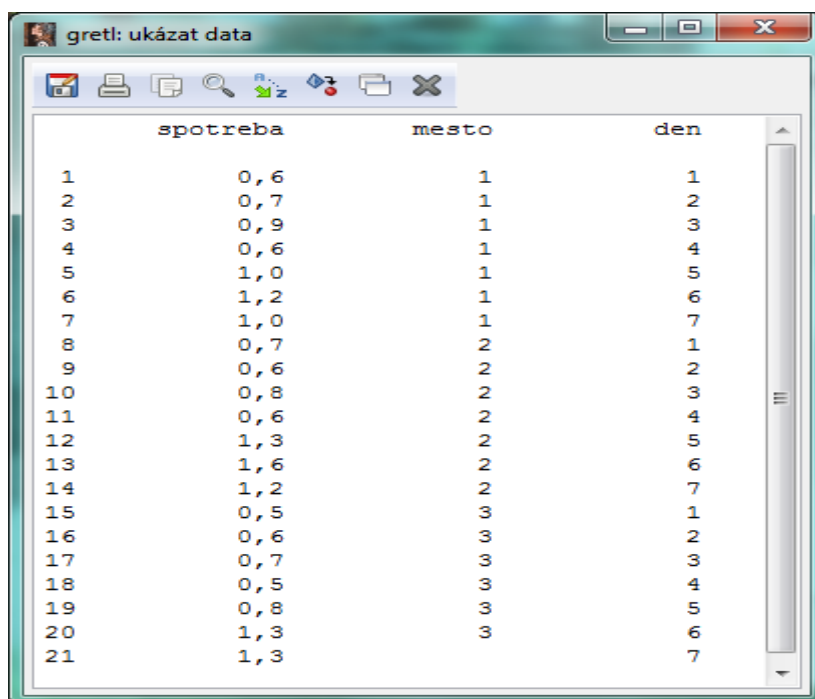
H_1 : spotřeba pitné vody závisí na dnu v týdnu.

Formulace druhé dvojice hypotéz:

H_0 : spotřeba pitné vody nezávisí na městě,

H_1 : spotřeba pitné vody závisí na městě.

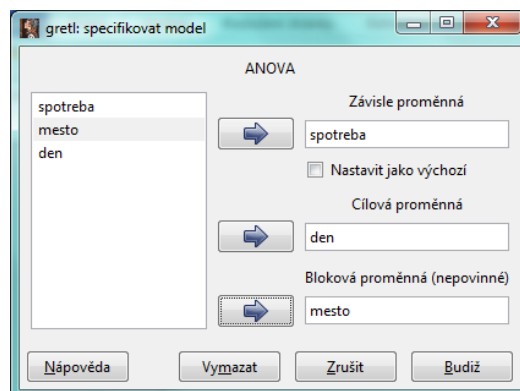
Obrázek 7 zachycuje zadávání hodnot do programu GRETL. V prvním sloupci je kvantitativní proměnná spotřeba vody, druhý sloupec zobrazuje město (1,2,3) a třetí sloupec je proměnná den (1,2,3,4,5,6,7). Kvalitativní proměnné musí být přiřozená čísla.



	spotreba	mesto	den
1	0,6	1	1
2	0,7	1	2
3	0,9	1	3
4	0,6	1	4
5	1,0	1	5
6	1,2	1	6
7	1,0	1	7
8	0,7	2	1
9	0,6	2	2
10	0,8	2	3
11	0,6	2	4
12	1,3	2	5
13	1,6	2	6
14	1,2	2	7
15	0,5	3	1
16	0,6	3	2
17	0,7	3	3
18	0,5	3	4
19	0,8	3	5
20	1,3	3	6
21	1,3	3	7

Obrázek 7: Zadávání hodnot do programu GRETL

Testování první dvojice hypotéz.



Obrázek 8: Testování první dvojice hypotéz

gretl: ANOVA

Analýza rozptylu, reakce = spotreba, úprava = den:

	Součet čtverců	df	Střední kvadrát
Úprava	1,75905	6	0,293175
Blok	0,092381	2	0,0461905
Residuum	0,220952	12	0,0184127
Úplné	2,07238	20	0,103619

$F(6, 12) = 0,293175 / 0,0184127 = 15,9224$ [p-hodnota 4,42e-005]

Obrázek 9: Výsledek testování první dvojice hypotéz

Výsledek: p -hodnota = $4,42 \cdot 10^{-5}$ a tato hodnota je menší než hladina významnosti 0,05, proto nulovou hypotézu o nezávislosti spotřeby pitné vody na dnu v týdnu zamítáme. Můžeme tedy tvrdit, že spotřeba pitné vody z 95% závisí na dnu v týdnu.

Testování druhé dvojice hypotéz.

gretl: specifikovat model

ANOVA

Závisle proměnná: spotreba

Nastavit jako výchozí

Cílová proměnná: mesto

Bloková proměnná (nepovinné): den

Nápověda Vymazat Zrušit Budiz

Obrázek 10: Testování druhé dvojice hypotéz

gretl: ANOVA

Analýza rozptylu, reakce = spotreba, úprava = mesto:

	Součet čtverců	df	Střední kvadrát
Úprava	0,092381	2	0,0461905
Blok	1,75905	6	0,293175
Residuum	0,220952	12	0,0184127
Úplné	2,07238	20	0,103619

$F(2, 12) = 0,0461905 / 0,0184127 = 2,50862$ [p-hodnota 0,1230]

Obrázek 11: Výsledek testování druhé dvojice hypotéz

Výsledek: p -hodnota = 0,123 a tato hodnota není menší než hladina významnosti 0,05, proto nulovou hypotézu o nezávislosti spotřeby pitné vody na městě nelze zamítnout. Z 95% nebylo prokázáno, že by spotřeba pitné vody závisela na městě.

2.3 Kruskal – Wallisova analýza rozptylu

Analýza rozptylu předpokládá ve své parametrické podobě normalitu rozdělení a homoskedasticitu (identické rozptyly). Pokud tyto podmínky nejsou splněny, je třeba použít neparametrický Kruskal-Wallisův test, který je obdobou jednofaktorového třídění v analýze rozptylu. Na rozdíl od parametrického testu nepředpokládá normalitu rozdělení, jeho nevýhodou je pak menší citlivost. Kruskal-Wallisův test je vícevýběrovým testem mediánů.

Nechť tyto náhodné výběry pochází ze spojitých rozdělení stejného typu a stejných rozptylů (homoskedasticita): $(X_{11}, X_{12}, \dots, X_{1n_1}); (X_{21}, X_{22}, \dots, X_{2n_2}); \dots; (X_{k1}, X_{k2}, \dots, X_{kn_k})$; kde n_i je rozsah jednotlivých výběrů.

Testujeme nulovou hypotézu: $H_0: \tilde{x}_1 = \tilde{x}_2 = \dots = \tilde{x}_k$, proti alternativní hypotéze: H_1 : neplatí H_0 .

Všechny veličiny X_{ij} tvoří dohromady sdružený náhodný výběr o rozsahu $N = \sum_{i=1}^k n_i$. Z tohoto výběru vytvoříme uspořádaný výběr (rostoucí posloupnost), a určí se pořadí R_{ij} každé veličiny X_{ij} . Tato pořadí uspořádáme do tabulky a určíme tzv. součty pořadí pro jednotlivé výběry T_i , kde $T_i = \sum_{j=1}^{n_j} R_{ij}$.

$$\text{Testová statistika je: } Q = \frac{12}{N \cdot (N + 1)} \cdot \sum_{i=1}^k \frac{T_i^2}{n_i} \quad 3 \cdot (N + 1).$$

Hodnotu Q porovnááme s kritickou hodnotou $\chi_{\alpha}^2(k - 1)$.



ŘEŠENÁ ÚLOHA 2.5

V následující tabulce jsou uvedeny ceny bytů v závislosti na počtu pokojů. Pomocí Kruskal-Wallisovy analýzy rozptylu zjistíte, zda je cena bytu závislá na počtu pokojů v bytě. Uvažujte hladinu významnosti 0,05.

Počet pokojů	Cena bytu v tis.Kč			
1	200	210	220	
2	320	310	330	340
3	500	520	540	510
4	600	620	610	

Řešení:

V další tabulce se zapíše pořadí R_{ij} každé veličiny X_{ij} a dále určíme tzv. součty pořadí pro jednotlivé výběry T_i .

Počet pokojů	R _{ij}				T _i	n _i
1	1	6	3		6	3
2	5	22	6	7	22	4
3	8	38	11	9	38	4
4	12	39	13		39	3

Tabulka pro výpočet testového kritéria

T _i	T _i ²	T _i ² / n _i
6	36	12
22	484	121
38	1444	361
39	1521	507
SUMA		1001

Dosadíme do testové statistiky $Q = \frac{12}{14 \cdot (14 + 1)} \cdot 1001 = 3 \cdot (14 + 1) = 12,2$.

Kritická hodnota $\chi_{0,05}^2(3) = CHINV(0,05; 3) = 7,81$.

Protože hodnota testové statistiky $Q = 12,2$ leží v kritickém oboru, tak nulovou hypotézu o nezávislosti znaků zamítáme. Můžeme tedy z 95 % tvrdit, že cena bytu závisí na počtu pokojů v bytě.

SAMOSTATNÉ ÚKOLY



Řešte v Excelu.

2.1 Ve čtyřech městech okresu Karviná jsme v jednotlivých dnech sledovali průměrnou spotřebu pitné vody (v m³) na jednoho obyvatele. Zjistěte, zda je průměrná spotřeba vody závislá na dni v týdnu, a je-li spotřeba v různých městech různá. Uvažujte hladinu významnosti 0,01. Zjištěné údaje jsou uvedeny v tabulce.

	Karviná	Orlová	Bohumín	Český Těšín
Po	0,64	0,75	0,54	0,76
Út	0,78	0,63	0,61	0,83
St	0,93	0,82	0,7	0,91
Čt	0,66	0,62	0,56	0,62
Pá	0,99	1,3	0,79	0,99
So	1,22	1,65	1,3	0,98
Ne	1,05	1,3	1,24	1,1

2.2 Výroba součástek může v podniku probíhat na jednom ze čtyř rozdílných strojů. I když každý stroj provádí stejné operace, má svá specifika. U každého stroje pracuje jeden dělník. Na hladině významnosti $\alpha = 0,01$ testujte hypotézu o tom, že počet vyrobených součástek není ovlivněn volbou stroje ani dělníkem, který na něm pracuje.

Dělník	Stroj			
	A	B	C	D
1	93	108	123	133
2	98	153	143	163
3	80	123	150	168
4	88	158	165	145
5	60	143	140	130

2.3 V následující tabulce jsou uvedeny průměrné bodové výsledky z matematiky na šesti vybraných školách v členských státech Víšegrádské skupiny (Česká republika, Slovensko, Maďarsko, Polsko). Pomocí Kruskal-Wallisovy analýzy rozptylu zjistěte, zda se vědomostní úroveň v matematice liší v jednotlivých státech V4. Uvažujte hladinu významnosti 0,05.

Česká republika	Slovensko	Maďarsko	Polsko
55,4	68,4	52,1	62,3
61,2	57,9	58,9	61,2
65,8	56,2	63,4	51,6
59,3	54,3	54,2	54,7
62,5	52,6	56,8	61,5
58,4	61,2	42,6	66,1



ODPOVĚDI

- 2.1 DNY: $F = 12,95$ $F_{\text{krit}} = 4,01$ $p\text{-hodnota} = 0,000$ – H_0 zamítáme (průměrná spotřeba pitné vody závisí na dnu v týdnu)
 MĚSTO: $F = 2,07$ $F_{\text{krit}} = 5,1$ $p\text{-hodnota} = 0,14$ – H_0 přijímáme (nebyla prokázána závislost průměrné spotřeby pitné vody na městě).
- 2.2 DĚLNÍK: $F = 2,45$ $F_{\text{krit}} = 5,41$ $p\text{-hodnota} = 0,1$ – H_0 přijímáme (nebyla prokázána závislost počtu součástek na dělníkovi, který na stroji pracuje).
 STROJ: $F = 20,47$ $F_{\text{krit}} = 5,95$ $p\text{-hodnota} = 0,000$ – H_0 zamítáme (počet vyrobených součástek závisí na stroji).
- 2.3 $N = 24$; $T = (92; 70; 53; 85)$; statistika $Q = 2,99$; kritická hodnota = 7,81; nulovou hypotézu o nezávislosti bodového výsledku na státu nezamítáme (soubory, z nichž pocházejí výběry jsou shodné)

SHRNUTÍ KAPITOLY



V této kapitole jsme uvažovali situaci, kdy se kromě třídění do skupin, vyskytovaly další faktory, říkáme jim bloky. Když kromě prvního faktoru uvažujeme ještě faktor druhý, říkáme pak, že je třídíme do bloků a v takovém případě se jedná o dvoufaktorovou ANOVA. Formálně vzato je ANOVA, ať jednofaktorová, dvoufaktorová nebo vícefaktorová, parametrickým testem statistické hypotézy, s nímž jste se seznámili v základním kurzu statistiky. V této kapitole jste se také seznámili s Kruskal-Wallisovou verzí ANOVA, která využívá Chi-kvadrát test statistické hypotézy.

3 REGRESNÍ ANALÝZA – JEDNOROZMĚRNÁ LINEÁRNÍ REGRESE



RYCHLÝ NÁHLED KAPITOLY

Analýzu rozptylu z první kapitoly je možné chápat jako analýzu závislosti kvantitativního znaku (proměnné) na kvalitativním znaku (proměnné). Naproti tomu závislostí kvantitativního znaku na kvantitativním znaku (nebo více kvantitativních znacích) se zabývá *regresní analýza*. V případě závislosti dvou znaků mluvíme o *jednorozměrné regresi* (případně *jednoduché regresi*), u znaku závislém na více kvantitativních veličinách hovoříme o *vícerozměrné regresi* (*vícenásobné regresi*). V této kapitole budeme vyšetřovat nejprve nejjednodušší *lineární* závislost dvou znaků, v další kapitole se budeme zabývat i nelineárními závislostmi dvou znaků důležitých z hlediska ekonomických aplikací.



CÍLE KAPITOLY

Po prostudování této kapitoly budete umět:

- vypočítat regresní koeficienty a vysvětlit metodu nejmenších čtverců,
 - vypočítat koeficient determinace a koeficient korelace,
 - vyjmenovat podmínky klasického lineárního regresního modelu.
-



ČAS POTŘEBNÝ KE STUDIU

K prostudování této kapitoly budete potřebovat asi 90 minut.



KLÍČOVÁ SLOVA KAPITOLY

Regresní přímka, metoda nejmenších čtverců, koeficient determinace, koeficient korelace.

3.1 Regresní analýza

V regresní analýze studujeme vztah mezi jedinou proměnnou (hodnotami statistického znaku) nazývanou *závisle proměnnou* (někdy *vysvětlovanou proměnnou*), označujeme ji Y , a obecně několika proměnnými (hodnotami statistických znaků), které nazýváme *nezávisle proměnné* (někdy *vysvětlující proměnné*), a označujeme je symboly X_1, X_2, \dots . Pokud se zabýváme jedinou nezávisle proměnnou X , hovoříme o *jednoduché regresi*, pokud je nezávisle proměnných více než jedna, mluvíme o *vícerozněrné (vícenásobné) regresi* (někdy též mnohonásobné regresi). V této a následující kapitole se věnujeme jednoduché regresi.

Závisí-li veličina Y na veličině X , pak to matematicky vyjadřujeme zápisem

$$Y = f(X). \quad (3.1)$$

V našem případě jsou Y a X *statistické znaky* (náhodné veličiny), pak hovoříme o *statistické závislosti*, funkční vztah (3.1) přejde v *regresní vztah (regresní model)*

$$y = f(x) + \varepsilon, \quad (3.2)$$

kde y , resp. x , představují hodnoty znaku Y , resp. X , ε je *náhodná složka*, funkci f nazýváme *regresní funkce*.

Jestliže je regresní funkce f lineární, což značí, že má tvar regresní přímky

$$f(x) = \beta_0 + \beta_1 x, \quad (3.3)$$

potom hovoříme o *jednoduché lineární regresi*, nemá-li regresní funkce lineární tvar, hovoříme o *jednoduché nelineární regresi*. Ve vzorci (3.3) jsou β_0, β_1 *parametry regresní funkce* neboli *regresní koeficienty*.

Mezi nejpoužívanější nelineární regresní funkce patří:

$$\text{regresní parabola:} \quad f(x) = \beta_0 + \beta_1 x^2, \quad (3.4)$$

$$\text{regresní hyperbola:} \quad f(x) = \beta_0 + \beta_1 \frac{1}{x}, \quad (3.5)$$

$$\text{regresní logaritmická funkce:} \quad f(x) = \beta_0 + \beta_1 \log x. \quad (3.6)$$

$$\text{regresní mocninná funkce:} \quad f(x) = \beta_0 x^{\beta_1}, \quad (3.7)$$

$$\text{regresní exponenciální funkce:} \quad f(x) = \beta_0 \beta_1^x. \quad (3.8)$$

Výše uvedené nelineární regresní funkce lze převést na lineární vhodnou transformací, jak uvidíme v následující kapitole.

Kromě výše uvedených příkladů nelineárních regresních funkcí existuje celá řada dalších významných nelineárních funkcí, např. Törnquistovy funkce, které nelze na lineární funkci jednoduše převést. Budeme se jimi zabývat v následující kapitole.

3.2 Jednoduchá regresní analýza

Představte si výběr párových hodnot $(y_1, x_1), (y_2, x_2), (y_3, x_3), \dots, (y_n, x_n)$, získaných (např. změřených) na statistických jednotkách základního souboru. Zde jsou y_i hodnotami závisle proměnné Y a x_i jsou hodnotami nezávisle proměnné X . Zmíněné párové hodnoty můžeme získat zejména dvojím způsobem:

- (A) Hodnoty nezávisle proměnné x_i jsme předem pevně zvolili a k nim jsme „změřili“ příslušné hodnoty y_i . V této situaci jsou hodnoty znaku X pevné (nenáhodné), zatímco hodnoty znaku Y považujeme za náhodné veličiny.
- (B) Párové hodnoty (y_i, x_i) „změříme“ na n náhodně zvolených jednotkách základního souboru. V této situaci jak hodnoty znaku X , tak hodnoty znaku Y považujeme za náhodné veličiny.

Výše uvedený datový soubor párových hodnot můžeme geometricky znázornit v rovině *bodovým grafem*, kde na vodorovnou osu „ x “ nanášíme hodnoty nezávisle proměnné a na svislou osu „ y “ příslušné hodnoty závisle proměnné. Výsledkem je geometrické znázornění n bodů v rovině, z jejichž vzájemné polohy můžeme soudit na regresní závislost znaku Y na X . Úkolem jednoduché lineární regrese je „proložit“ danými body přímkou (tj. nalézt lineární regresní funkci), která nejlépe charakterizuje polohu daných n bodů. Z předchozího odstavce víme, že tato regresní funkce má tvar $f(x) = \beta_0 + \beta_1 x$, kde β_0, β_1 jsou zatím neznámé hodnoty parametrů regresní přímky. Regresní model (3.2) má nyní tvar

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n. \quad (3.9)$$

Odhady b_0, b_1 těchto neznámých parametrů – *regresní koeficienty* získáme *metodou nejmenších čtverců*. Této metodě, která patří mezi nejdůležitější metody používané ve statistice, bude věnován následující odstavec.

3.3 Metoda nejmenších čtverců

Uvažujte data ve formě párových hodnot – bodů: $(y_1, x_1), (y_2, x_2), (y_3, x_3), \dots, (y_n, x_n)$. Úkolem jednoduché regrese je najít regresní funkci, která „nejlépe charakterizuje polohu“ daných n bodů. Nejprve budeme uvažovat obecný tvar regresní funkce $f(x; \beta_0, \beta_1)$ se dvěma parametry β_0, β_1 (nemusí to být nutně regresní přímka). Speciálními případy této regresní funkce je lineární funkce (3.3) a také nelineární funkce (3.4) – (3.8). Postup metody nejmenších čtverců bude vždy stejný, tj. nezávislý na konkrétním tvaru regresní funkce. Odhady b_0, b_1 neznámých parametrů β_0, β_1 získáme tak, že nalezneme hodnoty b_0, b_1 , pro něž nabývá své minimální hodnoty *reziduální součet čtverců* odchylek hodnot závisle proměnné y_i od teoretické hodnoty $Y_i = f(x_i; b_0, b_1)$, tj.

$$S_R = \sum_{i=1}^n (y_i - Y_i)^2 = \sum_{i=1}^n (y_i - f(x_i, b_0, b_1))^2. \quad (3.10)$$

Jak je známo z matematické analýzy, své minimum funkce S_R (zde je to funkce proměnných b_0, b_1) vždy nabývá pro ty hodnoty b_0, b_1 , pro něž se anulují její parciální derivace:

$$\frac{\partial S_R}{\partial b_0} = 0, \quad \frac{\partial S_R}{\partial b_1} = 0. \quad (3.11)$$

Vztahy (3.11) představují soustavu 2 rovnic o 2 neznámých b_0, b_1 , která se nazývá *soustava normálních rovnic*. Jejím řešením získáme hledané odhady regresních parametrů zvolené regresní funkce.

Vyřešíme nyní soustavu (3.11) pro speciální případ, který nás zejména zajímá, totiž pro lineární regresní funkci $f(x; \beta_0, \beta_1) = \beta_0 + \beta_1 x$. Dosadíme-li tuto funkci do vztahu (3.10), vypočteme příslušné parciální derivace, které položíme rovny 0, získáme konkrétní soustavu normálních rovnic

$$\begin{aligned} \sum_{i=1}^n y_i &= b_0 n + b_1 \sum_{i=1}^n x_i, \\ \sum_{i=1}^n x_i y_i &= b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2. \end{aligned} \quad (3.12)$$

Z těchto rovnic již snadno vypočteme hledané odhady b_0, b_1 takto:

$$b_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} \quad b_0 = \bar{y} - b_1 \bar{x}. \quad (3.13)$$

Z analytické geometrie si připomeňte, že regresní koeficient b_0 představuje průsečík regresní přímky s osou „y“, tedy hodnotu Y_0 pro $x = 0$, tento regresní koeficient se někdy nazývá *úrovňová konstanta*. Regresní koeficient b_1 vyjadřuje směrnici přímky, tedy sklon přímky k ose „x“, tj. změnu funkční hodnoty Y při změně nezávisle proměnné x o jednotku.

Pro jiné, než lineární tvary regresní funkce je postup metody nejmenších čtverců obdobný. Výsledkem je rovněž soustava 2 normálních rovnic, tyto rovnice však již nemusí být lineární, a proto soustavu již obvykle nelze snadno vyřešit. K řešení pak používáme *iterační numerické metody*, které zde nejsou předmětem našeho zájmu. V řešených úlohách jsou uvedeny způsoby nalezení odhadů regresních koeficientů metodou linearizace exponenciální a mocninné regresní funkce pomocí logaritmické transformace.

Na tomto místě bychom chtěli zvýraznit jeden důležitý fakt, který budeme v následujícím výkladu neustále využívat. Data pro regresní analýzu jsou výsledkem náhodného výběru, ať již jsme použili při jejich získání postup (A), nebo (B). Proto také výsledek jednoduché lineární regresní analýzy – odhady neznámých parametrů β_0, β_1 , tj. regresní koeficienty b_0, b_1 , budou náhodné veličiny. Při každém dalším náhodném výběru dat bude výsledek, tj. odhad b_0, b_1 , obecně jiný! Má proto význam hovořit dále o statistických charakteristikách těchto odhadnutých parametrů, jako např. střední hodnota, rozptyl.

3.4 Míra variability, koeficient determinace

Metoda nejmenších čtverců nás nyní přivedla k postupu, který jsme již použili v předchozí kapitole při analýze rozptylu. V ANOVA se jednalo o rozklad celkové variability znaku Y , vyjádřené jako celkový součet čtverců, na meziskupinový a vnitroskupinový (reziduální) součet čtverců. V analýze rozptylu jsme pracovali se znakem X , který měl kvalitativní povahu, a proto nebylo možné vyjádřit závislost regresním modelem. V regresní analýze má znak X – nezávisle proměnná – kvantitativní povahu, a proto je regresní model závislosti Y na X možný. Použijeme analogii s ANOVA v tom, že znak X zde bude nabývat hodnot x_1, x_2, \dots, x_n a i -tá skupina bude nyní charakterizována teoretickou hodnotou $Y_i = f(x_i; b_0, b_1)$, namísto skupinového průměru \bar{y}_i v ANOVA. Potom celkovou variabilitu vysvětlované proměnné charakterizuje *celkový součet čtverců*:

$$S_y = \sum_{i=1}^n (y_i - \bar{y})^2. \quad (3.14)$$

Část celkové variability vysvětlenou regresním modelem charakterizuje *teoretický součet čtverců*:

$$S_T = \sum_{i=1}^n (Y_i - \bar{y})^2, \quad (3.15)$$

nevysvětlenou část celkové variability představuje *reziduální součet čtverců* (3.10):

$$S_R = \sum_{i=1}^n (y_i - Y_i)^2, \quad (3.16)$$

kde $e_i = y_i - Y_i$ nazýváme *reziduum*.

Lze dokázat, že mezi jednotlivými součty čtverců platí základní vztah:

$$S_y = S_T + S_R. \quad (3.17)$$

Obdobně jako v analýze rozptylu jsme zavedli k vyjádření těsnosti vztahu Y a X poměr determinace, nyní zavedeme analogický pojem charakterizující přiléhavost dat k regresnímu modelu. Tímto pojmem je *koeficient determinace*, který definujeme vztahem

$$R^2 = 1 - \frac{S_R}{S_y}. \quad (3.18)$$

Ze vztahu (3.17) vyplývá, že koeficient determinace nabývá hodnoty z intervalu $[0,1]$ a určuje tu část celkové variability pozorovaných hodnot S_y , kterou lze vysvětlit daným regresním modelem. Jinak řečeno, po vynásobení koeficientu determinace hodnotou 100 obdržíme, kolik procent celkové variability je vysvětlitelných regresním modelem. Koeficient determinace je proto důležitou charakteristikou vhodnosti zvoleného regresního modelu.

Vztah (3.18) vzniká podílem náhodných veličin, a proto jakožto náhodná veličina je odhadem koeficientu determinace R^2 . Pro malé rozsahy výběru n je odhad (3.18) *vychýlený*, viz Ramík (2003), tj. nadhodnocuje přiléhavost k regresnímu modelu. Proto se používá *nevychýlený odhad* koeficientu determinace R_{adj}^2 (z angl. *adjusted*), který nazýváme *korigovaný (upravený) koeficient determinace*:

$$R_{adj}^2 = 1 - (1 - R^2) \frac{n-1}{n-2}. \quad (3.19)$$

Pro velké hodnoty n je však zlomek ve vzorci (3.19) blízký k jedné a korigovaný koeficient se blíží k „nekorigovanému“.

3.5 Klasický lineární model

Klasickým jednoduchým lineárním regresním modelem se nazývá regresní model (3.9):

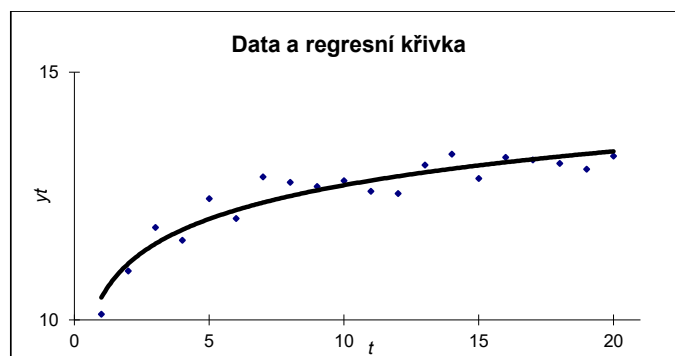
$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, i = 1, 2, \dots, n,$$

splňující následující podmínky:

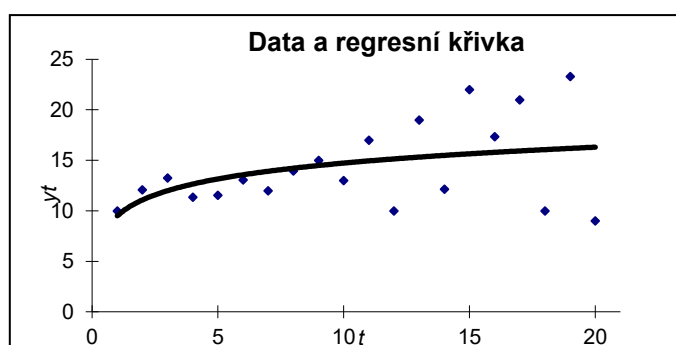
- (1) Hodnoty vysvětlující proměnné x_i se volí předem, viz (A) odstavec 3.2, nejsou to tedy náhodné veličiny.
- (2) Náhodné složky ε_i v modelu (3.9) mají *normální rozdělení* pravděpodobnosti se střední hodnotou 0 a (neznámým) rozptylem σ^2 . Konstantnost rozptylu nazýváme *homoskedasticita*.
- (3) Náhodné složky nejsou *korelované*, tj. $Cov(\varepsilon_i, \varepsilon_j) = 0$ pro každé $i \neq j$, $i, j = 1, 2, \dots, n$.

Podmínky (1) až (3) požadujeme tehdy, chceme-li zajistit splnění některých dalších vlastností: např. zjistit intervaly spolehlivosti koeficientů regresní funkce, interval spolehlivosti hodnoty regresní funkce, eventuálně chceme-li provádět testy hypotéz o některých prvcích regresního modelu. Těmito tématy se budeme zabývat v následujících odstavcích. Pokud totiž tyto podmínky splněny nejsou, nelze zajistit „spolehlivé předpovědi“.

V praxi jsou podmínky klasického modelu často splněny, nejsme-li si však jejich platností jisti, můžeme provést testy hypotéz jak o normalitě rozdělení náhodné složky (např. test dobré shody, viz např. Ramík (2003)), tak i testy o nekorelovanosti náhodných složek (např. t -test). Další testy uvedeme později v souvislosti s časovými řadami. Na Obrázku 12 je znázorněna situace, kdy podmínky klasického lineárního modelu jsou splněny, na Obrázku 13 je zachycena situace, kdy není splněna ani podmínka normality náhodných složek (na obrázku jsou všechny ε_i téměř stejné), ani podmínka nekorelovanosti (hodnoty y_i se nacházejí vedle sebe po jedné straně grafu regresní funkce).



Obrázek 12: Podmínky klasického modelu jsou splněny



Obrázek 13: Podmínky klasického modelu nejsou splněny

3.6 Diagnostická kontrola modelu

Kvalita každého sestaveného modelu se posuzuje pomocí diagnostických testů, kde jsou ověřovány vlastnosti náhodné složky, a to jsou heteroskedasticita, autokorelace a normalita reziduí. Pokud je model zvolen správně, pak má nesystematická (reziduální) složka modelu vlastnosti procesu bílého šumu.

3.6.1 HETEROSKEDASTICITA

Požadavkem na nesystematickou složku je její homoskedasticita čili konstantnost rozptylu. K posouzení homoskedasticity se využívá graf reziduí, v praxi pak také tzv. $ARCH(q)$ (AutoRegressive Conditional Heteroscedasticity) test, kde

- H_0 : reziduální složka vykazuje podmíněnou homoskedasticitu,
- H_1 : reziduální složka vykazuje podmíněnou heteroskedasticitu.

Pro posouzení efektu $ARCH(1)$ je konstruována tzv. umělá regrese, která je uvedena v rovnici (3.20), jak uvádí Arlt a Arltová (2007). Vysvětlovanou proměnnou je kvadrát reziduí a vysvětlující proměnnou kvadrát reziduí v prvním zpoždění.

$$\hat{a}_t^2 = \alpha_0 + \alpha_1 \hat{a}_{t-1}^2 + u_t \quad (3.20)$$

Metodou nejmenších čtverců jsou odhadnuty parametry a za předpokladu platnosti nulové hypotézy má statistika TR^2 rozdělení $\chi^2(1)$, T je počet měření, R^2 je index determinace. V případě vysokých hodnot statistiky TR^2 nulovou hypotézu zamítáme, a potvrzuje se, že n-systematická složka vykazuje podmíněnou heteroskedasticitu. Pro posouzení efektu $ARCH(q)$ je konstruována rovnice (3.21). V tomto případě za předpokladu platnosti nulové hypotézy má testové kritérium TR^2 rozdělení $\chi^2(q)$.

$$\hat{a}_t^2 = \alpha_0 + \alpha_1 \hat{a}_{t-1}^2 + \dots + \alpha_q \hat{a}_{t-q}^2 + u_t \quad (3.21)$$

Heteroskedasticita může být způsobena i přítomností odlehlých pozorování, a řešením může být jejich vypuštění z modelu. Heteroskedasticita je nežádoucí, protože způsobuje chybné testování parametrů v modelu.

3.6.2 AUTOKORELACE

Přítomnost autokorelace v modelu může znamenat, že nebyla odfiltrována veškerá systematická složka. V případě jednorozměrného procesu je pro zkoumání autokorelace používá výběrová autokorelační funkce (ACF) $\hat{r}_k = \frac{\sum_t \hat{a}_t \hat{a}_{t-k}}{\sum_t \hat{a}_t^2}$. V případě nekorelovanosti n-systematické složky leží hodnoty výběrové ACF uvnitř intervalu $(-2\sqrt{T}, 2\sqrt{T})$. Přítomnost autokorelace je ověřována na základě Portmanteau testu, kde

$$H_0: \rho_1 = \rho_2 = \dots = \rho_K = 0$$

$$H_1: \text{neplatí } H_0.$$

Je-li model správně konstruován, pak má statistika $Q = T \sum_{k=1}^K \hat{\rho}_k^2$ pro vysoká T a K přibližně rozdělení $\chi^2(K - p - q)$. Pro malé výběry se používá statistika označovaná jako modifikovaná Portmanteau statistika (Arlt a Arltová, 2007).

3.6.3 NORMALITA

Normalitu lze sledovat pomocí χ^2 testu dobré shody, nejčastěji je však využíván Jarque-Bera test, který je založený na testování šikmosti a špičatosti rozdělení. Hypotézy jsou

$$H_0: \text{normální rozdělení,}$$

H_1 : jiné než normální rozdělení.

Testové kritérium $JB = SK^2 + K^2$, kde SK je šikmost rozdělení a K je špičatost rozdělení, má za předpokladu platnosti nulové hypotézy rozdělení $\chi^2(2)$.



ŘEŠENÁ ÚLOHA 3.1

Společnost na výrobu bytového textilu zkoumala, jak souvisí zisk z prodeje s výdaji na reklamu. Tabulka 7 uvádí údaje obdržené v deseti náhodně vybraných firmách. Načrtněte bodový graf a určete typ regresní funkce popisující danou závislost. Stanovte koeficienty regresní funkce. Vypočítejte koeficient determinace a zhodnoťte těsnost závislosti vyjádřenou regresním modelem.

Tabulka 7: Zisk z prodeje a výdaje na reklamu

Pozorování	Výdaje na reklamu (tis. Kč)	Zisk z prodeje (10 tis. Kč)
1	6	5
2	8	8
3	9	9
4	9	12
5	12	21
6	15	25
7	16	32
8	20	36
9	22	51
10	23	59

Řešení („ruční“ výpočet):

Z grafu vidíte, že jde o přímou závislost, kterou je možné popsat regresní přímkou

$$Y = \beta_0 + \beta_1 x.$$

Máte za úkol stanovit hodnoty koeficientů b_0 , b_1 , neboli na základě dat odhadnout hodnoty parametrů β_1 , β_2 . Využijeme výsledků metody nejmenších čtverců, nebudete však dosazovat přímo do soustavy rovnic (3.12), ale použijete vztahy pro b_0 , b_1 , tj. (3.13), které je možné z dané soustavy vyjádřit, a to v numericky výhodném a snadno zapamatovatelném tvaru:

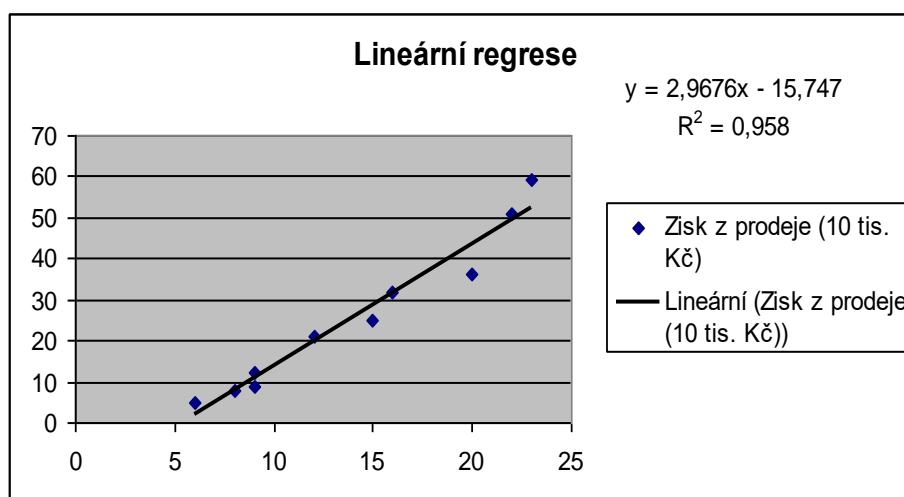
$$b_1 = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2} = \frac{462,1 - 14 \cdot 25,8}{230 - 14^2} = \frac{100,9}{34} = 2,97$$

$$b_0 = \bar{y} - b_1 \bar{x} = 25,8 - 2,97 \cdot 14 = -15,75.$$

Výpočty potřebných hodnot pomocí kalkulačky jsou uvedeny v Tabulce 8.

Tabulka 8: Výpočty

i	x_i	y_i	x_i^2	$x_i y_i$	Y_i	$(Y_i - \bar{y})^2$	$(y_i - \bar{y})^2$
1	6	5	36	0	2,04	565,21	432,64
2	8	8	64	64	7,98	318,22	316,84
3	9	9	81	81	10,95	221,15	282,24
4	9	12	81	108	10,95	221,15	190,44
5	12	21	144	252	19,86	35,62	23,04
6	15	25	225	375	28,77	8,61	0,64
7	16	32	256	512	31,74	34,84	38,44
8	20	36	400	720	43,62	315,88	104,04
9	22	51	484	1122	49,56	562,08	635,04
10	23	59	529	1357	52,53	711,60	1102,24
Součet	140	258	2300	4621	258	2994,3	3125,6
Průměr	14	25,8	230	462,1			



Obrázek 14: Graf regresní přímky

Hledaná regresní přímka má tvar: $Y = -15,75 + 2,97x$.

- a. K tomu, abychom vypočítali determinační koeficient, musíme znát hodnotu součtu S_T a součtu S_y . Tyto součty vypočítáme podle vztahů (3.14), (3.15). Pro výpočet teoretického součtu musíme pro každé x_i , $i = 1, \dots, 10$, znát teoretickou hodnotu Y_i .

$$Y_1 = -15,75 + 2,97 \cdot x_1 = -15,75 + 2,97 \cdot 6 = 2,04.$$

Tato hodnota udává, jaký by měl být zisk při výdajích $x = 6$. Protože však jde o stochastickou závislost mezi společenskými veličinami, může se tato hodnota lišit od skutečně zjištěné hodnoty $y = 5$. Všechny teoretické hodnoty Y_i i hodnoty součtů S_y a S_T jsou uvedeny v Tabulce 8. Koeficient determinace vypočítáme dosazením součtů S_y , S_T do vztahu (3.18).

$$R^2 = \frac{S_T}{S_y} = \frac{2994,3}{3125,6} = 0,958.$$

Tato hodnota znamená, že pomocí regresní přímky $Y = -15,78 + 2,97x$ je vysvětleno 95,8 % chování proměnné Y .

Řešení (výpočet v Excelu):

V Excelu využijeme graf funkce s funkcí Přidat spojnici trendu. Po volbě položky Vložit graf → XY bodový..., se otevře zadávací okno, kde zadáte:

Oblast dat: \$A\$1:\$B\$11

Sloupce: ✓ (zakliknout)

Potvrdíte OK

Obdržíte bodový graf, viz Obrázek 14 (ještě bez regresní přímky). Poklepem pravým tlačítkem myši na některý z bodů grafu obdržíte nabídku menu, kde zvolíte: Přidat spojnici trendu, Typ trendu regrese: zvolíte Lineární

Dále otevřete záložku *Možnosti*, kde zakliknete:

Zobrazit rovnici regrese (rovnice regresní přímky) a

Zobrazit hodnotu spolehlivosti R (hodnotu koeficientu determinace R^2).

Potvrdíte OK.

Obdržíte výsledek téměř takový, jaký je na Obrázku 14. K původním bodům se zobrazí regresní přímka, dále rovnice regresní přímky a hodnotu koeficientu determinace R^2 .



ŘEŠENÁ ÚLOHA 3.2

Společnost Air-Ostrava, zajišťující lety na trase Ostrava-Praha, sleduje při plánování letů také na hmotnost užitečného zatížení letadla, jehož významnou část tvoří pasažéři a jejich zavazadla. Zjistilo se, že hmotnost zavazadel cestujících souvisí s dobou, na kterou odcestovali. Výsledky průzkumu zachycuje Tabulka 9.

- a. Najděte rovnici regresní přímky popisující danou závislost.
- b. S jakou hmotností zavazadel lze počítat, bude-li na palubě 15 cestujících vracejících se za 2 dny, 7 cestujících vracejících se za 5 dnů, 5 cestujících vracejících se za 6 dnů a 1 cestující vracející se za 14 dní.

Tabulka 9: Výsledky průzkumu

Pozorování	Dny	Hmotnost
1	13	46
2	12	43
3	9	29
4	16	52
5	10	31
6	5	18
7	2	11
8	3	12
9	8	25
10	2	10
11	14	48
12	19	60
13	3	15
14	5	20
15	2	12

Řešení:

Prezentujeme zde pouze „ruční“ výpočet řešení (s kalkulačkou), řešení pomocí Excelu s využitím funkce Přidat spojnicí trendu v bodovém grafu ponecháváme na čtenáři.

a. K výpočtu regresních koeficientů b_0 , b_1 použijeme opět vztahů (3.13):

$$b_1 = \frac{\overline{x \cdot y} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2} = \frac{324,4 - 8,2 \cdot 28,8}{96,73 - 8,2^2} = 2,99, \quad b_0 = \bar{y} - b_1 \bar{x} = 28,8 - 2,99 \cdot 8,2 = 4,27$$

Regresní přímka má tedy tvar $Y = 4,27 + 2,99x$.

Tabulka 10: Výpočty

i	x_i	y_i	$x_i y_i$	x_i^2
1	13	46	598	169
2	12	43	516	144
3	9	29	261	81
4	16	52	832	256
5	10	31	310	100
6	5	18	90	25
7	2	11	22	4
8	3	12	36	9
9	8	25	200	64
10	2	10	20	4
11	14	48	672	196
12	19	60	1140	361
13	3	15	45	9
14	5	20	100	25
15	2	12	24	4
Součet	123	432	4866	1451
Průměr	8,2	28,8	324,4	96,73

b. Vypočítáme hodnotu Y pro $x = 2$: $Y(2) = 4,27 + 2,99 \cdot 2 = 10,25$,

$$x = 5: Y(5) = 4,27 + 2,99 \cdot 5 = 19,22,$$

$$x = 6: Y(6) = 4,27 + 2,99 \cdot 6 = 22,21,$$

$$x = 14: Y(14) = 4,27 + 2,99 \cdot 14 = 46,13.$$

Potom hmotnost zavazadel m , se kterou lze počítat, snadno zjistíte, uvážíte-li počty příslušných cestujících:

$$m = 15 \cdot Y(2) + 7 \cdot Y(5) + 5 \cdot Y(6) + 1 \cdot Y(14) = 153,75 + 134,54 + 111,05 + 46,13 = 445,47 \text{ kg.}$$



ŘEŠENÁ ÚLOHA 3.3 – SPOTŘEBNÍ FUNKCE KEYNESIÁNSKÉHO TYPU

Tato řešená úloha prezentuje ekonometrické modelování pro jednoduchou spotřební funkci keynesiánského typu pro české domácnosti v roce 2023. Predikujte vývoj spotřeby pro domácnost s měsíčním důchodem 55tis.Kč. Tato úloha bude řešena pomocí programu GRETl.

(Keynes: Lidé jsou v průměru ochotni zvyšovat svou spotřebu při rostoucích příjmech, ale ne v takové výši, jak rostou příjmy. Jedná se o přímou závislost reálné spotřeby především na reálném důchodu, přičemž spotřeba roste pomaleji než důchod.

Vymezení ekonomického modelu:

- Stanovení předmětu zkoumání – keynesiánská jednoduchá spotřební funkce
- Klasifikace ekonomických veličin – C_i (reálná spotřeba i -té domácnosti), Y_i (příjem domácnosti)
- Vymezení a verbální popis vazeb a vztahů mezi veličinami (přímá závislost reálné spotřeby především na reálném důchodu)
- Formulace výchozí základní hypotézy či tvrzení o chování ekonomických veličin (spotřeba roste pomaleji než důchod)

Vymezení matematického modelu:

Jednorovnicový lineární model: $C_i = \beta_1 + \beta_2 \cdot Y_i, \quad i = 1, 2, \dots, n,$

Kde β_1 je regresní parametr úrovně konstanty a β_2 je regresní parametr sklonu, který se očekává $0 < \beta_2 < 1$.

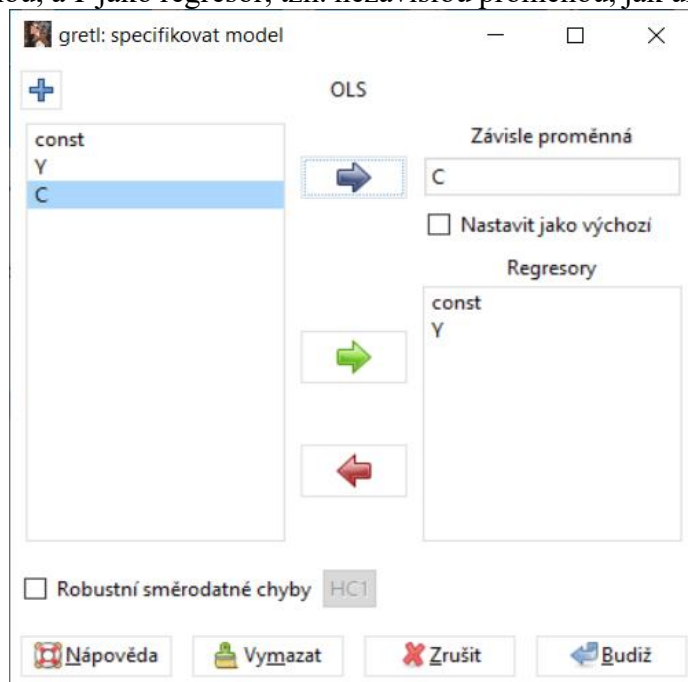
Formulace stochastického ekonometrického modelu:

Předpokládá zavedení náhodné složky u_i do rovnice: $C_i = \beta_1 + \beta_2 \cdot Y_i + u_i, \quad i = 1, 2, \dots, n,$ přičemž se předpokládá, že náhodná složka bude mít normální rozdělení s střední hodnotou nula, konstantním rozptylem, a nebude sériově závislá na svých zpožděných hodnotách.

	Y (příjem v Kč)	C (spotřeba v Kč)
1	46 995	33 907
2	46 644	34 717
3	46 295	34 363
4	46 847	35 365
5	48 695	35 799
6	49 887	36 722
7	50 801	37 993
8	52 631	40 389
9	54 906	40 750
10	58 410	41 826
11	62 493	43 496
12	66 326	45 079
13	67 850	46 323
14	65 124	46 636
15	66 550	46 662

Řešení:

Prezentujeme zde řešení pomocí programu GRETL. Nejprve do programu zadáme obě proměnné (C spotřeba domácností, Y příjem domácností). V hlavním menu vybereme MODEL→Ordinary Least Squares a objeví se následující dialogové okno, kde doplníme C jako závislou proměnnou, a Y jako regresor, tzn. nezávislou proměnnou, jak ukazuje Obrázek 15.



Obrázek 15: Dialogové okno – specifikace modelu

Po potvrzení dostáváme výsledek, který zachycuje Obrázek 16. Z toho vidíme, že regresní koeficient $b_2 = 0,568$ je statisticky významný na hladině významnosti 0,01 (p -hodnota je menší než 0,05), rovnice modelu je $C = 8539,27 + 0,568 \cdot Y$, koeficient determinace $R^2 = 0,96$. Predikce pro $Y = 55$ je $C = 8539,27 + 0,568 \cdot 55000 = 39\,779$ Kč.

Model 3: OLS, za použití pozorování 1-15
Závisle proměnná: C

	koeficient	směr. chyba	t-podíl	p-hodnota	
const	8539,27	1662,00	5,138	0,0002	***
Y	0,568289	0,0297110	19,13	6,66e-011	***

Střední hodnota závisle proměnné 40001,80
 Sm. odchylka závisle proměnné 4789,727
 Součet čtverců reziduí 11021087
 Sm. chyba regrese 920,7475
 Koeficient determinace 0,965686
 Adjustovaný koeficient determinace 0,963046
 F(1, 13) 365,8511
 P-hodnota (F) 6,66e-11
 Logaritmus věrohodnosti -122,5886
 Akaikovo kritérium 249,1772
 Schwarzovo kritérium 250,5933
 Hannan-Quinnovo kritérium 249,1621

zde je poznámka o zkratkách statistik modelu

Obrázek 16: Odhad koeficientů metodou nejmenších čtverců

Dále ověříme předpoklady modelu: heteroskedasticitu, normalitu a autokorelaci reziduí.

Pro testování heteroskedasticity vybereme ve výstupu modelu záložku TESTY→Heteroskedasticita→Whiteův test. A dostaneme výsledek, který zachycuje Obrázek 17.

Whiteův test heteroskedasticity
OLS, za použití pozorování 1-15
Závisle proměnná: uhat^2

	koeficient	směr. chyba	t-podíl	p-hodnota
const	-1,72312e+07	2,23342e+07	-0,7715	0,4553
Y	640,642	800,385	0,8004	0,4390
sq_Y	-0,00559329	0,00703687	-0,7949	0,4421

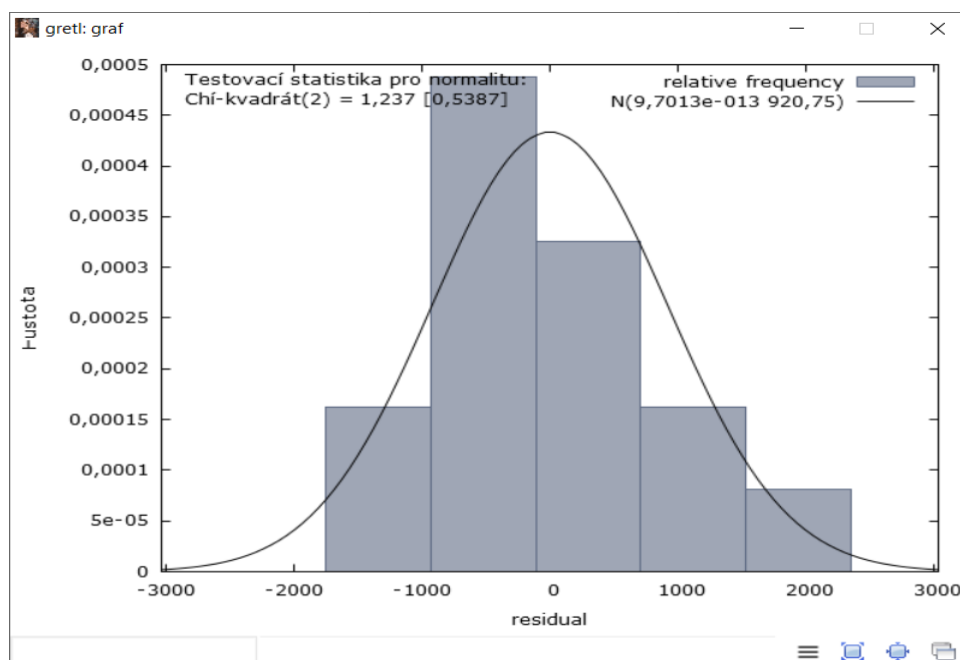
Neadjustovaný koeficient determinace = 0,051647

Testovací statistika: $TR^2 = 0,774702$,
 s p-hodnotou = $P(\text{Chi-kvadrát}(2) > 0,774702) = 0,678853$

Obrázek 17: Test heteroskedasticity

Vyhodnocení testu hetrosdedaticity provedeme na základě vypočtené p -hodnoty. Testuje se nulová hypotéza H_0 : homoskedasticita reziduí (tj. konstantní rozptyl reziduí), oproti alternativní hypotéze H_1 : heteroskedasticita reziduí. P -hodnota = 0,678 je větší než zvolené $\alpha = 0,05$; proto H_0 nelze zamítnout, nebylo tedy prokázáno, že by rezidua neměla konstantní rozptyl.

Pro testování normality vybereme ve výstupu modelu záložku TESTY→Normalita reziduí. A dostaneme výsledek, který zachycuje Obrázek 18. Vyhodnocení provedeného testu normality je pravděpodobně nejsnazší odvodit z průběhu grafu předpokládaného normálního rozdělení v porovnání se skutečným rozdělením reziduí a analýzou p -hodnoty Chí-kvadrát testu. Testuje se nulová hypotéza H_0 : Rezidua mají normální rozdělení, oproti H_1 : Rezidua nemají normální rozdělení. P -hodnota = 0,5387 je větší než zvolené $\alpha = 0,05$; proto H_0 nelze zamítnout, nebylo tedy prokázáno, že by rezidua neměla normální rozdělení.



Obrázek 18: Test normality

Pokud chceme pomocí programu GRETL testovat autokorelaci, musíme vstupní data uložit jako časovou řadu. Testuje se, zda je u_t závislé na u_{t-1} . Vybereme ve výstupu modelu záložku TESTY→Autokorelace. A dostaneme výsledek, který zachycuje Obrázek 19.

```

gretl: autokorelace
Breusch-Godfreyův test pro autokorelaci prvního řádu
OLS, za použití pozorování 1960-1974 (T = 15)
Závisle proměnná: uhat

-----
                koeficient    směr. chyba    t-podíl    p-hodnota
-----
const          148,475         1576,83         0,09416    0,9265
Y              -0,00252950         0,0281839      -0,08975    0,9300
uhat_1         0,417326                 0,264272         1,579      0,1403

Neadjustovaný koeficient determinace = 0,172055

Testovací statistika: LMF = 2,493721,
s p-hodnotou = P(F(1,12) > 2,49372) = 0,14

Alternativní statistika: TR^2 = 2,580829,
s p-hodnotou = P(Chi-kvadrát(1) > 2,58083) = 0,108

Ljung-Box Q' = 3,09597,
s p-hodnotou = P(Chi-kvadrát(1) > 3,09597) = 0,0785
    
```

Obrázek 19: Test autokorelace

Testuje se nulová hypotéza H_0 : Rezidua nejsou autokorelována, oproti H_1 : Rezidua jsou autokorelována. P -hodnota = 0,14 je větší než zvolené $\alpha = 0,05$; proto H_0 nelze zamítnout, nebylo tedy prokázáno, že by rezidua byla autokorelována.



SAMOSTATNÉ ÚKOLY

3.1 Personální ředitel firmy shromáždil údaje o věku (X) a době pracovní neschopnosti (Y) dvaceti náhodně vybraných stálých zaměstnanců. Zjištěné údaje jsou zaznamenány v tabulce.

X	Y	X	Y
20	4	58	20
35	14	46	13
35	15	43	16
34	10	33	10
32	10	29	10
28	9	36	11
25	12	48	14
46	15	55	15
38	15	36	14
50	16	19	6

Načrtněte bodový graf a najděte rovnici regresní funkce vyjadřující danou závislost. Zhodnoťte výstižnost (přiléhavost) regresní funkce vzhledem k datům.

3.2 Bylo sledováno, jak souvisí množství vadných výrobků (v % z vyrobených výrobků) s výkonem soustružníka (v % z předepsané normy). Bylo vybráno deset pracovníků, naměřené údaje jsou uvedeny v tabulce.

Výkon	56	68	72	85	92	102	107	111	123	142
Vadné výrobky	5,2	3,9	3,5	2,4	2,04	2	2,2	2,24	2,4	2,51

Stanovte regresní model a určete přiléhavost regresní přímky k datům.

3.3 Tabulka zachycuje stáří (v letech) osmi vybraných strojů v potravinářském závodě a týdenní náklady (v Kč) na provoz těchto strojů.

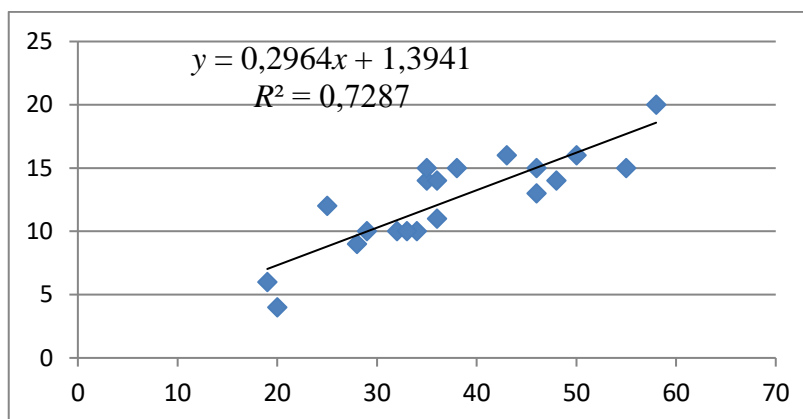
Stáří stroje	1	2	3	4	5	6	7	8
Náklady	44	52	61	80	94	108	111	116

- Odhadněte parametry lineární regresní funkce, která by měla vystihovat průběh závislosti nákladů na stáří.
- Určete koeficient determinace R^2 a interpretujte jej.
- Jaké týdenní náklady můžeme očekávat u stroje starého 4 roky?

ODPOVĚDI



3.1



3.2 $Y = -0,0285x + 5,56$; $R^2 = 0,53$.

3.3 a) $Y = 32,14 + 11,36x$

b) $R^2 = 0,97$ tzn. modelem je vysvětleno 97 % celkové variability.

c) $Y(4) = 32,14 + 11,36 \cdot 4 = 77,58$ Kč.



SHRNUTÍ KAPITOLY

Tato kapitola se zabývala jednoduchou regresní analýzou, byl zde formulován model jednoduché lineární regresní analýzy. Dále zde byla vysvětlena metoda nejmenších čtverců k nalezení „nejlepších“ hodnot regresních koeficientů v regresním modelu. Míra přiléhavosti dat k regresní křivce byla stanovena pomocí koeficientu determinace a jeho odmocniny – koeficientu korelace. Nakonec jste se seznámili s tzv. klasickým jednoduchým regresním modelem, který stanovuje 3 základní podmínky, kterým by měl vyhovovat regresní model vzhledem k existujícím datům.

4 REGRESNÍ ANALÝZA – JEDNOROZMĚRNÁ: INTERVALY SPOLEHLIVOSTI, TESTY HYPOTÉZ, NELINEÁRNÍ REGRESE

RYCHLÝ NÁHLED KAPITOLY



Tato kapitola vám rozšíří znalosti v jednorozměrné regresní analýze. Za předpokladů jednorozměrného klasického regresního modelu se budete zabývat stanovením intervalů spolehlivosti a dále testy hypotéz regresních koeficientů a testem nulovosti koeficientu determinace. Další odstavce se zabývají jednorozměrnou nelineární regresí. Nejprve budou vyšetřovány regresní funkce, které lze s pomocí vhodné transformace převést na funkce lineární dále parabolická regresní funkce, a nakonec nelineární regresní funkce tzv. Tornquiustova typu. Pro výpočet parametrů těchto funkcí, jež mají uplatnění především v marketingu, poznáte novou metodu tzv. metodu vybraných bodů.

CÍLE KAPITOLY



Po prostudování této kapitoly budete umět:

stanovit intervaly spolehlivosti pro regresní koeficienty,

testovat statistickou významnost regresních koeficientů,

testovat koeficient determinace a transformovat funkci na funkci lineární.

ČAS POTŘEBNÝ KE STUDIU



K prostudování této kapitoly budete potřebovat asi 90 minut.

KLÍČOVÁ SLOVA KAPITOLY



Intervaly spolehlivosti, testování regresních koeficientů, test koeficientu determinace.

4.1 Intervaly spolehlivosti

Jsou-li splněny předpoklady klasického lineárního modelu (3.9), tj. modelu

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n,$$

potom pro rozdělení odhadů regresních koeficientů b_0, b_1 jakožto náhodných veličin platí toto: Regresní koeficient b_j má normální rozdělení pravděpodobnosti se střední hodnotou β_j a rozptylem $\sigma^2 h_j$, kde $j = 0$ nebo 1 , čísla h_j jsou definována následujícími vztahy:

$$h_0 = \frac{\sum x_i^2}{n \sum x_i^2 - (\sum x_i)^2}, \quad (4.1)$$

$$h_1 = \frac{n}{n \sum x_i^2 - (\sum x_i)^2}. \quad (4.2)$$

V klasickém lineárním modelu předpokládáme, že náhodné složky mají konstantní rozptyl σ^2 , jeho hodnotu však neznáme. Neznámý rozptyl σ^2 můžeme nahradit jeho bodovým odhadem

$$s_R^2 = \frac{S_R}{n-2}, \quad (4.3)$$

který nazýváme *reziduální rozptyl*. Jak je vidět, v reziduálním rozptylu vystupuje v čitateli reziduální součet čtverců (3.16) dělený číslem $n-2$, což je *počet stupňů volnosti*, tj. rozsah dat n mínus počet regresních parametrů v modelu: 2. Odmocninu reziduálního rozptylu s_R nazýváme *směrodatná chyba*.

Oboustranný interval spolehlivosti pro regresní koeficient b_j , při zadaném koeficientu spolehlivosti $(1 - \alpha)$, je následující interval:

$$[b_j - t_{1-\alpha/2}(n-2) s_R \sqrt{h_j}, b_j + t_{1-\alpha/2}(n-2) s_R \sqrt{h_j}], \quad j = 0 \text{ nebo } 1. \quad (4.4)$$

Připomínáme, že zde $t_{1-\alpha/2}(n-2)$ je příslušný *kvantil* Studentova t -rozdělení, podrobnosti, viz Ramík (2003), h_j jsou dány vztahy (4.1), (4.2).

Bodový odhad regresních koeficientů b_j neříká nic o eventuální variabilitě tohoto koeficientu. Tuto informaci doplňuje směrodatná chyba (4.3) a zejména interval spolehlivosti (4.4), který informuje, v jakém rozmezí se regresní koeficient může pohybovat v rámci zadané spolehlivosti.

Odhadnutý lineární regresní model (3.1), který má tvar

$$y = b_0 + b_1 x + e, \quad (4.5)$$

resp. regresní funkce

$$Y = b_0 + b_1 x, \quad (4.6)$$

má praktický význam zejména při odhadu chování modelu v případě, že nezávisle proměnná nabývá nějakou v datech se nevyskytující hodnotu, označme ji např. x_0 . Model (4.5),

resp. regresní funkce (4.6), pak slouží k *předpovědi (predikci, prognóze, extrapolaci) hodnoty závisle proměnné*. Bodový odhad předpovědi získáme dosazením x_0 do (4.5), resp. (4.6), neboť predikovaná hodnota chyby (rezidua) e je 0, tedy

$$Y_0 = b_0 + b_1 x_0. \quad (4.7)$$

Informaci o tom, v jakém rozmezí se predikovaná hodnota závisle proměnné y může pohybovat, poskytne oboustranný interval spolehlivosti:

$$[Y_0 - t_{1-\alpha/2}(n-2) s_R \sqrt{H}, Y_0 + t_{1-\alpha/2}(n-2) s_R \sqrt{H}], \quad (4.8)$$

kde $H = 1 + \frac{1}{n} \left[1 + \frac{(n x_0 - \sum x_i)^2}{n \sum x_i^2 - (\sum x_i)^2} \right]$. Ostatní symboly v (4.8) mají stejný význam, jako

v intervalu (4.4).

4.2 Testy hypotéz

Metodou nejmenších čtverců lze zjistit, zda regresní koeficienty b_j jsou nenulová čísla, musíme mít však stále na paměti, že se jedná o realizace náhodných veličin, a tudíž má smysl testovat, zda naše původní parametry β_j jsou přesto nulové. Za předpokladů klasického lineárního modelu je možno testovat nulovou hypotézu:

$$H_0: \beta_j = 0, j = 0 \text{ nebo } 1 \quad (4.9)$$

proti oboustranné alternativní hypotéze

$$H_1: \beta_j \neq 0, j = 0 \text{ nebo } 1. \quad (4.10)$$

Při tomto testu použijeme testové kritérium

$$T = \frac{b_j}{\sqrt{\frac{S_R}{n-2} h_j}}, \quad (4.11)$$

kteří má při platnosti H_0 t -rozdělení s $n-2$ stupni volnosti, S_R je reziduální součet čtverců, h_j je dáno vztahy (4.1), (4.2), přičemž $j = 0$ nebo 1 .

Na hladině významnosti α (viz Ramík (2003)) je kritický obor vymezen nerovností

$$|T| > t_{1-\alpha/2}(n-2),$$

kde $t_{1-\alpha/2}(n-2)$ je příslušný kvantil Studentova t -rozdělení, který lze nalézt v tabulkách, nebo v Excelu pomocí funkce TINV.

Přijmete-li např. na dané hladině významnosti α nulovou hypotézu $H_0: \beta_1 = 0$, pak to znamená, že y *nezávisí* na x , jinak řečeno, pro libovolnou hodnotu nezávisle proměnné x nabývá závisle proměnná y neustále stejné hodnoty β_0 .

Vypočítaná hodnota koeficientu determinace je prakticky vždy kladná. Musíme však mít stále na paměti, že u hodnot vstupujících do výpočtu koeficientu determinace se jedná o realizace náhodných veličin, a tudíž má smysl testovat, zda teoretický koeficient determinace R^2 není přesto nulový. Za předpokladů klasického lineárního modelu je možno testovat nulovou hypotézu:

$$H_0: R^2 = 0,$$

proti oboustranné alternativní hypotéze

$$H_1: R^2 \neq 0.$$

Při tomto testu použijeme testové kritérium

$$T = \sqrt{\frac{R^2(n-2)}{1-R^2}}, \quad (4.11^*)$$

kteřé má při platnosti H_0 t -rozdělení, $n-2$ stupňů volnosti, R^2 je vypočítaný koeficient determinace.

Na hladině významnosti α je kritický obor vymezen nerovností $T > t_{1-\alpha}(n-2)$, (viz Rámík (2003)), kde $t_{1-\alpha}(n-2)$ je příslušný kvantil Studentova t -rozdělení, který lze nalézt v tabulkách, nebo v Excelu pomocí funkce TINV.

4.3 Nelineární regresní analýza

V tomto odstavci si povšimneme jednoduchého regresního modelu s nelineární regresní funkcí, který se však dá pouhou substitucí na lineární model převést. Konkrétně se jedná o dvě regresní funkce zmíněné již v kapitole 3:

$$\text{regresní mocninná funkce:} \quad f(x) = \beta_0 x^{\beta_1}, \quad (4.12)$$

$$\text{regresní exponenciální funkce:} \quad f(x) = \beta_0 \beta_1^x. \quad (4.13)$$

Regresní model s regresní funkcí (4.12) má tvar:

$$y = \beta_0 x^{\beta_1} + \varepsilon, \quad (4.14)$$

avšak namísto něj uvažujeme model, jež vznikne logaritmováním (4.12), kde položíme $y = f(x)$, tj. $\ln y = \ln \beta_0 + \beta_1 \ln x + \varepsilon'$, přitom \ln označuje přirozený logaritmus o základu $e = 2,718\dots$ Jestliže nyní položíte substituce

$$y' = \ln y, \quad x' = \ln x, \quad (4.15)$$

$$\beta'_0 = \ln \beta_0, \quad \beta'_1 = \beta_1, \quad (4.16)$$

pro transformaci (4.15) původních dat y_i, x_i obdržíte „čárkovaný“ jednoduchý lineární regresní model

$$y' = \beta'_0 + \beta'_1 x' + \varepsilon', \quad (4.17)$$

jehož parametry β'_0, β'_1 (regresní koeficienty) lze odhadnout metodou nejmenších čtverců aplikovanou na lineární model (4.17), a obdržíte tak jejich odhady b'_0, b'_1 . S použitím

vztahů (4.15) a (4.16) dostanete nazpět odhady b_0, b_1 původního nelineárního regresního modelu (4.12): $b_0 = e^{b'_0}, b_1 = e^{b'_1}$.

Analogickým postupem lze linearizovat jednoduchý nelineární regresní model s exponenciální regresní funkcí (4.13), která je v ekonomii známa jako *Cobb-Douglasova jednofaktorová produkční funkce*:

$$y = \beta_0 \beta_1^x + \varepsilon, \quad (4.18)$$

který substitucemi

$$y' = \ln y, x' = x, \quad (4.19)$$

$$\beta'_0 = \ln \beta_0, \beta'_1 = \ln \beta_1, \quad (4.20)$$

lze rovněž transformovat na „čárkovaný“ lineární model (4.17), jehož parametry β'_0, β'_1 odhadneme metodou nejmenších čtverců, a obdržíme tak jejich odhady b'_0, b'_1 . S použitím vztahů (4.20) vypočteme nazpět odhady b_0, b_1 původního nelineárního regresního modelu (4.18):

$$b_0 = e^{b'_0}, b_1 = e^{b'_1}. \quad (4.21)$$

Je však třeba upozornit, že na intervalové odhady, resp. testy hypotéz, regresních koeficientů b'_0, b'_1 lze použít postup z počátku této kapitoly pouze tehdy, když transformovaná, tj. „čárkovaná“ data y'_i, x'_i , splňují podmínky klasického regresního modelu z kapitoly 3. Meze intervalových odhadů, tedy krajní body intervalů spolehlivosti pak vypočítáme s použitím zpětných transformací (4.21).

Dalšími užitečnými nelineárními regresními funkcemi s uplatněním především v marketingu a výzkumu trhu (logistické funkce, Gompertzovy funkce, aj.) se budete zabývat v kapitole věnované analýze časových řad. Tam se budete zabývat i problémem výběru vhodného typu regresní funkce. V následujících odstavcích se ještě věnujeme známé parabolické regresní funkci a dále Törnquistovým funkcím, které nelze převést jednoduše na lineární tvar, jak tomu bylo v tomto odstavci.

4.4 Parabolická regrese

V kapitole 3.1. jsme označili parabolickou regresní funkci (3.4) za regresní funkci, kterou lze substitucí $x' = x^2$ převést na lineární tvar. V tomto případě se však jednalo pouze o speciální tvar paraboly (s vrcholem na ose y) se dvěma parametry. Obecný tvar paraboly však má parametry tři a vypadá takto:

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2. \quad (4.22)$$

Jednoduchý regresní model s parabolickou regresní funkcí pak má tvar

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon. \quad (4.23)$$

Máme-li tedy k dispozici data, tj. dvojice hodnot $(y_1, x_1), (y_2, x_2), (y_3, x_3), \dots, (y_n, x_n)$, pak lze odhady b_0, b_1, b_2 regresních parametrů $\beta_0, \beta_1, \beta_2$ získat metodou nejmenších čtverců, přičemž je zapotřebí řešit soustavu 3 normálních rovnic o 3 neznámých:

$$\begin{aligned}\sum y_i &= nb_0 + b_1 \sum x_i + b_2 \sum x_i^2, \\ \sum y_i x_i &= b_0 \sum x_i + b_1 \sum x_i^2 + b_2 \sum x_i^3, \\ \sum y_i x_i^2 &= b_0 \sum x_i^2 + b_1 \sum x_i^3 + b_2 \sum x_i^4.\end{aligned}\tag{4.24}$$

Uvědomte si, že neznámé jsou v této soustavě rovnic b_0, b_1, b_2 , zatímco y_i, x_i jsou známé hodnoty, které se dosadí do sum \sum v soustavě (4.24). Tuto soustavu 3 lineárních rovnic o 3 neznámých je snadné vyřešit např. známou Gaussovou eliminační metodou.

4.5 Törnqvistovy funkce

Zejména v marketingu se využívají Törnqvistovy regresní funkce (též Törnqvistovy křivky), což jsou regresní funkce s více parametry, které podle použití rozdělujeme na tři typy:

Törnqvistovy křivky I. typu vyjadřují závislosti *poptávky* po spotřebním zboží $f(x)$ na *výši příjmů* x ekonomických subjektů (např. rodin). Tyto křivky mají tvar:

$$f(x) = \frac{\beta_0 x}{x + \beta_1}.\tag{4.25}$$

Křivky tohoto typu se používají například při plánování a prognózování ve spotřebním průmyslu. Regresní funkce (4.25) slouží k modelování poptávky po zboží *nezbytného charakteru* (mléko, pečivo, obuv apod.).

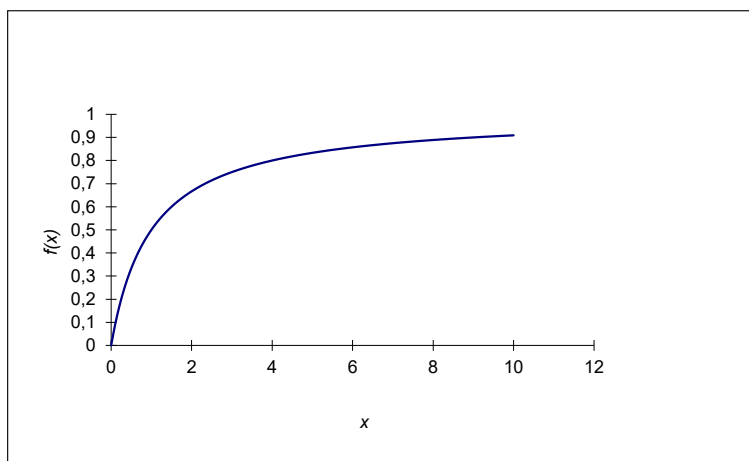
Při modelování poptávky po zboží *relativně nezbytného charakteru* (elektrospotřebiče, maso a uzeniny apod.) se používají *Törnqvistovy křivky II. typu*, které mají tvar:

$$f(x) = \frac{\beta_0(x - \beta_1)}{x + \beta_2}.\tag{4.26}$$

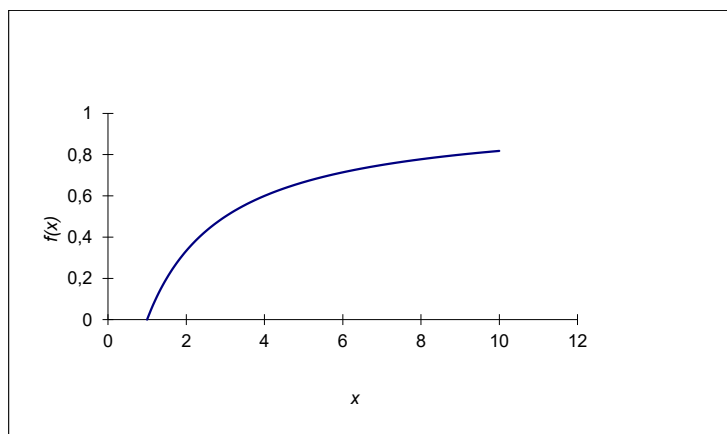
Törnqvistovy křivky III. typu se používají při modelování poptávky po zboží *zbytného charakteru* (auta, šperky, umělecká díla apod.). Tyto regresní funkce se třemi parametry mají tvar:

$$f(x) = \frac{\beta_0 x(x - \beta_1)}{x + \beta_2}.\tag{4.27}$$

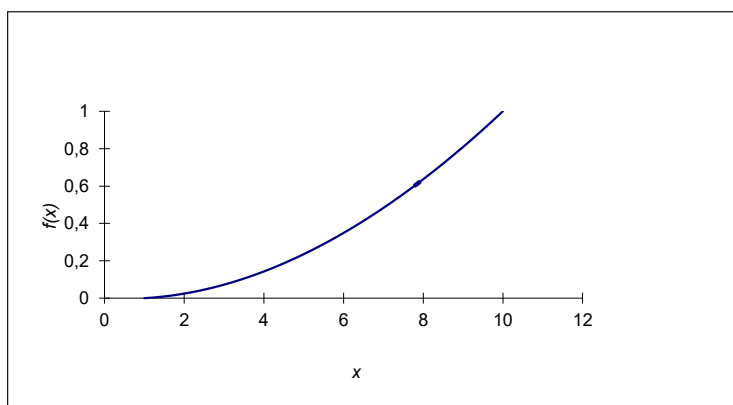
Odhady regresních parametrů funkcí (4.25) - (4.27) lze získat opět metodou nejmenších čtverců, avšak s použitím PC a Excelu, neboť soustava 3 normálních rovnic o 3 neznámých je nelineární, a proto se k řešení používají iterační numerické metody. Pro ruční výpočet můžeme alternativně využít i *metodu vybraných bodů*.



Obrázek 20: Törnqvistova křivka I. typu, $\beta_0 = \beta_1 = \beta_2 = 1$



Obrázek 21: Törnqvistova křivka II. typu, $\beta_0 = \beta_1 = \beta_2 = 1$



Obrázek 22: Törnqvistova křivka III. typu, $\beta_0 = \beta_1 = \beta_2 = 1$

4.6 Metoda vybraných bodů

Ukážeme si zde jinou metodu výpočtu neznámých parametrů, která sice nevede z teoretického pohledu k nejlepším odhadům, avšak její výhoda spočívá ve výpočetní nenáročnosti umožňující „ruční“ výpočet. Tato metoda se nazývá *metoda vybraných bodů* a spočívá v tom, že z daných údajů (Y_i, x_i) vybereme 3 charakteristické hodnoty, kterými necháme Törnquistovu křivku procházet, jinými slovy, položíme empirické hodnoty rovny hodnotám teoretickým. Jestliže charakteristické hodnoty poptávky Y_1, Y_2, Y_3 odpovídají hodnotám výše příjmů x_1, x_2, x_3 , pak ze vztahu (4.26) obdržíte soustavu 3 rovnic o 3 neznámých b_0, b_1, b_2 :

$$Y_1 = \frac{b_0(x_1 - b_1)}{x_1 + b_2}, \quad Y_2 = \frac{b_0(x_2 - b_1)}{x_2 + b_2}, \quad Y_3 = \frac{b_0(x_3 - b_1)}{x_3 + b_2}, \quad (4.28)$$

jejichž řešením např. postupným dosazováním získáme odhady neznámých parametrů b_0, b_1, b_2 .



ŘEŠENÁ ÚLOHA 4.1

Data v tabulce představují ceny brožovaných knih a k nim příslušné počty jejich stran.

- Určete lineární regresní model popisující závislost ceny knih na počtu stran.
- Určete interval, ve kterém bude s pravděpodobností 95 % ležet regresní koeficient b_1 .
- Na hladině významnosti 5 % testujte, zda je regresní koeficient b_1 statisticky významný.
- Vypočítejte koeficient determinace a na hladině významnosti 5 % testujte, zda je statisticky významný.
- V jakém rozmezí se bude pohybovat cena knihy s 250 stranami? Uvažujte hladinu významnosti 0,01.

Měření č.	1	2	3	4	5	6	7
Počet stran	20	35	48	50	130	200	86
Cena knihy	40	50	70	106	118	179	100

Řešení:

- Koeficienty regresní přímky $Y = b_0 + b_1x$ určíte pomocí vztahů (3.13):

$$b_1 = \frac{\bar{x} \cdot \bar{y} - \bar{\bar{x}} \cdot \bar{\bar{y}}}{\bar{x}^2 - \bar{\bar{x}}^2} = \frac{10135,71 - 81,29 \cdot 94,71}{10103,57 - 81,29^2} = \frac{2436,73}{3495,51} = 0,70$$

$$b_0 = \bar{y} - b_1 \cdot \bar{x} = 94,71 - 0,7 \cdot 81,29 = 37,81.$$

Hledaná regresní přímka má tvar $Y = 37,81 + 0,7x$.

- Úkolem je najít 95 % oboustranný interval spolehlivosti pro koeficient b_1 . Obecný tvar tohoto intervalu je následující (viz (4.4)):

$$[b_1 - t_{1-\alpha/2}(n-2) s_R \sqrt{h_1}, b_1 + t_{1-\alpha/2}(n-2) s_R \sqrt{h_1}],$$

kde s_R je odmocnina z reziduálního rozptylu $s_R^2 = \frac{S_R}{n-2}$, h_1 je definováno vztahem (4.2).

i	x_i	y_i	x_i^2	$x_i y_i$	Y_i	$(y_i - Y_i)^2$	$(y_i - \bar{y})^2$
1	20	40	400	800	51,81	139,48	2993,18
2	35	50	1225	1750	62,31	151,54	1998,98
3	48	70	2304	3360	71,41	1,99	610,58
4	50	106	2500	5300	72,81	1101,58	127,46
5	130	118	16900	15340	128,81	116,86	542,42
6	200	179	40000	35800	177,81	1,42	7104,80
7	86	100	7396	8600	98,01	3,96	27,98
Součet	569	663	70725	70950		1516,83	13405,43
Průměr	81,29	94,71	10103,57	10135,7			

Nejprve se vypočítá reziduální součet čtverců S_R (v tabulce výpočtů je to hodnota v předposledním sloupci dole):

$$S_R = \sum_{i=1}^7 (y_i - Y_i)^2 = 1516,83.$$

Teoretické hodnoty Y_i obdržíme postupným dosazováním hodnot x_i do rovnice regresní přímky. Hodnoty Y_i , jednotliví sčítanci i součet S_R jsou uvedeni v tabulce. Nyní můžeme vypočítat hodnotu reziduálního rozptylu s_R^2 .

$$s_R^2 = \frac{1516,83}{7-2} = 303,37.$$

Potom

$$s_R = \sqrt{s_R^2} = \sqrt{303,37} = 17,42.$$

Dále stanovíme hodnotu h_1 .

$$h_1 = \frac{n}{n \sum x_i^2 - (\sum x_i)^2} = \frac{7}{7 \cdot 70725 - 569^2} = \frac{7}{171314} = 0,00004.$$

V tabulkách Studentova rozdělení nalezneme $(1 - \alpha/2) = 97,5$ % kvantil t -rozdělení o $n-2 = 7 - 2 = 5$ stupních volnosti, tj. $t_{0,975}(5) = 2,57$.

Dosazením výše vypočítaných hodnot do vztahu pro interval spolehlivosti určíme jeho pravou a levou stranu:

$$L = 0,7 - 2,57 \cdot 17,42 \cdot \sqrt{0,00004} = 0,42.$$

$$P = 0,7 + 2,57 \cdot 17,42 \cdot \sqrt{0,00004} = 0,98.$$

Regresní koeficient b_1 bude s 95 % pravděpodobností ležet v intervalu $[0,42; 0,98]$.

c. Ačkoliv je hodnota koeficientu $b_1 = 0,7$, nesmíte zapomínat na to, že pracujete s náhodným výběrem a že teoretická hodnota parametru β_1 přesto může být nulová. Bude se proto testovat nulová hypotéza

$$H_0: \beta_1 = 0$$

proti oboustranné alternativní hypotéze

$$H_1: \beta_1 \neq 0.$$

K ověření nulové hypotézy vypočítáme hodnotu testového kritéria (4.11)

$$T = \frac{b_1}{\sqrt{\frac{S_R}{n-2} h_1}} = \frac{0,7}{\sqrt{\frac{1516,8}{7-2} \cdot 0,00004}} = \frac{0,7}{0,11} = 6,35.$$

V tabulkách t -rozdělení nalezneme $t_{0,975}(5) = 2,57$. Protože $6,35 > 2,57$, zamítáme nulovou hypotézu ve prospěch hypotézy alternativní, což znamená, že na zvolené hladině významnosti je parametr β_1 nenulový, a tedy statisticky významný.

d. Koeficient determinace R^2 vypočítáme podle vztahu

$$R^2 = 1 - \frac{S_R}{S_y} = 1 - \frac{1516,83}{13405,43} = 0,89.$$

Testové kritérium stanovíte podle vztahu (4.11*)

$$T = \sqrt{\frac{R^2(n-2)}{1-R^2}} = \sqrt{\frac{0,89 \cdot 5}{1-0,89}} = 6,35.$$

Protože $6,35 > 2,57$, zamítá se nulová hypotéza ve prospěch hypotézy alternativní, což znamená, že na zvolené hladině významnosti je koeficient determinace R^2 nenulový, a tedy statisticky významný.

e. Určete 99 % interval spolehlivosti pro predikovanou hodnotu Y , je-li $x_0 = 250$.

Podle (4.8) je tvar tohoto intervalu

$$[Y_0 - t_{1-\alpha/2}(n-2) s_R \sqrt{H}, Y_0 + t_{1-\alpha/2}(n-2) s_R \sqrt{H}],$$

kde

$$Y_0 = b_0 + b_1 x = 37,81 + 0,7 \cdot 250 = 212,81$$

$$t_{1-\alpha/2}(n-2) = 4,032$$

$$s_R = 17,42$$

$$H = 1 + \frac{1}{n} \left[1 + \frac{(n x_0 - \sum x_i)^2}{n \sum x_i^2 - (\sum x_i)^2} \right] = 1 + \frac{1}{7} \left[1 + \frac{(7 \cdot 250 - 569)^2}{7 \cdot 70725 - 569^2} \right] = 1 + \frac{1}{7} \left(1 + \frac{1394761}{171314} \right) = 1 + \frac{1}{7} \cdot 9,14 = 2,31.$$

Meze hledaného intervalu jsou:

$$L = 212,81 - 4,032 \cdot 17,42 \cdot \sqrt{2,31} = 106,06.$$

$$P = 212,81 + 4,032 \cdot 17,42 \cdot \sqrt{2,31} = 319,56.$$

Cena knihy se bude s 99 % pravděpodobností pohybovat v intervalu [106,06;319,56].

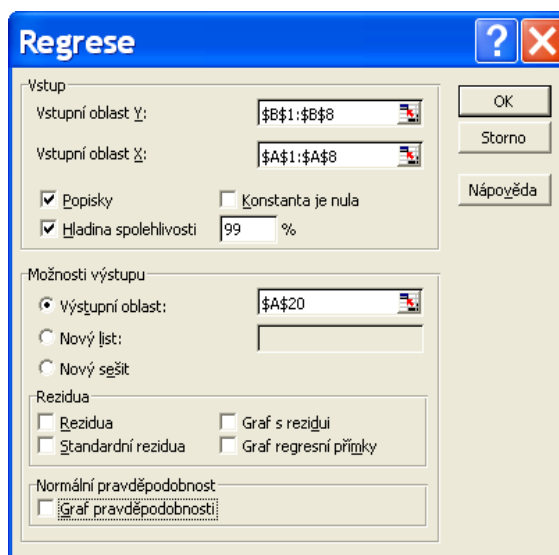
Nakonec si ukážeme řešení pomocí Excelu. Na tomto místě to bude další možnost řešení úlohy jednoduché (i vícenásobné) regrese s využitím menu:

Data → Analýza dat... → Regrese.

Data jsou uspořádána ve worksheetu ve 2 sloupcích:

	A	B	C
1	Počet stran	Cena knihy	
2	20	40	
3	35	50	
4	48	70	
5	50	106	
6	130	118	
7	200	179	
8	86	100	
9			

Otevře se okno regrese, které vyplníte takto:



VÝSLEDEK

Regresní statistika	
Násobné R	0,942
Hodnota spolehlivosti R	0,887
Nastavená hodnota spolehlivosti R	0,864
Chyba stř. hodnoty	17,416
Pozorování	7

ANOVA					
	Rozdíl	SS	MS	F	znamnost F
Regrese	1	11888,84	11888,84	39,19608	0,001525
Rezidua	5	1516,586	303,3172		
Celkem	6	13405,43			

	Koeficienty	ba stř. hod	t stat	Hodnota P	Dolní 95%	Horní 95%	Dolní 99,0%	Horní 99,0%
Hranice	38,059	11,19022	3,401	0,019	9,294	66,825	-7,061	83,180
Počet stran	0,697	0,111327	6,261	0,002	0,411	0,983	0,248	1,146

V první části výstupu jsou popisky s nepřesnými překlady do češtiny, správně má být:

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)

V této části výstupu je důležitá druhá hodnota – koeficient determinace $R^2 = 0,887$, který odpovídá ručně získanému výsledku z části d.

Druhá tabulka ve výstupu – ANOVA není v pravém slova smyslu metoda ANOVA, jak jsme se jí zabývali v kapitolách 1 a 2, jde tu o analogii využívající podobnosti vztahů (1.5) a (3.17). Analogicky jako v metodě ANOVA je zde výsledek F-testu statistické významnosti celého regresního modelu: Významnost $F = 0,001525$. Tato hodnota je menší než 0,05 a proto je celý regresní model statisticky významný.

Ve třetí – poslední tabulce jsou uvedeny relevantní informace k vypočítanému regresnímu modelu. Nejprve jsou uvedeny odhady regresních koeficientů:

Hranice = úrovněová konstanta = b_0

Počet stran = sklon regresní přímky = koeficient u nezávisle proměnné „počet stran“ = b_1

Ve sloupci Hodnota P jsou uvedeny p -hodnoty (signifikance) testů nulovosti příslušných regresních koeficientů:

Pro regresní koeficient b_0 je tato hodnota $0,019 < 0,05$; b_0 je statisticky významný tj. $\beta_0 \neq 0$.

Pro regresní koeficient b_1 je tato hodnota $0,002 < 0,05$; b_1 je statisticky významný tj. $\beta_1 \neq 0$.

Intervaly spolehlivosti regresních koeficientů jsou uvedeny ve sloupcích:

Dolní 95 %, Horní 95 %, resp. Dolní 99,0 %, Horní 99,0 %.

Konkrétně 95 % interval spolehlivosti koeficientu β_1 je [0,411; 0,983], což je stejný výsledek, jaký jsme obdrželi předtím ručním výpočtem.



ŘEŠENÁ ÚLOHA 4.2

Při sledování závislosti vlastních nákladů na skladování zahrnující i ztráty způsobené zastavením výroby z nedostatku součástek (Y) na velikosti dodávek (X) v 18 obuvnických závodech jsme obdrželi následující údaje - viz. tabulka.

a. Nalezněte regresní funkci popisující závislost Y na X a určete její rovnici.

b. Stanovte optimální velikost dodávky.

Podnik	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Dodávka	28	32	35	40	42	45	49	51	53	56	57	60	61	64	69	72	75	77
Náklady	62	59	58	53	50	46	44	42	40	41	38	35	36	36	38	40	42	46

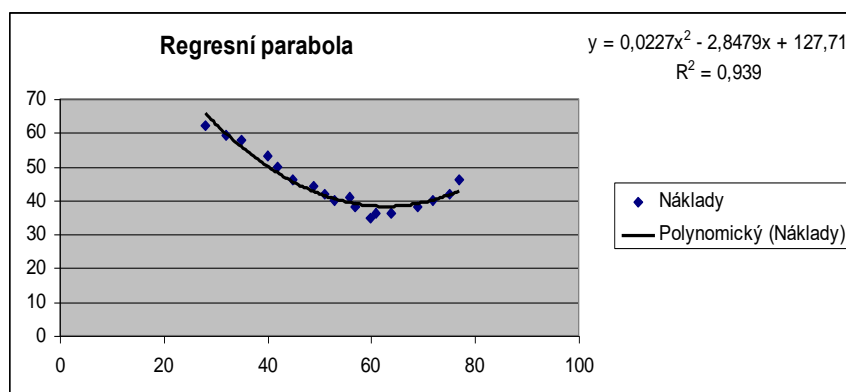
Řešení:

a. Jak z průběhu bodového diagramu, tak i rozboru empirických údajů plyne, že závislost mezi velikostí dodávek a náklady na skladování dobře vystihuje parabolická regresní funkce $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$.

Náklady na skladování mají zpočátku klesající tendenci, malá dodávka způsobuje vysoké náklady na převzetí připadající na jednu součástku a způsobuje výpadky ve výrobě. Tuto tendenci později vystřídá vzestup – příliš velká dodávka zvyšuje stav zásob, prodlužuje skladovací dobu a vyvolává nutnost úvěrového krytí – viz Obrázek 23.

Odhady hodnot parametrů parabolické regrese obdržíme řešením soustavy normálních rovnic

$$\begin{aligned} \sum_i y_i &= nb_0 + b_1 \sum_i x_i + b_2 \sum_i x_i^2 \\ \sum_i y_i x_i &= b_0 \sum_i x_i + b_1 \sum_i x_i^2 + b_2 \sum_i x_i^3 \\ \sum_i y_i x_i^2 &= b_0 \sum_i x_i^2 + b_1 \sum_i x_i^3 + b_2 \sum_i x_i^4 \end{aligned}$$



Obrázek 23: Parabolická regrese

Dosazením hodnot ze součtového řádku tabulky do těchto rovnic dostaneme:

$$\begin{aligned}806 &= 18b_0 + 966b_1 + 55534b_2 \\41618 &= 966b_0 + 55534b_1 + 3372084b_2 \\2330182 &= 55534b_0 + 3372084b_1 + 213664858b_2.\end{aligned}$$

Řešením této soustavy rovnic (např. Cramerovým pravidlem) získáme regresní koeficienty

$$b_0 = 127,71; b_1 = -2,8479; b_2 = 0,0227.$$

Hledaná parabola má tvar $Y = 127,71 - 2,8479x + 0,0227x^2$.

b. Optimální velikost objednávky zjistíme jako minimum funkce

$$Y = 127,71 - 2,8479x + 0,0227x^2$$

tak, že položíme její první derivaci rovnu nule, tj.

$$Y' = -2,8479 + 0,0454x = 0, \text{ tudíž } x = 62,7.$$

Optimální velikost dodávky je 62 nebo 63 kusů.

Nakonec provedeme výpočet pomocí Excelu s využitím funkce Přidat spojnicí trendu v bodovém grafu.

Po zobrazení dat pomocí grafu XY bodový poklepete pravým tlačítkem myši,

zvolíte položku Typ trendu a regrese: Polynomický (stupeň 2),

Dále otevřete záložku Možnosti, kde zakliknete:

Zobrazit rovnici regrese (rovnice regresní přímky) a současně zakliknete

Zobrazit hodnotu spolehlivosti R (hodnotu koeficientu determinace R^2).

Potvrdíte OK.

i	x_i	y_i	x_i^2	x_i^3	x_i^4	$x_i y_i$	$x_i^2 y_i$
1	28	62	784	21952	614656	1736	48608
2	32	59	1024	32768	1048576	1888	60416
3	35	58	1225	42875	1500625	2030	71050
4	40	53	1600	64000	2560000	2120	84800
5	42	50	1764	74088	3111696	2100	88200
6	45	46	2025	91125	4100625	2070	93150
7	49	44	2401	117649	5764801	2156	105644
8	51	42	2601	132651	6765201	2142	109242
9	53	40	2809	148877	7890481	2120	112360
10	56	41	3136	175616	9834496	2296	128576
11	57	38	3249	185193	10556001	2166	123462
12	60	35	3600	216000	12960000	2100	126000
13	61	36	3721	226981	13845841	2196	133956
14	64	36	4096	262144	16777216	2304	147456
15	69	38	4761	328509	22667121	2622	180918
16	72	40	5184	373248	26873856	2880	207360
17	75	42	5625	421875	31640625	3150	236250
18	77	46	5929	456533	35153041	3542	272734
Součet	966	806	55534	3372084	213664858	41618	2330182

Obdržíte výsledek téměř takový, jaký je na následujícím obrázku. K původním bodům se zobrazí regresní parabola, dále rovnice regresní paraboly a hodnotu koeficientu determinace R^2 . Výsledek je stejný, jako při ručním výpočtu, viz výše.



ŘEŠENÁ ÚLOHA 4.3

V jisté firmě zkoumali, jak závisí vlastní náklady na jednotku produkce (Y) na objemu produkce (X). Následující tabulka uvádí zjištěné údaje v různých obdobích.

- Najděte regresní hyperbolický model popisující danou závislost.
- Pomocí koeficientu determinace zhodnoťte přiléhavost regresní funkce k datům.

Řešení:

a. Dosadíte potřebné údaje do normálních rovnic, které získáte z hyperbolické regresní funkce (3.5) tak, že k nalezení minima součtu čtverců odchylek:

$$F(b_0, b_1) = \sum \left(y_i - \left(b_0 + b_1 \frac{1}{x_i} \right) \right)^2 \text{ se anulují parciální derivace, tj. } \frac{\partial F}{\partial b_0} = 0 \text{ a } \frac{\partial F}{\partial b_1} = 0.$$

Tím obdržíte následující normální rovnice:

$$\sum y_i = n \cdot b_0 + b_1 \sum \frac{1}{x_i}$$

$$\sum \frac{y_i}{x_i} = b_0 \sum \frac{1}{x_i} + b_1 \sum \frac{1}{x_i^2}$$

a obdržíme soustavu 2 rovnic o 2 neznámých

$$1574 = 13 \cdot b_0 + b_1 \cdot 7,13$$

$$1812,19 = b_0 \cdot 7,13 + b_1 \cdot 8,33.$$

Řešením této soustavy získáte odhady regresních parametrů:

$$b_0 = 3,32; b_1 = 214,71.$$

Hledaná regresní hyperbola má tvar: $Y = 3,32 + \frac{214,71}{x}$.

b. Nejdříve vypočítáte teoretické hodnoty Y_i postupným dosazením hodnot x_i do rovnice regresní hyperboly

$$Y_1 = 3,32 + \frac{214,71}{x_1} = 3,32 + \frac{214,71}{0,5} = 432,74.$$

Všechny hodnoty Y_i jsou uvedeny v tabulce, viz níže.

Dále vypočítáte součty S_T , S_y

$$\begin{aligned} S_T &= \sum_{i=1}^{13} (Y_i - \bar{y})^2 = (432,74 - 121,08)^2 + (310,05 - 121,08)^2 + \dots + (24,58 - 121,08)^2 = \\ &= 203722,02. \end{aligned}$$

$$S_y = \sum_{i=1}^{13} (y_i - \bar{y})^2 = (456 - 121,08)^2 + (297 - 121,08)^2 + \dots + (14 - 121,08)^2 = 2060,97.$$

i	x_i	y_i	$1/x_i$	$1/x_i^2$	y_i/x_i	Y_i	$(Y_i - \bar{y})^2$	$(y_i - \bar{y})^2$
1	0,5	456	2,00	4,00	912,00	432,74	97131,96	112171,41
2	0,7	297	1,43	2,04	424,29	310,05	35709,66	30947,85
3	0,9	206	1,11	1,23	228,89	241,89	14595,06	7211,41
4	1,4	165	0,71	0,51	117,86	156,68	1267,36	1928,97
5	1,9	118	0,53	0,28	62,11	116,33	22,56	9,49
6	3,2	79	0,31	0,10	24,69	70,42	2566,44	1770,73
7	4,2	57	0,24	0,06	13,57	54,44	4440,89	4106,25
8	4,8	54	0,21	0,04	11,25	48,05	5333,38	4499,73
9	6,9	40	0,14	0,02	5,80	34,44	7506,49	6573,97
10	7,9	35	0,13	0,02	4,43	30,50	8204,74	7409,77
11	8,8	30	0,11	0,01	3,41	27,72	8716,09	8295,57
12	9,2	23	0,11	0,01	2,50	26,66	8915,14	9619,69
13	10,1	14	0,10	0,01	1,39	24,58	9312,25	11466,13
Součet	60,5	1574	7,13	8,33	1812,19		203722,02	206010,97
Průměr	4,65	121,08	0,55	0,64	139,40			

Hodnoty jednotlivých sčítanců i součtů S_T , S_y jsou uvedeny v tabulce.

Koeficient determinace R^2 vypočítáte podle vztahu (3.18).

$$R^2 = \frac{S_T}{S_y} = \frac{203722,02}{206011,97} = 0,99.$$

Hodnota koeficientu determinace 0,99 je vysoká, což znamená, že daným regresním modelem s vysvětlující proměnnou „objem produkce“ je vysvětleno 99 % variability znaku Y . Pouze 1 % chování proměnné Y je ovlivněno jinými faktory.

ŘEŠENÁ ÚLOHA 4.4



Data v tabulce ukazují poptávku po určitém druhu zboží (v tis. ks) při různých cenách (v Kč). Popište závislost poptávky na ceně mocninnou regresní funkcí.

Pozorování	1	2	3	4	5	6
Cena	8,5	40	92	180	200	250
Poptávka	200	140	80	45	42	18

Řešení:

Úkolem je nalézt odhady parametrů β_1 , β_0 regresní funkce $Y = \beta_0 x^{\beta_1}$.

Použijete linearizující transformace, a to tak, že obě strany rovnice zlogaritmujete a použijete vhodnou substituci (viz odstavec 4.3), čímž získáte rovnici

$$Y' = \beta'_0 + \beta'_1 x',$$

kde $Y' = \ln Y$, $x' = \ln x$, $\beta'_0 = \ln \beta_0$, $\beta'_1 = \beta_1$, což je rovnice regresní přímky.

Regresní koeficienty b'_0 , b'_1 určíme pomocí známých vztahů takto:

$$b'_1 = \frac{\overline{x'y'} - \overline{x'} \cdot \overline{y'}}{\overline{x'^2} - \overline{x'}^2} = \frac{17,49 - 4,39 \cdot 4,18}{20,7 - 4,39 \cdot 4,39} = \frac{-0,86}{1,43} = -0,6$$

$$b'_0 = \overline{y'} - b'_1 \overline{x'} = 4,18 - (-0,6 \cdot 4,39) = 6,8.$$

<i>i</i>	<i>x</i>	<i>y</i>	<i>x'</i>	<i>y'</i>	<i>x'y'</i>	<i>x'²</i>
1	8,5	200	2,14	5,30	11,34	4,58
2	40	140	3,69	4,94	18,23	13,61
3	92	80	4,52	4,38	19,81	20,45
4	180	45	5,19	3,81	19,77	26,97
5	200	42	5,30	3,74	19,80	28,07
6	250	18	5,52	2,89	15,96	30,49
Průměr			4,39	4,18	17,49	20,70

Odhady b_0, b_1 původního modelu snadno vypočítáte zpětnou transformací

$$b'_1 = b_1, b_0 = e^{b'_0}.$$

Proto bude $b_1 = -0,6; b_0 = 897,85$.

Hledaná mocninná regresní funkce má tvar $Y = 897,85 \cdot x^{-0,6}$.

Nakonec provedeme výpočet pomocí Excelu s využitím funkce Přidat spojnicí trendu v bodovém grafu.

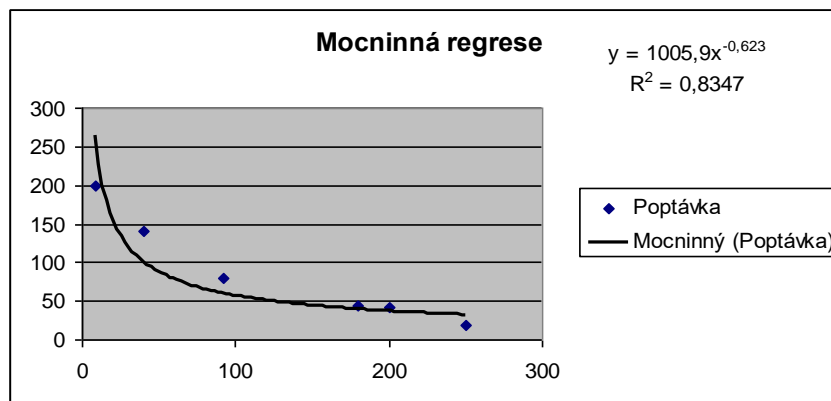
Po zobrazení dat pomocí grafu XY bodový poklepete pravým tlačítkem myši, zvolíte položku Typ trendu a regrese: Mocninný,

Dále otevřete záložku Možnosti, kde zakliknete:

Zobrazit rovnici regrese (rovnice regresní přímky) a současně zakliknete

Zobrazit hodnotu spolehlivosti R (hodnotu koeficientu determinace R^2).

Obdržíte výsledek, jaký je na Obrázku 24. K původním bodům se zobrazí regresní mocninná funkce, dále její rovnice a hodnotu koeficientu determinace R^2 . Výsledek je poněkud odlišný od výsledku, který jsme získali při ručním výpočtu, viz výše. Tato odlišnost je způsobena tím, že Excel počítá koeficienty přímo metodou nejmenších čtverců bez použití linearizace s logaritmickou transformací. Metoda použita Excelem je přesnější než metoda linearizace, a proto bychom ji dali při aplikaci přednost. Metoda linearizace je zase výpočetně jednodušší, je ji možno provést ručně, v době počítačů však tato výhoda ztrácí na významu.



Obrázek 24: Mocninná regrese

ŘEŠENÁ ÚLOHA 4.5

Tabulka uvádí stáří pletacích strojů (X) v letech a náklady na jejich údržbu (Y) v tis. Kč. Popište závislost Y na X exponenciální regresní funkcí.

Měření	1	2	3	4	5	6	7	8	9	10	11	12
Stáří	14	0,8	3	7,5	8,4	14,8	4,5	15,6	17,3	11,5	13,2	1,5
Náklady	47,5	8	10	17	22	76,4	12,5	76	94,5	25	30,6	12

Řešení:

Úkolem je nalézt odhady regresních parametrů exponenciální regresní funkce

$$y = \beta_0 \beta_1^x.$$

Pomocí logaritmické transformace převedeme tuto funkci na funkci lineární:

$$\ln y = \ln \beta_0 + x \ln \beta_1.$$

Použitím substitute

$$y' = \ln Y, x' = x, \beta'_0 = \ln \beta_0, \beta'_1 = \ln \beta_1$$

obdržíte regresní přímku $y' = \beta'_0 + \beta'_1 x'$.

Odhady parametrů β'_0, β'_1 této přímky určíme použitím známých vztahů

$$b'_1 = \frac{\overline{x'y'} - \overline{x'} \cdot \overline{y'}}{\overline{x'^2} - \overline{x'}^2} = \frac{34,8 - 9,34 \cdot 3,25}{118,59 - 9,34^2} = \frac{4,45}{31,35} = 0,14$$

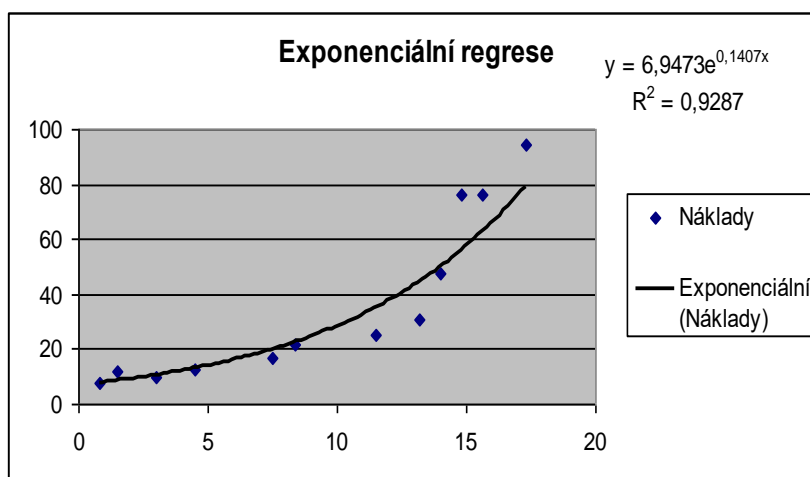
$$b'_0 = \overline{y'} - b'_1 \overline{x'} = 3,25 - (0,14 \cdot 9,34) = 1,94.$$

Regresní koeficienty původní funkce snadno vypočítáme zpětnou transformací:

$b_0 = e^{b'_0} = 6,96$; $b_1 = e^{b'_1} = 1,15$. Hledaná exponenciální regresní funkce má tvar:

$$y = 6,96 \cdot 1,15^x = 6,96 \cdot e^{0,14x}.$$

i	$x_i = x'_i$	y_i	y'_i	$x'_i y'_i$	x'^2
1	14	47,5	3,86	54,04	196,00
2	0,8	8	2,08	1,66	0,64
3	3	10	2,30	6,90	9,00
4	7,5	17	2,83	21,23	56,25
5	8,4	22	3,09	25,96	70,56
6	14,8	76,4	4,34	64,23	219,04
7	4,5	12,5	2,53	11,39	20,25
8	15,6	76	4,33	67,55	243,36
9	17,3	94,5	4,55	78,72	299,29
10	11,5	25	3,22	37,03	132,25
11	13,2	30,6	3,42	45,14	174,24
12	1,5	12	2,48	3,72	2,25
Průměr	9,34		3,25	34,80	118,59



Obrázek 25: Exponenciální regrese

Nakonec provedeme výpočet pomocí Excelu s využitím funkce Přidat spojnicí trendu v bodovém grafu. Po zobrazení dat pomocí grafu XY bodový poklepete pravým tlačítkem myši, zvolíte položku Typ trendu a regrese: Exponenciální,

Dále otevřete záložku Možnosti, kde zakliknete:

Zobrazit rovnici regrese (rovnice regresní přímky) a současně zakliknete

Zobrazit hodnotu spolehlivosti R (hodnotu koeficientu determinace R^2). Potvrdíte OK.

Obdržíte výsledek, jaký je na Obrázku 25. K původním bodům se zobrazí regresní exponenciální funkce, dále její rovnice a hodnotu koeficientu determinace R^2 . Výsledek je prakticky stejný jako výsledek, který jsme získali při ručním výpočtu, viz výše.



SAMOSTATNÉ ÚKOLY

4.1 Tabulka zachycuje stáří (v letech) osmi vybraných strojů v potravinářském závodě a týdenní náklady (v Kč) na provoz těchto strojů.

Stáří stroje	1	2	3	4	5	6	7	8
Náklady	44	52	61	80	94	108	111	116

- Odhadněte parametry regresní funkce $f(x) = \beta_0 + \beta_1 \ln x$, která by měla vystihovat průběh závislosti nákladů na stáří.
- Jaké týdenní náklady můžeme očekávat u stroje starého 4 roky?
- Určete koeficient determinace a interpretujte jej.

4.2 V tenisovém zápase má významný vliv na vítězství hráče úspěšnost jeho prvního podání. Data v tabulce představují počet úspěšných prvních podání (X) a počet vyhraných bodů při úspěšném prvním podání (Y) deseti vybraných hráčů z předních míst žebříčku ATP.

X	31	42	39	41	50	38	33	49	37	46
Y	22	31	29	26	33	26	23	30	29	31

Zvolte nejprve lineární a potom parabolický typ regresní funkce popisující závislost Y na X .

- Určete regresní parametry obou zvolených regresních funkcí.
- Stanovte 95 % interval spolehlivosti pro regresní koeficient b_1 u lineární regrese.
- Zhodnoťte výstižnost obou zvolených regresních funkcí. Která z nich lépe vystihuje data?

4.3 Ve výzkumu účinnosti léku se zkoumalo procento zlepšení účinnosti daného léku (Y) v závislosti na přidaném množství nové látky v mg (X).

Zvolte exponenciální typ regresní funkce popisující závislost Y na X .

X	1	1,5	2	2,5	3	3,5	4	4,5	5	5,5
Y	1,304	1,108	1,09	1,36	1,547	2,011	2,024	2,052	2,428	3,648

- Určete regresní parametry exponenciální regresní funkce.
- Sestrojte graf regresní funkce.
- Zhodnoťte výstižnost exponenciální regresní funkce.

ODPOVĚDI



4.1 a) $Y = 32,29 + 38,44 \cdot \ln x$ b) $Y(4) = 32,29 + 38,44 \cdot \ln 4 = 85,58Kč$

c) $R^2 = 0,92$

4.2 lineární regresní funkce

a) $Y = 7,95 + 0,49x$

b) $b_1 \in \langle 0,26; 0,73 \rangle$

c) $R^2 = 0,75$

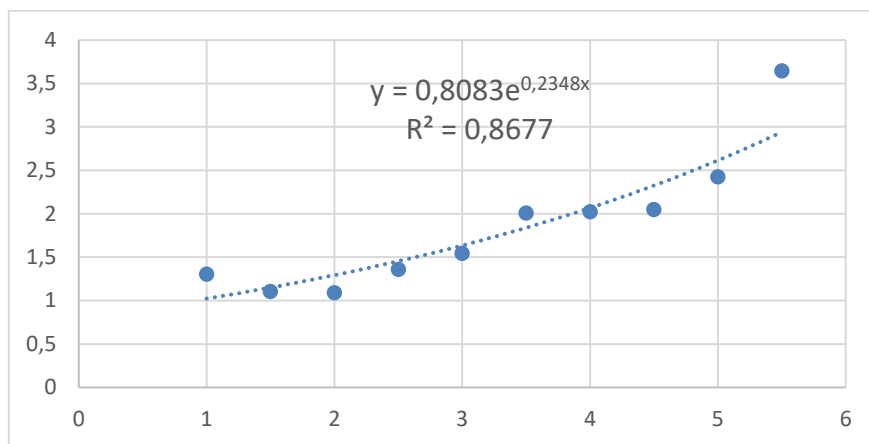
kvadratická regresní funkce

$Y = -25,94 + 2,19x - 0,02x^2$

$R^2 = 0,79$

Model lépe vystihuje kvadratická regresní funkce.

4.3 Rovnice exponenciální regresní funkce a koeficient determinace je v následujícím grafu.



SHRNUTÍ KAPITOLY

Tato kapitola přinesla rozšíření znalostí v jednorozměrné regresní analýze. Kapitola se zabývala stanovením intervalů spolehlivosti, testováním hypotéz regresních koeficientů a testem nulovosti koeficientu determinace. Dále zde byla představena jednorozměrná nelineární regrese. Byly zde vyšetřovány regresní funkce, které lze s pomocí vhodné transformace převést na funkce lineární, dále parabolická regresní funkce, a nakonec nelineární regresní funkce tzv. Tornqvistova typu. V této kapitole jste se seznámili s tzv. metodou vybraných bodů.

5 REGRESNÍ ANALÝZA – VÍCEROZMĚRNÁ

RYCHLÝ NÁHLED KAPITOLY



V této kapitole navážete na jednoduchou regresi vyšetřovanou v předchozí kapitole. Nyní budeme předpokládat, že vysvětlovaná proměnná závisí na několika (více než jedné) vysvětlujících proměnných. Vícenásobný lineární regresní model je zobecněním jednoduchého lineárního regresního modelu. Lineární regresní model bude rozšířen na vícenásobný regresní model lineární v parametrech, který předpokládá lineární vztah pouze v regresních koeficientech, nikoliv nutně v nezávisle proměnných. Odhady regresních koeficientů se stanoví opět metodou nejmenších čtverců, přitom lze využít maticové symboliky, která usnadňuje práci s vektory a maticemi. Podobně jako v případě jednoduché regrese budou formulovány předpoklady klasického regresního modelu, přičemž obdržíte analogické výsledky pro intervaly spolehlivosti regresních koeficientů a odpovídající testy hypotéz jako v případě jednoduché regrese. Nejprve budeme předpokládat, že vysvětlovaná proměnná Y závisí na několika vysvětlujících proměnných X_1, X_2, \dots, X_k .

CÍLE KAPITOLY



Po prostudování této kapitoly budete umět:

- napsat rovnici vícenásobného regresního modelu,
- vypočítat odhady regresních koeficientů pomocí maticové symboliky,
- vypočítat odhady regresních koeficientů v EXCELU a v GRETLU,
- interpretovat hodnotu koeficientu determinace a koeficientu korelace.

ČAS POTŘEBNÝ KE STUDIU



K prostudování této kapitoly budete potřebovat asi 120 minut.



KLÍČOVÁ SLOVA KAPITOLY

Vícenásobná regresní analýza, koeficient determinace, koeficient korelace.

5.1 Vícerozměrná regresní analýza

Na rozdíl od předchozích dvou kapitol, kde jsme předpokládali, že vysvětlovaná proměnná Y závisí na jediné vysvětlující proměnné X , budeme nyní předpokládat, že vysvětlujících proměnných je několik (tj. alespoň 2), řekněme k , kde $k \geq 2$, přitom k je celé číslo. Vysvětlující statistické znaky (proměnné) označíme X_1, X_2, \dots, X_k , i -tému pozorování (i -té realizaci) hodnot vysvětlujících znaků $x_{i1}, x_{i2}, \dots, x_{ik}$ odpovídá hodnota vysvětlovaného znaku y_i . Vícenásobný lineární regresní model je zobecněním jednoduchého lineárního regresního modelu (4.9) a má následující tvar:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad i = 1, 2, \dots, n. \quad (5.1)$$

Jak jste viděli v předchozí kapitole při aplikaci metody linearizace, bylo pro použití metody nejmenších čtverců podstatné, že regresní funkce byla lineární v parametrech β_i , nikoliv v proměnné x . Tohoto důležitého faktu využijeme nyní a formulujeme poněkud obecnější model, než (5.1), totiž vícenásobný regresní model *lineární v parametrech*. Ten vypadá takto

$$y_i = \beta_0 + \beta_1 f_1(x_{i1}, x_{i2}, \dots, x_{ik}) + \beta_2 f_2(x_{i1}, x_{i2}, \dots, x_{ik}) + \dots + \beta_k f_k(x_{i1}, x_{i2}, \dots, x_{ik}) + \varepsilon_i, \quad i = 1, 2, \dots, n. \quad (5.2)$$

kde $f_j(x_1, x_2, \dots, x_k)$, $j = 1, 2, \dots, k$, jsou funkce proměnných x_1, x_2, \dots, x_k , nezávislé na parametrech β_i .

5.2 Metoda nejmenších čtverců

Odhady regresních koeficientů b_0, b_1, \dots, b_k lze stanovit metodou nejmenších čtverců, která spočívá v minimalizaci součtu kvadrátů (tj. druhých mocnin) odchylek skutečných hodnot dat y_i od teoretických hodnot $Y_i = b_0 + b_1 f_1(x_{i1}, x_{i2}, \dots, x_{ik}) + \dots + b_k f_k(x_{i1}, x_{i2}, \dots, x_{ik})$. Podobně, jako u jednoduchého modelu, vypočteme odhady ze soustavy normálních rovnic:

$$\frac{\partial S_R}{\partial b_0} = 0, \quad \frac{\partial S_R}{\partial b_1} = 0, \quad \dots, \quad \frac{\partial S_R}{\partial b_k} = 0. \quad (5.3)$$

V (5.3) se jedná o parciální derivace funkce S_R podle proměnných b_i . Označení

$$F_{ij} = f_i(x_{j1}, x_{j2}, \dots, x_{jk}), \quad i = 1, 2, \dots, k, \quad j = 1, 2, \dots, n, \quad (5.4)$$

umožní využít maticovou symboliku. Soustavu rovnic (5.2) lze maticově zapsat takto:

$$\mathbf{y} = \mathbf{F}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (5.5)$$

kde matice:

$$\mathbf{F} = \begin{bmatrix} 1 & F_{11} & \cdots & F_{k1} \\ 1 & F_{12} & \cdots & F_{k2} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & F_{1n} & \cdots & F_{kn} \end{bmatrix} \text{ se nazývá matice regresorů,}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \text{ je vektor pozorování vysvětlované proměnné } Y,$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}, \text{ resp. } \mathbf{b} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_k \end{bmatrix}, \text{ je vektor regresních koeficientů, resp. vektor jejich od-}$$

hadů. Dále

$$\boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}, \text{ je vektor náhodných složek.}$$

Při výpočtu vektoru odhadů \mathbf{b} regresních koeficientů metodou nejmenších čtverců obdržíte soustavu normálních lineárních rovnic, které lze maticově vyjádřit. Pozor, používáte přitom pravidla pro sčítání a násobení matic, tzn. pravidlo „řádek krát sloupec“. Toho lze dosáhnout tak, že regresní rovnici $\mathbf{y} = \mathbf{F}\mathbf{b}$, vynásobíte zleva transponovanou maticí \mathbf{F}^T , takže obdržíte

$$\mathbf{F}^T\mathbf{y} = \mathbf{F}^T\mathbf{F}\mathbf{b}, \quad (5.6)$$

a za předpokladu, že matice $\mathbf{F}^T\mathbf{F}$ je regulární, a tedy existuje k ní matice inverzní $(\mathbf{F}^T\mathbf{F})^{-1}$, lze nalézt řešení soustavy, tj. vektor odhadů regresních koeficientů modelu (5.5), a to po vynásobení (5.6) zleva maticí $(\mathbf{F}^T\mathbf{F})^{-1}$, ve tvaru:

$$\mathbf{b} = (\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\mathbf{y}. \quad (5.7)$$

Ve speciálním případě jednoduché lineární regrese je $k = 1$, pak matice regresorů a další prvky z (5.6) mají tvar:

$$\mathbf{F} = \begin{bmatrix} 1 & x_{11} \\ 1 & x_{21} \\ \vdots & \vdots \\ 1 & x_{n1} \end{bmatrix}, \mathbf{F}^T\mathbf{F} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}, \mathbf{F}^T\mathbf{y} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix},$$

a soustava normálních rovnic (5.6) je následující:

$$\begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}, \quad (5.8)$$

což je tvar ekvivalentní rovnicím (3.12), (3.13).

5.3 Náhodný vektor a jeho charakteristiky

Nyní ještě rozšíříme pojmy střední hodnoty a rozptylu používané doposud pro náhodnou veličinu (skalár), a to pro náhodný vektor:

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}, \quad (5.9)$$

kde složky X_i jsou náhodné veličiny. *Střední hodnota* $\mathbf{E}(\mathbf{X})$ *vektorové náhodné veličiny* \mathbf{X} je vektor středních hodnot jednotlivých složek, tj.:

$$\mathbf{E}(\mathbf{X}) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_n) \end{bmatrix}. \quad (5.10)$$

Rozptyl (variance) $\mathbf{Var}(\mathbf{X})$ *vektorové náhodné veličiny* \mathbf{X} je matice:

$$\mathbf{Var}(\mathbf{X}) = \mathbf{E}((\mathbf{X} - \mathbf{E}(\mathbf{X}))^T (\mathbf{X} - \mathbf{E}(\mathbf{X}))). \quad (5.11)$$

Rozptyl náhodného vektoru (5.11) je čtvercová matice typu $(n \times n)$.

5.4 Klasický lineární model

O *klasickém (vícerozměrném) lineárním regresním modelu* hovoříme tehdy, když matice regresorů má nejjednodušší tvar, tj. když je matice tvořena danými hodnotami pozorování vysvětlujících proměnných:

$$F_{ij} = x_{ij}, \quad i = 1, 2, \dots, k, \quad j = 1, 2, \dots, n. \quad (5.12)$$

V tom případě má matice regresorů tvar:

$$\mathbf{F} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{k1} \\ 1 & x_{12} & \cdots & x_{k2} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & \cdots & x_{kn} \end{bmatrix}. \quad (5.13)$$

U klasického lineárního modelu požadujeme splnění podmínek 1. až 3. z minulé kapitoly, přitom u těchto podmínek nebylo důležité, zda jde o jednoduchý nebo vícerozměrný regresní model:

1. Hodnoty vysvětlujících proměnných X_1, X_2, \dots, X_k , tvořící matici regresorů \mathbf{F} podle (5.13) se volí předem, nejsou to tedy náhodné veličiny.

2. Reziduum $\boldsymbol{\varepsilon}$ v modelu (3.5) má *normální rozdělení* pravděpodobnosti s nulovou střední hodnotou a (neznámým) rozptylem σ^2 , tj.:

$$\mathbf{E}(\boldsymbol{\varepsilon}) = \mathbf{0}, \quad (5.14)$$

$$\mathbf{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}, \quad (5.15)$$

kde symbol \mathbf{I} označuje jednotkovou matici.

Vztah (5.15) zahrnuje zároveň podmínku 3. z klasického lineárního modelu, viz kapitola 3.5, neboť na diagonále matice $\mathbf{Var}(\boldsymbol{\varepsilon})$ jsou rozptyly σ^2 jednotlivých složek náhodného vektoru $\boldsymbol{\varepsilon}$ a mimo diagonálu vystupují nulové kovariance těchto složek. V tom případě hovoříme o *homoskedasticitě*. V opačném případě hovoříme o přítomnosti *heteroskedasticity*.

3. Vysvětlující proměnné X_1, X_2, \dots, X_k , nejsou kolineární, tj. sloupcové vektory matice regresorů (5.13) jsou nekorelované. V opačném případě hovoříme o přítomnosti multikolinearity.

5.5 Míry variability a koeficient determinace

Podobně jako u jednoduché regrese, zajímáme se nyní o *celkovou variabilitu* vysvětlované proměnné, kterou charakterizuje celkový součet čtverců:

$$S_y = \sum_{i=1}^n (y_i - \bar{y})^2. \quad (5.16)$$

Část celkové variability vysvětlenou regresním modelem charakterizuje *teoretický součet čtverců*:

$$S_T = \sum_{i=1}^n (Y_i - \bar{y})^2, \quad (5.17)$$

kde $Y_i = b_0 + b_1 f_1(x_{i1}, x_{i2}, \dots, x_{ik}) + \dots + b_k f_k(x_{i1}, x_{i2}, \dots, x_{ik})$, b_i jsou odhady regresních parametrů získané MNČ. Nevysvětlenou část celkové variability představuje reziduální součet čtverců:

$$S_R = \sum_{i=1}^n (y_i - Y_i)^2, \quad (5.18)$$

kde $e_i = y_i - Y_i$ je *reziduum*, tj. odhad náhodné složky ε_i .

Mezi jednotlivými součty čtverců platí základní vztah:

$$S_y = S_T + S_R \quad (5.19)$$

Obdobně, jako v případě jednoduché regrese, zavedeme analogický pojem, charakterizující přiléhavost dat k regresnímu modelu, *koeficient determinace*, který definujeme vztahem:

$$R^2 = \frac{S_T}{S_y} = 1 - \frac{S_R}{S_y}. \quad (5.20)$$

Koeficient determinace nabývá hodnoty z intervalu $[0, 1]$ a určuje tu část celkové variability pozorovaných hodnot y_i , kterou lze vysvětlit daným regresním modelem. Jinak řečeno, po vynásobení koeficientu determinace stem obdržíme, kolik procent celkové variability je vysvětlitelných regresním modelem.

Nevychýlený odhad koeficientu determinace R_{adj}^2 , který nazýváme *korigovaný (upravený) koeficient determinace*, definujeme takto:

$$R_{adj}^2 = 1 - \left(1 - R^2\right) \frac{n-1}{n-p}, \quad (5.21)$$

kde $p = k+1$ označuje počet parametrů v regresním modelu (5.2).

5.6 Intervaly spolehlivosti a testy hypotéz

Tento odstavec je přirozeným rozšířením kapitoly 4 pro jednoduchý klasický lineární model, tj. model (3.9) se dvěma parametry β_0, β_1 . Nyní máme analogický model, avšak s $k+1$ parametry $\beta_0, \beta_1, \dots, \beta_k$.

Jsou-li splněny předpoklady klasického lineárního modelu (5.5), tj. modelu:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad i = 1, 2, \dots, n, \quad (5.22)$$

potom pro rozdělení odhadů regresních koeficientů b_0, b_1, \dots, b_k , jakožto náhodných veličin, platí toto: Regresní koeficient b_j má normální rozdělení pravděpodobnosti se střední hodnotou β_j a rozptylem $\sigma^2 h_{jj}$, kde $j = 0, 1, \dots, k$, čísla h_{jj} jsou diagonálními prvky matice:

$$\mathbf{H} = (\mathbf{F}^T \mathbf{F})^{-1}, \quad (5.23)$$

kde matice \mathbf{F} je definována vztahem (5.13).

V klasickém lineárním modelu předpokládáme, že reziduální složky mají konstantní rozptyl σ^2 , jeho hodnotu však zpravidla neznáme. Neznámý rozptyl σ^2 můžeme nahradit jeho bodovým odhadem:

$$s_R^2 = \frac{S_R}{n-p}, \quad (5.24)$$

který nazýváme v souladu s (5.22) *reziduální rozptyl*. V reziduálním rozptylu vystupuje v čitateli reziduální součet čtverců (5.18) dělený číslem $n-p$, což je *počet stupňů volnosti*, tj. rozsah dat n mínus počet regresních koeficientů v modelu: $p = k + 1$. Odmocninu reziduálního rozptylu s_R nazýváme *směrodatná chyba*.

Oboustranný interval spolehlivosti pro regresní koeficient b_j , při zadaném koeficientu spolehlivosti $(1 - \alpha)$, je následující interval:

$$\left[b_j - t_{1-\alpha/2}(n-p) \sqrt{\frac{S_R h_{jj}}{n-p}}, b_j + t_{1-\alpha/2}(n-p) \sqrt{\frac{S_R h_{jj}}{n-p}} \right], \quad j = 0, 1, \dots, k. \quad (5.25)$$

Zde $t_{1-\alpha/2}(n-p)$ je příslušný kvantil Studentova t -rozdělení, h_{jj} diagonální prvky matice (5.23). Interval (4.23) je speciálním případem intervalu (5.25) v případě $k = 1$.

Bodový odhad regresních koeficientů b_j , vypočtený metodou nejmenších čtverců, doplňuje interval spolehlivosti (5.25), který informuje, v jakém rozmezí se regresní koeficient může pohybovat v rámci zadané spolehlivosti v případě jiného náhodného výběru dat (ze stejného základního souboru). Odhadnutý lineární regresní model (3.9), který má tvar:

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_kx_k + e, \quad (5.26)$$

kde e je reziduum, tj. odhad náhodné složky ε , resp. regresní funkce:

$$Y = b_0 + b_1x_1 + b_2x_2 + \dots + b_kx_k, \quad (5.27)$$

má praktický význam zejména při odhadu chování modelu pro nezávisle proměnné nevyskytující se v datech, např. hodnoty $x_{01}, x_{02}, \dots, x_{0k}$. Model (5.26), resp. regresní funkce (5.27), pak slouží k predikci hodnoty závisle proměnné. Bodový odhad předpovědi získáme dosazením $\mathbf{x}_0 = (x_{01}, x_{02}, \dots, x_{0k})'$ do (5.27):

$$Y_0 = b_0 + b_1x_{01} + b_2x_{02} + \dots + b_kx_{0k}. \quad (5.28)$$

Informaci o tom, v jakém rozmezí se predikovaná hodnota vysvětlované proměnné může pohybovat, poskytuje oboustranný interval spolehlivosti:

$$[Y_0 - t_{1-\alpha/2}(n-p) s_R \sqrt{1 + \mathbf{x}_0^T \mathbf{H} \mathbf{x}_0}, Y_0 + t_{1-\alpha/2}(n-p) s_R \sqrt{1 + \mathbf{x}_0^T \mathbf{H} \mathbf{x}_0}], \quad (5.29)$$

kde $\mathbf{H} = (\mathbf{F}^T \mathbf{F})^{-1}$ a matice \mathbf{F} je definována vztahem (5.13). Ostatní symboly v (5.29) mají stejný význam, jako v intervalu spolehlivosti (5.25).

5.7 Individuální T-testy o hodnotách regresních koeficientů

Zjistíme-li metodou nejmenších čtverců, že regresní koeficienty b_j jsou nějaká nenulová čísla, musíme mít stále na paměti, že se jedná o realizace náhodných veličin, a tudíž má smysl testovat, zda naše původní parametry β_j nemohou být přesto nulové. Za předpokladů klasického lineárního modelu je možno pro $j = 0, 1, \dots, k$ testujeme nulovou hypotézu:

$$H_0: \beta_j = 0, \quad (5.30)$$

proti oboustranné alternativní hypotéze:

$$H_1: \beta_j \neq 0. \quad (5.31)$$

Při tomto testu použijeme testové kritérium:

$$t = \frac{b_j}{\sqrt{\frac{S_R}{n-p} h_{jj}}}, \quad (5.32)$$

kteří má při platnosti H_0 t -rozdělení s $n-p$ stupni volnosti, S_R je reziduální součet čtverců, h_{jj} jsou diagonální prvky matice \mathbf{H} z (5.23), přičemž $j = 0, 1, \dots, k$, $p = k + 1$.

Na hladině významnosti α je kritický obor vymezen nerovností:

$$|t| > t_{1-\alpha/2}(n-p),$$

kde $t_{1-\alpha/2}(n-p)$ je příslušný kvantil Studentova t -rozdělení, viz funkce v Excelu TINV. Nemůžeme-li např. na dané hladině významnosti α zamítnout nulovou hypotézu $H_0: \beta_j = 0$, pak to znamená, že y nezávisí na x_j , jinak řečeno, pro libovolnou hodnotu vysvětlující proměnné x_j nabývá vysvětlovaná proměnná y stále stejné hodnoty.

5.8 F-test hypotézy o hodnotách regresních koeficientů

V minulém odstavci jste individuálními t -testy zjišťovali vliv jednotlivých vysvětlujících proměnných na vysvětlovanou proměnnou. V tomto odstavci se budeme zabývat testem, který najednou odhalí, zda vůbec existuje nějaká vysvětlující proměnná, která má na vysvětlovanou proměnnou nějaký vliv. Testuje se nulová hypotéza:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0, \quad (5.33)$$

proti alternativní hypotéze, že pro alespoň jeden regresní koeficient platí $\beta_j \neq 0$.

$$\text{Testové kritérium: } T = \frac{\frac{S_T}{p-1}}{\frac{S_R}{n-p}} \quad (5.34)$$

má Fisherovo rozdělení F s $(p-1)$ a $(n-p)$ stupni volnosti. Na hladině významnosti α je kritický obor vymezen nerovností:

$$T > F_{1-\alpha}(p-1, n-p), \quad (5.35)$$

kde $F_{1-\alpha}(p-1, n-p)$ je příslušný kvantil rozdělení. Pokud hodnota testového kritéria padne do kritického oboru, tedy pokud platí (5.35), potom H_0 zamítáme, což znamená, že některá z vysvětlujících proměnných má statisticky významný efekt na vysvětlovanou proměnnou y . Pokud však nulovou hypotézu nelze na dané hladině významnosti zamítnout, pak vysvětlující proměnné x_i nemají statisticky významný efekt na y .



ŘEŠENÁ ÚLOHA 5.1

Při zjišťování vlivů na pracovní neschopnost zaměstnanců 10 podniků byly získány následující údaje:

Průměrný věk (roky)	Podíl žen v počtu pracovníků (%)	Pracovní neschopnost (%)
37	55	4,4
33	32	0,7
46	59	7,6
34	36	1,8
25	18	0,1
32	47	3,4
38	22	1,6
40	36	3,5
32	29	3,3
41	38	4,7

- Odhadněte parametry lineární regresní funkce popisující závislost pracovní neschopnosti na průměrném věku zaměstnanců a na podílu žen mezi zaměstnanci.
- Pomocí koeficientu determinace charakterizujte přiléhavost daného regresního modelu k datům.
- Jak se změní pracovní neschopnost zaměstnanců, zvýší-li se jejich průměrný věk o 2 roky při stejném podílu žen?
- Určete 95 % intervaly spolehlivosti pro regresní koeficienty b_0, b_1, b_2 .
Na hladině významnosti $\alpha = 0,01$ testujte hypotézu $\beta_1 = \beta_2 = 0$.

Řešení:

a. Naším úkolem je nalézt regresní koeficienty b_0, b_1, b_2 regresní funkce

$$Y = b_0 + b_1X_1 + b_2X_2,$$

kde X_1 je průměrný věk zaměstnanců,

X_2 je podíl žen v počtu zaměstnanců.

Regresní koeficienty b_0, b_1, b_2 vypočítáme pomocí metody nejmenších čtverců. Využijeme přitom nejprve maticové symboliky, kterou jsme použili v textu.

$$\mathbf{F} = \begin{bmatrix} 1 & 37 & 55 \\ 1 & 33 & 32 \\ 1 & 46 & 59 \\ 1 & 34 & 36 \\ 1 & 25 & 18 \\ 1 & 32 & 47 \\ 1 & 38 & 22 \\ 1 & 40 & 36 \\ 1 & 32 & 29 \\ 1 & 41 & 38 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 4,4 \\ 0,7 \\ 7,6 \\ 1,8 \\ 0,1 \\ 3,4 \\ 1,6 \\ 3,5 \\ 3,3 \\ 4,7 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}.$$

Vektor \mathbf{b} vypočítáme pomocí vztahu (5.7). Matice $\mathbf{F}^T\mathbf{F}$ a $\mathbf{F}^T\mathbf{y}$ mají obecně tvar:

$$\mathbf{F}^T\mathbf{F} = \begin{bmatrix} n & \sum x_{1i} & \sum x_{2i} \\ \sum x_{1i} & \sum x_{1i}^2 & \sum x_{1i}x_{2i} \\ \sum x_{2i} & \sum x_{1i}x_{2i} & \sum x_{2i}^2 \end{bmatrix}, \quad \mathbf{F}^T\mathbf{y} = \begin{bmatrix} \sum y_i \\ \sum x_{1i}y_i \\ \sum x_{2i}y_i \end{bmatrix}.$$

Hodnoty potřebné k výpočtu těchto matic jsou uvedeny v následující tabulce:

Pozorování	X_1	X_2	Y	X_1^2	X_2^2	X_1X_2	X_1Y	X_2Y
1	37	55	4,4	1369	3025	2035	162,8	242,0
2	33	32	0,7	1089	1024	1056	23,1	22,4
3	46	59	7,6	2116	3481	2714	349,6	448,4
4	34	36	1,8	1156	1296	1224	61,2	64,8
5	25	18	0,1	625	324	450	2,5	1,8
6	32	47	3,4	1024	2209	1504	108,8	159,8
7	38	22	1,6	1444	484	836	60,8	35,2
8	40	36	3,5	1600	1296	1440	140,0	126,0
9	32	29	3,3	1024	841	928	105,6	95,7
10	41	38	4,7	1681	1444	1558	192,7	178,6
Σ	358	372	31,1	13128	15424	13745	1207,1	1374,7

Potom

$$\mathbf{F}^T\mathbf{F} = \begin{bmatrix} 10 & 358 & 372 \\ 358 & 13128 & 13745 \\ 372 & 13745 & 15424 \end{bmatrix} \quad \mathbf{F}^T\mathbf{y} = \begin{bmatrix} 31,1 \\ 1207,1 \\ 1374,7 \end{bmatrix}.$$

K matici $\mathbf{F}^T\mathbf{F}$ musíme vypočítat matici inverzní:

$$(\mathbf{F}^T\mathbf{F})^{-1} = \begin{bmatrix} 4,355 & -0,131 & -0,012 \\ -0,131 & 0,005 & -0,001 \\ -0,012 & -0,001 & 0,001 \end{bmatrix}.$$

Vektor \mathbf{b} je výsledkem součinu matic $(\mathbf{F}^T\mathbf{F})^{-1}$ a $\mathbf{F}^T\mathbf{y}$:

$$(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\mathbf{y} = \begin{bmatrix} -6,59 \\ 0,18 \\ 0,09 \end{bmatrix}.$$

Hledaná regresní funkce má tvar: $Y = -6,59 + 0,18x_1 + 0,09x_2$.

b. K tomu, abychom vypočítali determinační koeficient, musíme znát hodnotu teoretického součtu čtverců S_T a celkového součtu čtverců S_Y . Tyto součty vypočítáme podle vztahů (5.17), (5.16). Pro výpočet teoretického součtu musíme pro každé $x_{1i}, x_{2i}, i = 1, \dots, 10$, znát teoretickou hodnotu $Y_i, i = 1, \dots, 10$, např. Y_1 vypočítáme takto:

$$Y_1 = -6,59 + 0,18x_{11} + 0,09x_{22} = -6,59 + 0,18 \cdot 37 + 0,09 \cdot 55 = 5,02$$

	X_1	X_2	y	Y	$(y - \bar{y})^2$	$(Y - \bar{Y})^2$
1	37	55	4,4	5,02	1,664	3,648
2	33	32	0,7	2,23	5,808	0,774
3	46	59	7,6	7,00	20,160	15,132
4	34	36	1,8	2,77	1,716	0,116
5	25	18	0,1	-0,47	9,060	12,816
6	32	47	3,4	3,40	0,084	0,084
7	38	22	1,6	2,23	2,280	0,774
8	40	36	3,5	3,85	0,152	0,548
9	32	29	3,3	1,78	0,036	1,769
10	41	38	4,7	4,21	2,528	1,210
Součet	358	372	31,1	32,02	43,489	36,872

Tato hodnota udává, jaká by měla být teoreticky pracovní neschopnost při průměrném věku zaměstnanců téměř 37 let a podílu žen v počtu pracovníků 55%. Protože však jde o stochastickou závislost, liší se tato hodnota od skutečně zjištěné hodnoty $y = 4,4$. Všechny teoretické hodnoty Y_i jsou uvedeny v následující tabulce. Jednotliví sčítanci i hodnoty součtů S_y a S_T jsou rovněž uvedeni v tabulce.

Koeficient determinace vypočítáme dosazením do vztahu (3.20):

$$R^2 = \frac{S_T}{S_y} = \frac{36,87}{43,49} = 0,848.$$

Tato hodnota znamená, že pomocí regresní funkce $Y = -6,59 + 0,18x_1 + 0,09x_2$ je vysvětleno 84,8% celkové variability proměnné Y .

c. Velikost změny znaku Y je při změně znaku X_1 o jednotku rovna b_1 . Má-li se tedy zvýšit průměrný věk o 2 roky při nezměněné zaměstnanosti žen X_2 , zvýší se pracovní neschopnost o $2b_1$, tj. o 0,36%.

d. Obecný tvar těchto intervalů je následující (viz (3.25)):

$$\left[b_i - t_{1-\alpha/2}(n-p) \sqrt{\frac{S_R h_{ii}}{n-p}}, b_i + t_{1-\alpha/2}(n-p) \sqrt{\frac{S_R h_{ii}}{n-p}} \right],$$

kde S_R je reziduální součet čtverců,

$t_{1-\alpha/2}(n-p)$ je kvantil t -rozdělení o $n-p$ stupních volnosti,

p je počet parametrů regresní funkce,

h_{ii} prvek matice $\mathbf{H} = (\mathbf{F}'\mathbf{F})^{-1}$.

Hodnotu S_R vypočítáme ze vztahu:

$$S_R = S_y - S_T = 43,49 - 36,87 = 6,62.$$

V tabulce t -rozdělení nalezneme $(1-\alpha/2) = 97,5$ % kvantil t -rozdělení o $n-p = 10-3 = 7$ stupních volnosti:

$$t_{0,975}(7) = 2,365,$$

$$h_{00} = 4,355; h_{11} = 0,0051; h_{22} = 0,001, \mathbf{H} = \{h_{ij}\}, i, j = 0, 1, 2.$$

Dosazením výše vypočítaných hodnot do vztahu pro interval spolehlivosti určíme jeho pravou a levou krajní hodnotu L a P :

Pro b_0 , tj. $i = 0$:

$$L = 6,59 - 2,365 \sqrt{\frac{6,62 \cdot 4,355}{7}} = 1,79,$$

$$P = 6,59 + 2,365 \sqrt{\frac{6,62 \cdot 4,355}{7}} = 11,39.$$

95 % interval spolehlivosti pro regresní koeficient b_0 je [1,79;11,39]. Pro b_1 , tj. $i = 1$:

$$L = 0,18 - 2,365 \sqrt{\frac{6,62 \cdot 0,0051}{7}} = 0,016,$$

$$P = 0,18 + 2,365 \sqrt{\frac{6,62 \cdot 0,0051}{7}} = 0,344.$$

Pak 95 % interval spolehlivosti pro regresní koeficient b_1 je [0,016; 0,344].

Pro b_2 , tj. $i = 2$:

$$L = 0,09 - 2,365 \sqrt{\frac{6,62 \cdot 0,001}{7}} = 0,017,$$

$$P = 0,09 + 2,365 \sqrt{\frac{6,62 \cdot 0,001}{7}} = 0,163.$$

Potom 95 % interval spolehlivosti pro regresní koeficient b_2 je [0,017; 0,163].

e. Pro ověření hypotézy použijeme F-test. Budeme testovat nulovou hypotézu:

$$H_0: \beta_1 = \beta_2 = 0$$

proti alternativní hypotéze

$$H_1: \text{alespoň jedno } \beta_i \text{ je různé od nuly.}$$

K ověření nulové hypotézy použijeme testové kritérium (3.34):

$$F = \frac{\frac{S_T}{p-1}}{\frac{S_R}{n-p}} = \frac{\frac{36,87}{2}}{\frac{6,62}{7}} = 19,49.$$

V tabulce F-rozdělení najdeme $(1-\alpha)\%$ kvantil F-rozdělení o $p-1$ a $n-p$ stupních volnosti:

$$F_{1-0,01}(2,7) = 9,55.$$

Protože je $19,49 > 9,55$, zamítáme nulovou hypotézu ve prospěch alternativní hypotézy, což znamená, že regresní parametry jsou vesměs nenulové, a tudíž existuje statisticky významná závislost Y na X_1 nebo X_2 .

Řešení v Excelu.

Regresní statistika	
Násobné R	0,912
Hodnota spolehlivosti R (koeficient determinace)	0,831
Nastavená hodnota spolehlivosti R	0,783
Chyba stř. hodnoty	1,024
Pozorování	10

ANOVA					
	Rozdíl	SS	MS	F	Významnost F
Regrese	2	36,155	18,078	17,255	0,002
Rezidua	7	7,334	1,048		
Celkem	9	43,489			

e) Protože hodnota Významnost F je menší než hladina významnosti 0,01; nulovou hypotézu zamítáme, tzn. že regresní parametry jsou vesměs nenulové.

	Koeficienty	Chyba stř. hodnoty	t Stat	Hodnota P	Dolní 95%	Horní 95%
Hranice	-6,595	2,136	-3,087	0,018	-11,645	-1,544
průměrný věk X1	0,178	0,073	2,441	0,045	0,006	0,351
podíl žen (%) X2	0,089	0,032	2,758	0,028	0,013	0,166

ŘEŠENÁ ÚLOHA 5.2



Následující tabulka obsahuje údaje o tržbách, velikosti výdajů na reklamu a o počtu obchodních zástupců pro 11 firem zabývajících se nákupem a prodejem:

Reklamní výdaje (tis. Kč)	Obchodní zástupci	Objem prodeje (mil. Kč)
180	35	260
230	38	310
260	33	280
240	40	300
280	38	340
300	32	380
340	42	410
320	49	440
360	53	400
380	55	430
260	33	310

- Popište závislost objemu produkce na reklamních výdajích a na počtu obchodních zástupců dvourozměrný lineárním regresním modelem.
- F-testem posuďte významnost tohoto regresního modelu. Uvažujte hladinu významnosti $\alpha = 0,01$.
- Na hladině významnosti $\alpha = 0,01$ testujte individuální významnost regresního parametru β_1 .
- Jaký objem produkce lze očekávat, vydá-li firma na reklamu 450 tis. Kč a současně bude mít 50 obchodních zástupců? Určete bodový odhad objemu produkce.

Řešení:

Regresní statistika	
Násobné R	0,916
Hodnota spolehlivosti R	0,839
Nastavená hodnota spolehlivosti R	0,799
Chyba stř. hodnoty	28,434
Pozorování	11

koeficient determinace

ANOVA

	Rozdíl	SS	MS	F	Významnost F
Regrese	2	33822,799	16911,399	20,917	0,001
Rezidua	8	6468,110	808,514		
Celkem	10	40290,909			

b) Hodnota Významnost F je menší než 0,01;
 model je zvolen správně,
 zamítáme nulovou hypotézu o nulovosti obou koeficientů

	Koeficienty	Chyba stř. hodnoty	t Stat	Hodnota P
Hranice	63,830	47,652	1,340	0,217
reklamní výdaje (tis. Kč)	0,849	0,224	3,789	0,005
obchodní zástupci	1,076	1,656	0,650	0,534

a) $Y = 63,83 + 0,85 \cdot x_1 + 1,08 \cdot x_2$

c) Koeficient $b_1 = 0,849$ je statisticky významný na hladině významnosti 0,01; protože Hodnota P je menší než 0,01.

d) 500,33 mil. Kč



SAMOSTATNÉ ÚKOLY

5.1 Firma sledovala, jak jsou její tržby ovlivněny výdaji na reklamu v různých sdělovacích prostředcích. Výsledky průzkumu jsou uvedeny v následující tabulce.

Rádio, TV (tis. Kč)	Noviny, časopisy (tis. Kč)	Tržby (tis. Kč)
0	16	254
22	29	765
28	30	864
33	35	1001
39	27	911
41	36	1121
49	0	856
55	12	932
60	23	1152

63	34	1403
68	54	1702

- Určete jednoduchý lineární regresní model popisující závislost obratu na velikosti prostředků vydaných na reklamu v novinách a časopisech.
- Určete dvourozměrný lineární regresní model popisující závislost obratu na velikosti prostředků vydaných na reklamu v novinách a časopisech a na velikosti prostředků vydaných na reklamu v rozhlasu a v televizi.
- Pomocí F-testu rozhodněte, je-li vhodné k popisu závislosti používat zvolený vícenásobný lineární model. Uvažujte hladinu významnosti $\alpha = 0,05$.
- Přispělo významně zavedení další vysvětlující proměnné k zlepšení výstižnosti modelu?
- Jaký obrat je možné očekávat, vydá-li se na reklamu v tisku 32 tis. Kč a na reklamu v rozhlasu a televizi 47 tis. Kč? Proveďte bodový odhad.

5.2 Mezinárodní organizace WHO zjistila údaje o dětské úmrtnosti (v promile) - DÚ, gramotnosti žen (v procentech) - GŽ a HDP na hlavu (v dolarech) - HDP u 64 rozvojových zemí:

DÚ	GŽ	HDP	DÚ	GŽ	HDP
128	37	1870	142	50	8640
204	22	130	104	62	350
202	16	310	287	31	230
197	65	570	41	66	1620
96	76	2050	312	11	190
209	26	200	77	88	2090
170	45	670	142	22	900
240	29	300	262	22	230
241	11	120	215	12	140
55	55	290	246	9	330
75	87	1180	191	31	1010
129	55	900	182	19	300
24	93	1730	37	88	1730
165	31	1150	103	35	780
94	77	1160	67	85	1300
96	80	1270	143	78	930
148	30	580	83	85	690
98	69	660	223	33	200
161	43	420	240	19	450
118	47	1080	312	21	280
269	17	290	12	79	4430
189	35	270	52	83	270
126	58	560	79	43	1340
12	81	4240	61	88	670
167	29	240	168	28	410
135	65	430	28	95	4370
107	87	3020	121	41	1310
72	63	1420	115	62	1470
128	49	420	186	45	300
27	63	19830	47	85	3630
152	84	420	178	45	220
224	23	530	142	67	560

- Určete lineární regresní model popisující závislost dětské úmrtnosti na gramotnosti žen a HDP v rozvojových zemích.
- Pomocí F–testu rozhodněte, je-li vhodné k popisu závislosti používat zvolený vícenásobný lineární model. Uvažujte hladinu významnosti $\alpha = 0,05$.
- Jsou regresní koeficienty modelu statisticky významné? Stanovte jejich intervaly spolehlivosti pro hladinu významnosti $\alpha = 0,10$.
- Pomocí koeficientu determinace určete přiléhavost dat k modelu. Jak se změní dětská úmrtnost při zvýšení HDP o 1000 USD při stejném stupni ngramotnosti žen? Naopak: jak se změní dětská úmrtnost při zvýšení gramotnosti žen o 1 procento při stejné úrovni HDP?



ODPOVĚDI

5.1 a) jednoduchý lineární regresní model

Regresní statistika	
Násobné R	0,658
Hodnota spolehlivosti R	0,433
Nastavená hodnota spolehlivosti R	0,370
Chyba stř. hodnoty	292,354
Pozorování	11

ANOVA					
	Rozdíl	SS	MS	F	Významnost F
Regrese	1	587103,478	587103,478	6,869	0,028
Rezidua	9	769235,250	85470,583		
Celkem	10	1356338,727			

	Koeficienty	Chyba stř. hodnoty	t Stat	Hodnota P
Hranice	538,482	195,714	2,751	0,022
Noviny, časopisy (tis. Kč)	17,019	6,494	2,621	0,028

$$Y = 539,5 + 17,2 \cdot x$$

b) dvourozměrný lineární regresní model

Regresní statistika	
Násobné R	0,992
Hodnota spolehlivosti R	0,985
Nastavená hodnota spolehlivosti R	0,981
Chyba stř. hodnoty	50,634
Pozorování	11

ANOVA					
	Rozdíl	SS	MS	F	Významnost F

Regrese	2	1335828,082	667914,041	260,514	0,000
Rezidua	8	20510,645	2563,831		
Celkem	10	1356338,727			

	Chyba stř. hod-			
	Koeficienty	noty	t Stat	Hodnota P
Hranice	87,214	42,969	2,030	0,077
Rádio, TV (tis. Kč)	13,905	0,814	17,089	0,000
Noviny, časopisy (tis. Kč)	12,275	1,158	10,596	0,000

$$Y = 87,21 + 13,9 \cdot x_1 + 12,27 \cdot x_2$$

c) Ano, hodnota Významnost F je menší než 0,05; proto vícenásobný lineární model je vhodný.

d) Ano, koeficient determinace se z hodnoty 0,43 zvýšil na hodnotu 0,98.

e) 1 133,15 tis. Kč = 1 133 150 Kč

5.2 a)

Regresní statistika	
Násobné R	0,841
Hodnota spolehlivosti R	0,708
Nastavená hodnota spolehlivosti R	0,698
Chyba stř. hodnoty	41,748
Pozorování	64

ANOVA

	Rozdíl	SS	MS	F	Významnost F
Regrese	2	257362,373	128681,187	73,833	0,000
Rezidua	61	106315,627	1742,879		
Celkem	63	363678,000			

	Chyba stř. hod-				Horní	
	Koeficienty	noty	t Stat	Hodnota P	Dolní 90,0%	90,0%
Hranice	263,642	11,593	22,741	0,000	244,278	283,005
GŽ	-2,232	0,210	-10,629	0,000	-2,582	-1,881
HDP	-0,006	0,002	-2,819	0,006	-0,009	-0,002

$$Y = 263,64 - 2,23 \cdot x_1 - 0,006 \cdot x_2$$

b) Ano, hodnota Významnost F je menší než 0,05; proto vícenásobný lineární model je vhodný.

c) Oba regresní koeficienty jsou statisticky významné, protože Hodnota P je menší než 0,1. Intervaly spolehlivosti: $b_1 \in (-2,5; -1,8)$; $b_2 \in (-0,009; -0,002)$

d) Koeficient determinace je roven 0,71; tzn., že 71 % celkové variability je vysvětleno modelem.

e) Při zvýšení HDP o 1000 USD při stejném stupni ngramotnosti žen klesne dětská úmrtnost o 5,6 promile. Při zvýšení gramotnosti žen o 1 %, při stejné úrovni HDP, klesne dětská úmrtnost o 0,22 promile.



SHRNUTÍ KAPITOLY

V této kapitole jste se seznámili s vícenásobným lineárním regresním modelem. Lineární regresní model byl rozšířen na vícenásobný regresní model lineární v parametrech. Odhady regresních koeficientů byly opět stanoveny metodou nejmenších čtverců, přitom bylo využito maticové symboliky, která usnadňuje práci s vektory a maticemi. Podobně jako v případě jednoduché regrese byly formulovány předpoklady klasického regresního modelu.

6 REGRESNÍ ANALÝZA – VÍCEROZMĚRNÁ: MULTIKOLINEARITA, HETEROSKEDASTICITA, AUTOKORELACE

RYCHLÝ NÁHLED KAPITOLY



V této kapitole se naučíte identifikovat, analyzovat a odstraňovat problémy, které způsobuje nesplnění hlavních předpokladů klasického vícerozměrného lineárního regresního modelu formulované v kapitole 5.4: multikolinearita, heteroskedasticita a autokorelace. Multikolinearitou tedy rozumíme vzájemnou statistickou závislost, tj. korelaci, mezi vysvětlujícími proměnnými ve vícenásobném lineárním regresním modelu. Další důležitou vlastností klasického lineárního regresního modelu je homoskedasticita. Jde o vlastnost (5.15), která spočívá v tom, že rozptyl poruchy ε_i v populačním lineárním regresním modelu je konstantní. Autokorelace je korelace mezi pozorováními uspořádanými v čase, (data jsou časové řady) nebo v prostoru (data jsou průřezová, tj. v jednom časovém okamžiku/intervalu). Říkáme, že v regresním modelu není přítomná autokorelace, jestliže náhodné veličiny jsou vzájemně nekorelované.

CÍLE KAPITOLY



Po prostudování této kapitoly budete umět:

- uvést předpoklady klasického vícerozměrného lineárního modelu,
- identifikovat multikolinearitu, heteroskedasticitu a autokorelaci v modelu,
- aplikovat Bartletův test heteroskedasticity v Excelu.

ČAS POTŘEBNÝ KE STUDIU



K prostudování této kapitoly budete potřebovat asi 90 minut.



KLÍČOVÁ SLOVA KAPITOLY

Multikolinearita, heteroskedasticita, autokorelace, Bartletův test heteroskedasticity.

6.1 Co je multikolinearita?

Multikolinearitou tedy rozumíme vzájemnou statistickou závislost, tj. korelaci, mezi vysvětlujícími proměnnými ve vícenásobném lineárním regresním modelu:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon. \quad (6.1)$$

Informaci o této vzájemné závislosti poskytuje matice výběrových korelačních koeficientů:

$$\mathbf{R} = \begin{bmatrix} 1 & r_{12} & \dots & r_{1k} \\ r_{21} & 1 & \dots & r_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ r_{k1} & r_{k2} & \dots & 1 \end{bmatrix}. \quad (6.2)$$

Zřejmě je matice (6.2) symetrická, tj. $r_{ij} = r_{ji}$ pro všechna i, j . Pokud jsou všechny dvojice vysvětlujících proměnných vzájemně nekorelované, potom platí, že $r_{ij} = r_{ji} = 0$, tj. $\mathbf{R} = \mathbf{I}$, čili \mathbf{R} je jednotkovou maticí.

Uvědomte si, že na diagonále matice \mathbf{R} musejí být všechny prvky rovny 1, neboť korelace vektoru dat se sebou samým je vždy rovna 1! Jsou-li však alespoň některé nediagonální prvky matice \mathbf{R} nenulové, hovoříme o *multikolinearitě*. Matice \mathbf{R} pak není jednotkovou maticí a její determinant je menší než 1. Je-li multikolinearita vysoká, hovoříme o *škodlivé multikolinearitě*, pak se determinant matice \mathbf{R} blíží k nule. V tom případě dává metoda nejmenších čtverců odhady regresních koeficientů s širokými intervaly spolehlivosti, takže výsledky jsou prakticky neupotřebitelné.

Na to, kdy je multikolinearita „škodlivá“, existují různé názory, opírající se víceméně o zkušenost. Někteří autoři považují za škodlivou multikolinearitu, když alespoň jeden nediagonální prvek matice \mathbf{R} je větší než 0,8.

Zjistí-li se škodlivá multikolinearita, je možno postupovat v zásadě dvojím způsobem. Buď vysvětlující proměnnou, která je zdrojem multikolinearity, vypustíme z modelu, nebo doplníme data, eventuálně získáme nový vzorek dat. Škodlivá multikolinearita je totiž často důsledkem „špatného“ vzorku dat. Projevuje se obvykle vysokým koeficientem determinace (blízkým k 1) a zároveň jsou individuální koeficienty statisticky nevýznamné (t-test), model jako celek je naopak statisticky významný (F-test), viz kap. 5.7 a 5.8.

Celou záležitost ilustrujeme na řešené úloze 6.1.

ŘEŠENÁ ÚLOHA 6.1



V následující Tabulce 11 jsou uvedeny měsíční výdaje, měsíční příjmy a majetek (v Kč) u 10 českých rodin. Provedte regresní analýzu měsíčních výdajů rodin v závislosti na měsíčních příjmech a majetku. Vysvětlete dosažené výsledky pomocí jednorozměrné regrese.

Tabulka 11: Měsíční výdaje, příjmy a majetek v Kč

Y výdaje	X1 příjmy	X2 majetek
8400	9600	100000
7800	12000	120000
10800	14400	150000
11400	16800	170000
13200	19200	200000
13800	21600	225000
14400	24000	246000
16800	26400	264000
18600	28800	392000
18000	31200	322000

Řešení:

Data z Tabulky 11 uložíme v excelovské tabulce. Známým postupem v menu: Data → Analýza dat... → Regrese, a získáme po vyplnění příslušných políček tento výsledek:

VYSLEDEK						
<i>Regresní statistika</i>						
Násobné R	0,981					
Hodnota spolehlivosti R	0,962					
Nastavená hodnota spolk	0,951					
Chyba stř. hodnoty	832,660					
Pozorování	10					
ANOVA						
	Rozdíl	SS	MS	F	ýznamnost F	
Regrese	2	1,23E+08	61581370	88,82062	1,06E-05	
Rezidua	7	4853260	693322,9			
Celkem	9	1,28E+08				
	Koeficienty	ba stř. hodí	t stat	Hodnota P	Dolní 95%	Horní 95%
Hranice	2943,676	832,579	3,536	0,010	974,940	4912,413
X1 příjmy	0,569	0,847	0,672	0,523	-1,433	2,571
X2 majetek	-0,006	0,083	-0,071	0,946	-0,203	0,191

V tomto výstupu se vyskytují zdánlivě paradoxní výsledky. Z Tabulky ANOVA vyplývá, že regresní model

$$y = 2943,676 + 0,569x_1 - 0,006x_2 + \varepsilon$$

je jako celek statisticky významný (F-test), zatímco individuální regresní koeficienty u proměnných „příjmy“ resp. „majetek“ jsou statisticky nevýznamné, neboť obě odpovídající p-hodnoty (signifikance) jsou větší než 0,05 (0,523 resp. 0,946). Koeficient determinace $R^2 = 0,962$ je vysoký – blízký k 1, což svědčí o vysoké přiléhavosti dat k modelu. Navíc je u regresního koeficientu u proměnné x_2 záporné znaménko, což je evidentně v rozporu s intuicí, která říká: čím je větší majetek, tím je vyšší spotřeba rodiny. Tento zdánlivý rozpor je způsoben kolinearitou regresorů, o čemž svědčí jejich korelační matice

$$\mathbf{R} = \begin{bmatrix} 1,000 & 0,999 \\ 0,999 & 1,000 \end{bmatrix},$$

kteřou lze snadno zjistit tak, že vypočítáte $r_{12} = r_{21} = 0,999012$ pomocí excelovské funkce =CORREL (B4:B13;C4:C13), za předpokladu, že data pro x_1 jsou uložena v oblasti B4:B13, data pro x_2 jsou uložena v oblasti C4:C13. Vysvětlující proměnné x_1 a x_2 jsou kolineární, neboť koeficient korelace $r_{12} = r_{21} = 0,999012$ je blízký k 1.

Vypustíme-li nyní jednu z vysvětlujících proměnných, např. x_2 – majetek, a provedeme-li (jednoduchou) regresi x_1 na y , obdržíme s analogickým využitím Excelu tento výsledek:

VÝSLEDEK						
<i>Regresní statistika</i>						
Násobné R		0,981				
Hodnota spolehlivosti R		0,962				
Nastavená hodnota spol.		0,957				
Chyba stř. hodnoty		779,160				
Pozorování		10				
ANOVA						
	Rozdíl	SS	MS	F	Významnost F	
Regrese	1	1,23E+08	1,23E+08	202,8679	5,75275E-07	
Rezidua	8	4856727	607090,9			
Celkem	9	1,28E+08				
	Koeficienty/ba stř. hod.	t stat	Hodnota P	Dolní 95%	Horní 95%	
Hranice	2934,545	769,658	3,813	0,005	1159,710	4709,381
X1 příjmy	0,509	0,036	14,243	0,000	0,427	0,592

Vidíte, že v novém regresním modelu je regresní koeficient statisticky významný, neboť odpovídající p-hodnota (signifikance) je menší než 0,05 (0,000...), což je ve shodě s tabulkou ANOVA.

Podobně, vypustíme-li nyní vysvětlující proměnnou x_1 – příjem, a provedeme-li (jednoduchou) regresi x_2 na y , obdržíme s analogickým využitím Excelu výsledek z následujícího výstupu. Opět vidíte, že v novém regresním modelu je regresní koeficient statisticky významný, neboť odpovídající p-hodnota (signifikance) je menší než 0,05 (0,000...), což je ve shodě s tabulkou ANOVA. Navíc je znaménko u regresního koeficientu 0,050 kladné, což je v souladu s intuicí, že totiž velikost spotřeby je přímo úměrná velikosti majetku.

VÝSLEDEK						
<i>Regresní statistika</i>						
Násobné R		0,979614				
Hodnota spolehlivosti R		0,959644				
Nastavená hodnota spol.		0,954599				
Chyba stř. hodnoty		803,6024				
Pozorování		10				
ANOVA						
	Rozdíl	SS	MS	F	Významnost F	
Regrese	1	1,23E+08	1,23E+08	190,2357	7,37266E-07	
Rezidua	8	5166214	645776,8			
Celkem	9	1,28E+08				
	Koeficienty	ba stř. hod.	t stat	Hodnota P	Dolní 95%	Horní 95%
Hranice	2880,627	798,404	3,608	0,007	1039,503	4721,750
X2 majetek	0,050	0,004	13,793	0,000	0,042	0,058

6.2 Co je heteroskedasticita?

Další důležitou vlastností klasického lineárního regresního modelu je *homoskedasticita*. Jde o vlastnost (5.15), která spočívá v tom, že rozptyl poruchy ε_i v populačním lineárním regresním modelu je konstantní, tj. v modelu

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad i = 1, 2, \dots, n, \quad (5.1)$$

platí podmínka

$$\text{Var}(\varepsilon) = \sigma^2 \mathbf{I}, \quad (5.15)$$

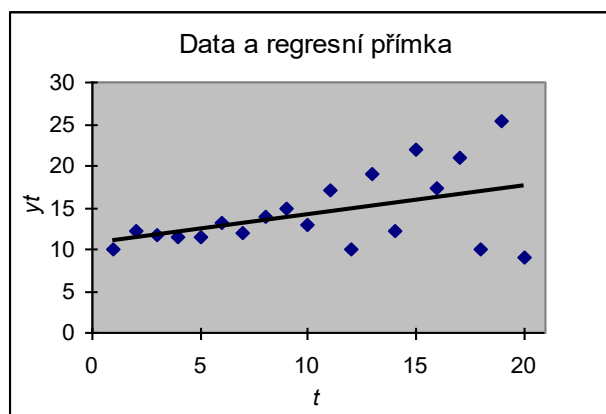
kde symbol \mathbf{I} označuje jednotkovou matici.

Podmínku (5.15) je možné ekvivalentně vyjádřit také takto

$$E(\varepsilon_i^2) = \sigma^2, \quad i = 1, 2, \dots, n, \quad (6.3)$$

kde E je známý operátor střední hodnoty.

Pokud podmínka (5.15) není splněna, potom hovoříme o *heteroskedasticitě*. Příklad heteroskedasticity v případě jednorozměrného lineárního regresního modelu je na Obrázku 26. Je zřejmé, že rozptyl hodnoty y se zvětšuje s rostoucí hodnotou x .



Obrázek 26: Heteroskedasticita v regresním modelu

Heteroskedasticita může být způsobena různými příčinami. Častou příčinou heteroskedasticity je fakt, že při postupném sběru dat se technika sběru postupně zlepšuje a chyba se proto zmenšuje. Naopak se chyba zvětšuje s přítomností odlehlých hodnot. Dalším zdrojem heteroskedasticity je nesprávná specifikace modelu, např. tím, že jsou opominuty důležité vysvětlující proměnné regresního modelu. Přítomnost heteroskedasticity v regresním modelu je silně nežádoucí, a to zejména z těchto důvodů:

- Přítomnost heteroskedasticity způsobuje neplatnost odhadů rozptylů regresních koeficientů, a tudíž také odhadů jejich intervalů spolehlivosti a testů hypotéz o jejich statistické významnosti atd., viz kap. 5.6.
- Prognózy s využitím regresního modelu obsahujícího heteroskedasticitu jsou často nespolehlivé a dokonce nerealistické.

6.2.1 JAK ZJISTIT HETEROSKEDASTICITU?

Jak poznáme, že v regresním modelu, který jsme sestavili na základě nějakých dat, je přítomna heteroskedasticita? Podobně jako v případě multikolinearity neexistují přesná pravidla, jak detekovat přítomnost heteroskedasticitu, pouze pár heuristických zásad.

Velmi často poznáme přítomnost heteroskedasticity z věcné povahy problému. Například je známo, že s rostoucím věkem zaměstnanců se zvětšuje rozptyl jejich platů. Ať je typ závislosti platu na věku lineární nebo ne, bude v modelu přítomna heteroskedasticita.

Pokud však nemáme podobné předběžné empirické informace o povaze problému, předpokládáme, že heteroskedasticita není přítomna, že tudíž je rozptyl náhodné složky modelu konstantní. Takové tvrzení pak můžeme podrobit zkoumání např. grafické analýze nebo statistickému testu reziduí e_i . S oběma postupy se zde seznámíte.

Grafická analýza

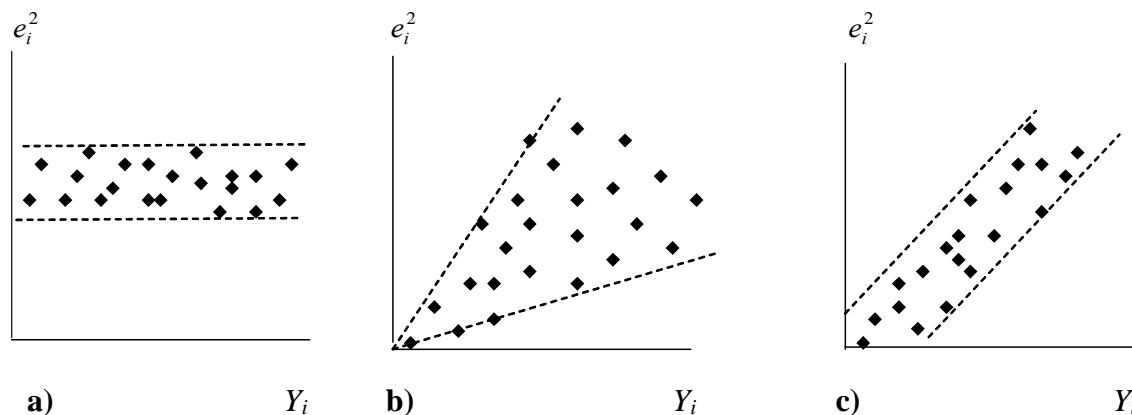
Zobrazíme si závislost kvadrátu reziduí e_i^2 na teoretické hodnotě Y_i . Na Obrázku 22 jsou zobrazeny 3 důležité případy tvaru, které mohou nastat, kde

$$Y_i = b_0 + b_1 f_1(x_{i1}, x_{i2}, \dots, x_{ik}) + \dots + b_k f_k(x_{i1}, x_{i2}, \dots, x_{ik}), \quad (6.4)$$

přičemž b_i jsou odhady regresních parametrů získané MNC,

$$e_i = y_i - Y_i \quad (6.5)$$

je reziduum, tj. odhad náhodné složky ε_i .



Obrázek 27: Závislost e_i^2 na Y_i

Na Obrázku 27 a) hodnota e_i^2 v zásadě nezávisí na Y_i , což naznačuje, že náhodná složka je konstantní, a tudíž heteroskedasticita není přítomna. Na druhou stranu Obr. 27 b) a c) hodnota e_i^2 v zřejmém závisí na Y_i , což naznačuje přítomnost heteroskedasticity. Konkrétní tvar závislosti vám dobře potvrdí zobrazení bodového diagramu závislosti y_i na vybrané datové hodnoty j -té vysvětlující proměnné x_{ji} .

Testy heteroskedasticity

Detekce heteroskedasticity s pomocí statistického testu hypotézy je obvykle založena na nulové hypotéze, že rozptyly náhodné složky ε_i^2 jsou konstantní, přičemž se analyzují jejich odhady, tj. rezidua e_i^2 . V literatuře můžete nalézt podrobné testy heteroskedasticity s názvy jako Parkův test, Glejserův test, Goldfeld-Quandtův test aj., viz např. Gujarati (2003). Tyto statistické testy lze provádět pomocí specializovaných statistických programů, např. SPSS, v Excelu specializované funkce na tyto testy bohužel chybí. My si zde proto ukážeme tzv. Bartletův test heteroskedasticity, který představuje zjednodušený Goldfeld-Quandtův test a lze k jeho provedení využít funkce Excelu.

Bartletův test

Test vychází z rozdělení dat podle velikosti (některé) vysvětlující proměnné – označíme ji X , do dvou částí: $x_i \leq \hat{x}$ a $x_i > \hat{x}$, přičemž jsou data uspořádána podle X , \hat{x} je medián z x_i .

- Testuje se hypotéza o rovnosti rozptylů reziduí v obou částech (v Excelu: Analýza dat, Dvouvýběrový F-test pro rozptyl)
- Pokud se hypotéza o rovnosti rozptylů reziduí (není přítomna heteroskedasticita) v obou částech zamítá, potom se hypotéza o přítomnosti heteroskedasticity, přijímá (a obráceně).

Použití Bartletova testu si ukážeme na příkladu. Ještě předtím se budeme zabývat otázkou, jak odstranit zjištěnou heteroskedasticitu, tj. jak modifikovat původní model, tak aby heteroskedasticitu neobsahoval.

6.2.2 JAK ODSTRANIT HETEROSKEDASTICITU?

Nejznámější metodou k odstranění heteroskedasticity je *metoda vážených nejmenších čtverců* MVNČ. V MVNČ předpokládáme určitý typ nekonstantního chování rozptylu náhodné složky.

Předpoklad 1: Rozptyl náhodné složky je přímo úměrný kvadrátu vysvětlující proměnné x , tj.

$$E(\varepsilon_i^2) = \sigma^2 x_i^2, \quad i = 1, 2, \dots, n. \quad (6.6)$$

Transformovaný regresní model získáme tak, že regresní rovnici

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n, \quad (6.7)$$

vydělíme hodnotou x_i , čímž obdržíme

$$\frac{y_i}{x_i} = \frac{\beta_0}{x_i} + \beta_1 + \frac{\varepsilon_i}{x_i} = \beta_0 \frac{1}{x_i} + \beta_1 + \delta_i, \quad i = 1, 2, \dots, n, \quad (6.8)$$

kde pro novou náhodnou chybu δ_i platí po dosazení z (6.6)

$$E(\delta_i^2) = E\left(\frac{\varepsilon_i^2}{x_i^2}\right) = \sigma^2, \quad i = 1, 2, \dots, n. \quad (6.9)$$

Provedením transformace $y_i' = \frac{y_i}{x_i}$, $x_i' = \frac{1}{x_i}$, $i = 1, 2, \dots, n$.

(6.10)

obdržíme z (6.8) nový regresní model

$$y_i' = \beta_1 + \beta_0 x_i' + \delta_i, \quad i = 1, 2, \dots, n. \quad (6.11)$$

což je nový lineární regresní model podle (6.9) však bez heteroskedasticity.

Uvažovali jsme jednoduchý regresní model, avšak rozšíření výše uvedeného postupu na vícerozměrný regresní model je snadné. Předpoklad 1 modifikujeme tak, že rozptyl náhodné složky je přímo úměrný kvadrátu vysvětlující proměnné x_j , tj.

$$E(\varepsilon_i^2) = \sigma^2 x_{ij}^2, \quad i = 1, 2, \dots, n. \quad (6.6)$$

Namísto modelu (6.7) uvažujeme model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \varepsilon_i, \quad i = 1, 2, \dots, n. \quad (6.7^*)$$

Pro nový vícerozměrný regresní model použijeme namísto transformace (6.10) nová transformovaná data

$$y_i' = \frac{y_i}{x_{ij}}, \quad x_{ij}' = \frac{1}{x_{ij}}, \quad x_{ik}' = \frac{x_{ik}}{x_{ij}}, \quad k \neq j, \quad i = 1, 2, \dots, n. \quad (6.10^*)$$

Předpoklad 2: Rozptyl náhodné složky je přímo úměrný vysvětlující proměnné x , tj.

$$E(\varepsilon_i^2) = \sigma^2 x_i, \quad i = 1, 2, \dots, n. \quad (6.12)$$

Transformovaný regresní model získáme tak, že regresní rovnici

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n, \quad (6.13)$$

vydělíme hodnotou $\sqrt{x_i}$, čímž obdržíme

$$\frac{y_i}{\sqrt{x_i}} = \frac{\beta_0}{\sqrt{x_i}} + \beta_1 \sqrt{x_i} + \frac{\varepsilon_i}{\sqrt{x_i}} = \beta_0 \frac{1}{\sqrt{x_i}} + \beta_1 \sqrt{x_i} + \delta_i, \quad i = 1, 2, \dots, n, \quad (6.14)$$

kde pro novou náhodnou chybu δ_i platí po dosazení z (6.12)

$$E\left(\frac{\varepsilon_i^2}{x_i}\right) = E\left(\frac{\varepsilon_i^2}{x_i}\right) = \sigma^2, \quad i = 1, 2, \dots, n. \quad (6.15)$$

Provedením transformace $y_i' = \frac{y_i}{\sqrt{x_i}}$, $x_i' = \frac{1}{\sqrt{x_i}}$, $x_i'' = \sqrt{x_i}$, $i = 1, 2, \dots, n$.

(6.16)

obdržíme z (6.16) nový regresní model

$$y_i' = \beta_0 x_i' + \beta_1 x_i'' + \vartheta_i, \quad i = 1, 2, \dots, n, \quad (6.17)$$

což je nový lineární regresní model bez úrovně konstanty podle (6.15) však bez heteroskedasticity. Rozšíření na vícerozměrný regresní model je možné udělat analogicky jako v případě Předpokladu 1. Odstranění heteroskedasticity si prakticky vyzkoušíte v následující řešené úloze.



ŘEŠENÁ ÚLOHA 6.2

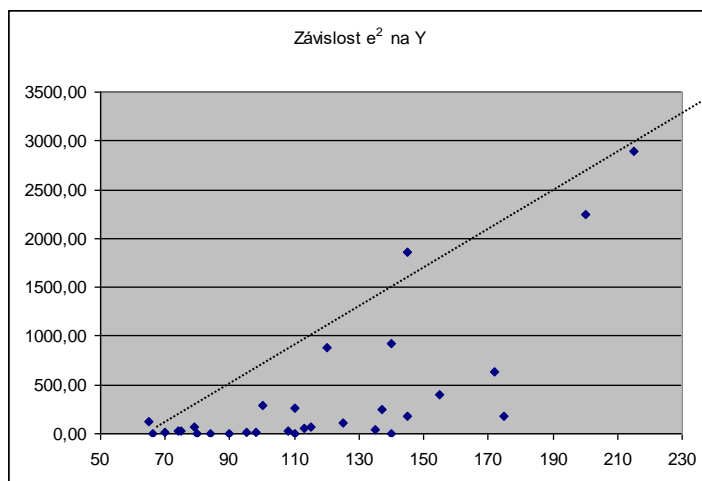
V následující tabulce jsou uvedeny příjmy a spotřební výdaje 30 rodin v tis. Kč/rok. Vytvořte lineární regresní model závislosti výdajů na příjmech, graficky a statistickým testem zjistěte přítomnost heteroskedasticity. Z původního modelu pak heteroskedasticitu odstraňte pomocí MVNČ. Použijte přitom Excel.

č.rodiny	Výdaje	Příjmy	č.rodiny	Výdaje	Příjmy
1	66	80	16	115	180
2	65	100	17	120	225
3	70	85	18	100	170
4	80	110	19	145	240
5	79	120	20	110	185
6	84	115	21	172	220
7	98	130	22	200	230
8	95	140	23	175	245
9	90	125	24	140	260
10	75	90	25	135	190
11	74	105	26	140	205
12	110	160	27	155	200
13	113	150	28	230	270
14	125	165	29	137	230
15	108	145	30	145	290

Řešení:

V Excelu vytvoříme z daných údajů graf: XY bodový a pomocí pravého tlačítka iniciujeme nabídku s volbou Přidat spojnicí trendu... V podnabídce Možnosti zaklikneme 2 položky: Zvolit rovnici regrese a Zvolit koeficient spolehlivosti (tj. koeficient determinace). Obdržíme výsledek, z něhož vyplývá lineární regresní model: $y = 9,29 + 0,64 \cdot x + \varepsilon$.

Dále vedle sloupce y_i vytvoříme pomocí vzorce regresní rovnice sloupec teoretických hodnot Y_i . Další sloupec vytvoříme jako rozdíl sloupců y_i a Y_i , což bude sloupec reziduí. Poslední sloupec bude druhá mocnina reziduí. Společně pak vytvoříme XY bodový graf mezi Y_i a e_i^2 . Výsledkem je následující graf na Obr. 28, který napovídá přítomnost heteroskedasticity, neboť body v grafu netvoří pás rovnoběžný s vodorovnou osou, jako na Obr. 27 a), ale spíše kužel, jako na Obr. 27 b).

Obrázek 28: Kužel závislosti e_i^2 na Y_i

K exaktnímu prokázání heteroskedasticity použijeme Bartletův test. Podle rostoucích hodnot X – Příjmů seřadíme hodnoty reziduí a z nich vytvoříme dva stejně velké soubory e_1 a e_2 :

Příjmy	e_1	Příjmy	e_2
80	1,99	170	-8,09
85	-10,83	180	-29,68
90	3,03	185	-17,19
100	-1,74	190	-13,54
105	-8,65	200	-16,05
110	-0,69	205	25,28
115	4,45	220	47,37
120	-4,46	225	13,51
125	-0,60	230	-30,36
130	5,08	230	6,00
140	-4,78	240	2,14
145	-1,28	245	20,09
150	7,63	260	53,74
160	10,77	270	-15,63
165	5,58	290	-43,08

Budeme testovat, zda rozptyly obou souborů jsou stejné pomocí F-testu z Excelu:

V menu: Data → Analýza dat → Dvouvýběrový F-test pro rozptyl zadáme umístění oblastí sloupců e_1 a e_2 , eventuální popisky a oblast výstupu. Obdržíme výstup:

Dvouvýběrový F-test pro rozptyl

	Soubor 1	Soubor 2
Stř. hodnota	0,366225	-0,366225
Rozptyl	35,88461	792,7791
Pozorování	15	15
Rozdíl	14	14
F	0,045264	
P(F<=f) (1)	3,89E-07	
F krit (1)	0,402621	

V tomto výstupu je důležitá P-hodnota: $P(F \leq f) (1) = 3,89 \text{ E-}07 = 0,000000389 < 0,05$.

Na hladině $\alpha = 0,05$ proto nulovou hypotézu H_0 : „Rozptyly obou uvažovaných souborů jsou stejné“ zamítáme. Uvažované soubory mají různý rozptyl, což znamená, že rozptyl náhodné složky regresního modelu není konstantní neboli, že heteroskedasticita je v modelu přítomna.

Nakonec ukážeme, jak přítomnou heteroskedasticitu odstranit. V Obr. 28 se body grafu nacházejí v „lineárním kuželu“, proto zvolíme pro transformaci Předpoklad 2.

Transformace podle (6.16): $y_i' = \frac{y_i}{\sqrt{x_i}}$, $x_i' = \frac{1}{\sqrt{x_i}}$, $x_i'' = \sqrt{x_i}$, $i = 1, 2, \dots, 30$.

obdržíme nový regresní model $y_i' = 16,75x_i' + 0,59_1x_i'' + \vartheta_i$, $i = 1, 2, \dots, 30$,

který je bez heteroskedasticity.

č.rodiny	y'	x'	x''	č.rodiny	y'	x'	x''
1	7,379	0,112	8,944	16	8,572	0,075	13,416
2	6,500	0,100	10,000	17	8,000	0,067	15,000
3	7,593	0,108	9,220	18	7,670	0,077	13,038
4	7,628	0,095	10,488	19	9,360	0,065	15,492
5	7,212	0,091	10,954	20	8,087	0,074	13,601
6	7,833	0,093	10,724	21	11,596	0,067	14,832
7	8,595	0,088	11,402	22	13,188	0,066	15,166
8	8,029	0,085	11,832	23	11,180	0,064	15,652
9	8,050	0,089	11,180	24	8,682	0,062	16,125
10	7,906	0,105	9,487	25	9,794	0,073	13,784
11	7,222	0,098	10,247	26	9,778	0,070	14,318
12	8,696	0,079	12,649	27	10,960	0,071	14,142
13	9,226	0,082	12,247	28	13,997	0,061	16,432
14	9,731	0,078	12,845	29	9,034	0,066	15,166
15	8,969	0,083	12,042	30	8,515	0,059	17,029

6.3 Co znamená autokorelace?

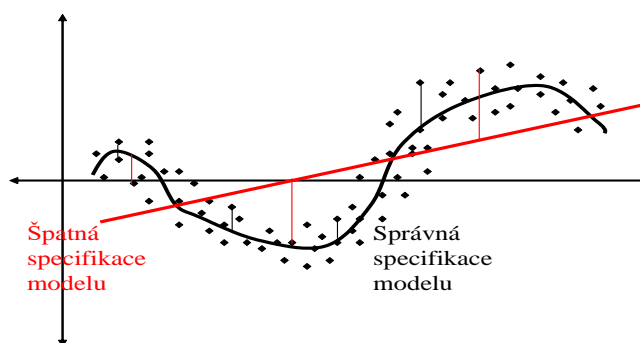
Autokorelace je korelace mezi pozorováními uspořádanými v čase (data jsou časové řady) nebo v prostoru (data jsou průřezová, tj. v jednom časovém okamžiku/intervalu). Říkáme, že v regresním modelu není přítomná autokorelace, jestliže náhodné veličiny jsou vzájemně nekorelované, symbolicky to lze vyjádřit takto

$$E(\varepsilon_i \varepsilon_j) = 0, \quad i \neq j, \quad i, j = 1, 2, \dots, n. \quad (6.18)$$

Jestliže naopak existuje dvojice indexů $i \neq j$, přičemž platí $E(\varepsilon_i \varepsilon_j) \neq 0$, řekneme, že v regresním modelu je *přítomna autokorelace*.

Autokorelace se nejčastěji vyskytuje v regresních modelech založených na datech ve formě časových řad. Potom indexy i , (resp. j) představují časové okamžiky t . Časovým řadám a jejich analýze se budou věnovat následující kapitoly 8 až 12, kde bude podrobněji pojednáno také o autokorelaci.

Následující Obrázek 29 dává příklad dvou regresních modelů dat, z nichž jeden je správně specifikován (nelineární regresní model – černá křivka), druhý je nesprávně specifikován (lineární regresní křivka – červená přímka). Nesprávná specifikace modelu způsobuje, že rezidua jsou vzájemně korelována, což se projevuje tak, že datové body leží vždy ve větší oblasti podél vodorovné osy na jedné straně regresní křivky, zatímco v případě nekorelovaných reziduí leží datové body rovnoměrně po obou stranách regresní křivky v celé oblasti vodorovné osy (tj. nezávisle proměnné).



Obrázek 29: Autokorelace: špatná a správná specifikace modelu

ŘEŠENÁ ÚLOHA 6.3 – GRETL



V následující tabulce jsou uvedena data týkající se největších nemocnic ve státech V4. X_1 značí průměrný denní počet pacientů, X_2 počet obsazených lůžek za měsíc, X_3 velikost populace (v tis.) ve spádové oblasti, X_4 průměrnou délku pobytu v nemocnici (ve dnech), Y počet pracovních hodin, vykázaných za měsíc. Úloha bude řešena pomocí softwaru GRETL. Testujte na hladině významnosti $\alpha = 0,05$.

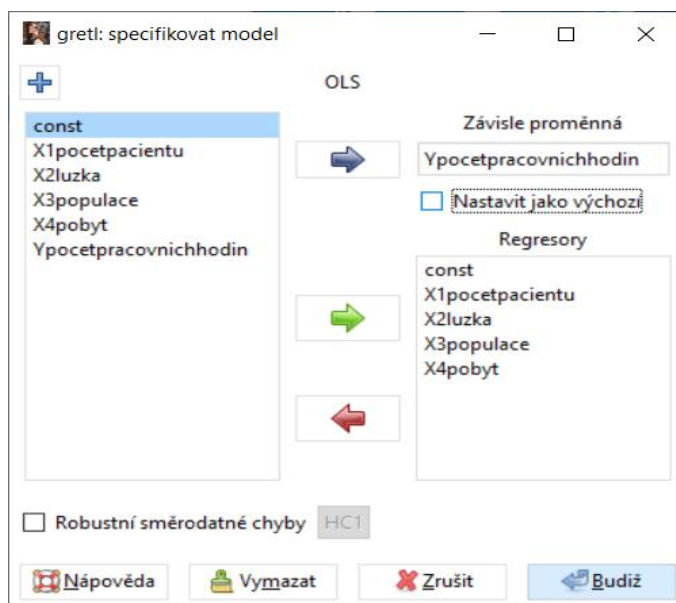
- a) Zkonstruuje regresní model závislosti počtu pracovních hodin na ostatních proměnných. Do modelu zahrňte všechny proměnné a interpretujte výsledky.
- b) Prozkoumejte závislosti mezi proměnnými pomocí párových korelačních koeficientů.
- c) Vyberte vhodnou podmnožinu vysvětlujících proměnných a sestrojte nový regresní model.
- d) Určete předpověď počtu pracovních hodin pro „průměrnou“ nemocnici s průměrným denním počtem pacientů 150, s 5000 obsazenými lůžky za měsíc, se 100 tis. obyvateli ve spádové oblasti a s průměrnou délkou pobytu 6 dní. Použijte jednak původní model se všemi zařazenými proměnnými, jednak redukovaný model, a předpovědi porovnejte.
- e) Proveďte diagnostiku modelu, tzn. kontrolu předpokladů o heteroskedasticitě, normalitě a autokorelaci reziduí.

X1 (počet pacientů)	X2 (lůžka)	X3 (populace)	X4 (pobyt)	Y (počet pracovních hodin)
15,57	472,92	18	4,45	566,52
44,02	1339,75	9,5	6,92	696,82
20,42	620,25	12,8	4,28	1033,15
18,74	568,33	36,4	3,9	1603,62
49,2	1497,6	35,7	5,5	1611,37
44,92	1365,83	24	4,6	1613,27
55,48	1687	43,3	5,62	1854,17
59,28	1639,92	46,7	5,15	2160,55
94,39	2872,33	78,7	6,18	2305,58
128,02	3655,08	180,5	6,15	3503,93
96	2912	60,9	5,88	3571,89
131,42	3921	103,7	4,88	3741,4
127,21	3865,67	126,8	5,5	4026,52
252,9	7684,1	157,7	7	10343,81
409,2	12446,33	169,4	10,78	11732,17
463,7	14098,4	331,4	7,05	15414,94
510,22	15524	371,6	6,35	18854,45

Řešení:

Prezentujeme zde řešení pomocí programu GRETL. Nejprve do programu zadáme všechny proměnné.

- a) V hlavním menu vybereme MODEL→Ordinary Least Squares a objeví se následující dialogové okno, kde doplníme Y (počet pracovních hodin) jako závislou proměnnou, a X₁, X₂, X₃, X₄ jako regresory, tzn. nezávislé proměnné, jak ukazuje Obrázek 30.



Obrázek 30: Dialogové okno – specifikace modelu

Potvrdíme tlačítkem budiž a dostáváme následující výstup, který je zobrazen na Obrázku 31.

Model 1: OLS, za použití pozorování 1-17
Závisle proměnná: Ypocetpracovnichhodin

	koeficient	směr. chyba	t-podíl	p-hodnota	
const	2789,60	1251,11	2,230	0,0456	**
X1pocetpacientu	-26,8215	119,247	-0,2249	0,8258	
X2luzka	2,15068	3,77026	0,5704	0,5789	
X3populace	-1,53312	8,67951	-0,1766	0,8627	
X4pobyt	-561,216	244,108	-2,299	0,0403	**

Střední hodnota závisle proměnné 4978,480
Sm. odchylka závisle proměnné 5560,534
Součet čtverců reziduí 7388576
Sm. chyba regrese 784,6749
Koefficient determinace 0,985065
Adjustovaný koeficient determinace 0,980087
F(4, 12) 197,8692
P-hodnota(F) 7,67e-11
Logaritmus věrohodnosti -134,4709
Akaikovo kritérium 278,9419
Schwarzovo kritérium 283,1079
Hannan-Quinnovo kritérium 279,3560

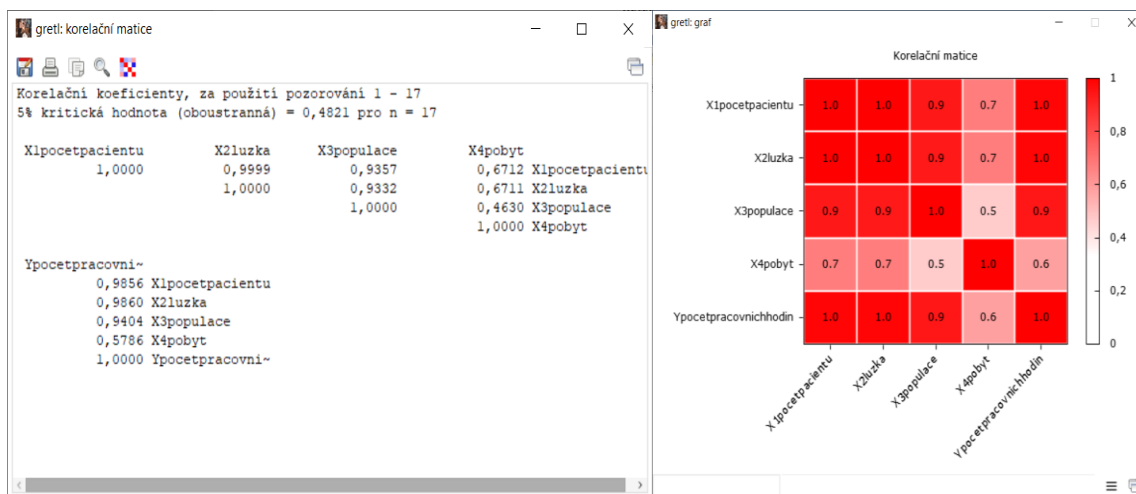
zde je poznámka o zkratkách statistik modelu

Pomine-li se konstanta, p-hodnota byla nejvyšší pro proměnnou 3 (X3populace)

Obrázek 31: Odhad koeficientů metodou nejmenších čtverců

Z tohoto výstupu vidíme, že rovnice modelu je $Y = 2789,6 - 26,8x_1 + 2,1x_2 - 1,5x_3 - 561,2x_4$; přičemž jediný statisticky významný koeficient je koeficient b_4 . To ukazuje na možnou multikolaritu mezi proměnnými. A dává to smysl, protože počet pacientů, počet obsazených lůžek a velikost populace jistě spolu souvisí. Proto v dalším kroku vypočteme korelační matici a podíváme se na vzájemné korelační koeficienty.

b) V nabídce ZOBRAZIT vybereme KORELAČNÍ MATICE a dostaneme výsledek na Obrázku 32.



Obrázek 32: Korelační matice

Z tohoto výsledků vidíme, že v modelu ponecháme proměnnou x_4 (doba pobytu v nemocnici) a z ostatních proměnných ponecháme proměnnou x_2 (počet obsazených lůžek), protože koeficient b_2 měl nejmenší p hodnotu ve výstupu regresního modelu.

Multikolinearita může mít několik negativních dopadů na regresní analýzu:

- Snížení interpretovatelnosti: Když jsou nezávislé proměnné mezi sebou silně korelované, je obtížnější určit, jaký vliv má každá z těchto proměnných na závislou proměnnou. Koeficienty regresního modelu mohou být nepřesné a záviset na konkrétní konfiguraci dat.
- Nestabilita koeficientů: Malé změny v datech mohou vést k velkým změnám v koeficientech regresního modelu, což ztěžuje stabilitu a spolehlivost interpretace.
- Zvýšení variancí koeficientů: Multikolinearita může způsobit zvýšení variancí odhadnutých koeficientů, což může vést k nižší přesnosti predikce a může také zhoršit schopnost modelu generalizovat na nová data.
- Nepřesnost významnosti proměnných: Multikolinearita může vést k nesprávným výsledkům týkajícím se statistické významnosti nezávislých proměnných. Proměnné, které by mohly být významné, mohou být považovány za nevýznamné kvůli vzájemné korelaci s jinými proměnnými.
- Zvýšení rozptylu reziduí: Multikolinearita může také způsobit, že rezidua (odchylky skutečných hodnot od předpovědaných) budou mít vyšší rozptyl, což může naznačovat, že model nepopisuje data efektivně.

Jak se s multikolinearitou zachází? Existuje několik postupů:

- Zpětná eliminace proměnných: Zkuste postupně vyřazovat nebo kombinovat silně korelované proměnné, abyste snížili kolinearitu a zlepšili stabilitu modelu.
- Získání více dat: Větší množství dat může pomoci rozptýlit vliv multikolinearity.
- Transformace proměnných: Transformace dat (např. normalizace, standardizace) může pomoci snížit multikolinearitu.

Celkově je důležité identifikovat a řešit multikolinearitu, aby byly výsledky regresní analýzy spolehlivé a interpretabilní.

- c) Nový model tedy sestavíme bez proměnných x_2 a x_4 . Výsledek vidíme na Obrázku 33, rovnice modelu je $Y = 2585,5 + 1,2x_2 - 531x_4$; a oba regresní koeficienty b_2, b_4 jsou statisticky významné, tzn. nenulové, a tedy proměnné x_2 a x_4 přispívají k vysvětlení nezávislé proměnné y (počet odpracovaných hodin v nemocnici). Koeficient determinace je vysoký (0,98), a tato hodnota říká, že 98% celkové variability je vysvětleno modelem. Také P -hodnota (F) = $1,91 \cdot 10^{-13}$ je menší než zvolená hladina významnosti $\alpha = 0,05$; proto nulovou hypotézu o nulovosti všech regresních koeficientů zamítáme, jinými slovy můžeme tvrdit, že model jako celek je zvolen správně.

	koeficient	směr. chyba	t-podíl	p-hodnota
const	2585,52	807,093	3,203	0,0064 ***
X2luzka	1,23243	0,0504440	24,43	7,02e-013 ***
X4pobyt	-530,933	156,249	-3,398	0,0043 ***

Střední hodnota závisle proměnné	4978,480
Sm. odchylka závisle proměnné	5560,534
Součet čtverců reziduí	7542086
Sm. chyba regrese	733,9758
Koeficient determinace	0,984755
Adjustovaný koeficient determinace	0,982577
F(2, 14)	452,1552
P-hodnota (F)	1,91e-13
Logaritmus věrohodnosti	-134,6457
Akaikovo kritérium	275,2914
Schwarzovo kritérium	277,7911
Hannan-Quinnovo kritérium	275,5399

zde je poznámka o zkratkách statistik modelu

Obrázek 33: Odhad nového (redukovaného) modelu

- d) Předpověď počtu odpracovaných hodin v nemocnici s průměrným denním počtem pacientů 150, s 5000 obsazenými lůžky za měsíc, se 100 tis. obyvateli ve spádové oblasti a s průměrnou délkou pobytu 6 dní.

$$\text{Původní model: } Y = 2789,6 - 26,8x_1 + 2,1x_2 - 1,5x_3 - 561,2x_4$$

Dosadíme hodnoty a dostáváme:

$$Y = 2789,6 - 26,8 \cdot 150 + 2,1 \cdot 5000 - 1,5 \cdot 100 - 561,2 \cdot 6 = 5752,4 \text{ hodin.}$$

$$\text{Redukovaný model: } Y = 2585,5 + 1,2x_2 - 531x_4$$

Dosadíme hodnoty a dostáváme:

$$Y = 2585,5 + 1,2 \cdot 5000 - 531 \cdot 6 = 5399,5 \text{ hodin.}$$

Nyní se podívejme na znaménka u jednotlivých regresních koeficientů, která říkají, že s rostoucím počtem lůžek roste počet odpracovaných hodin v dané nemocnici, a s rostoucím počtem dnů, které pacienti stráví v nemocnici počet odpracovaných hodin v dané nemocnici klesá.

- e) Diagnostika modelu: Testy → Heteroskedasticita → Whiteův test

Vyhodnocení testu heteroskedasticity provedeme na základě vypočtené p -hodnoty. Testuje se nulová hypotéza H_0 : homoskedasticita reziduí (tj. konstantní rozptyl reziduí), oproti alternativní hypotéze H_1 : heteroskedasticita reziduí. P -hodnota = 0,135 je větší než zvolené $\alpha = 0,05$; proto H_0 nelze zamítnout, nebylo tedy prokázáno, že by rezidua neměla konstantní rozptyl.

	koeficient	směr. chyba	t-podíl	p-hodnota
const	-1,58618e+06	5,74515e+06	-0,2761	0,7876
X2luzka	481,042	819,263	0,5872	0,5689
X4pobyt	548599	2,29624e+06	0,2389	0,8156
sq_X2luzka	-0,0426439	0,0151866	-2,808	0,0170 **
X2_X3	42,9138	126,293	0,3398	0,7404
sq_X4pobyt	-77643,6	220241	-0,3525	0,7311

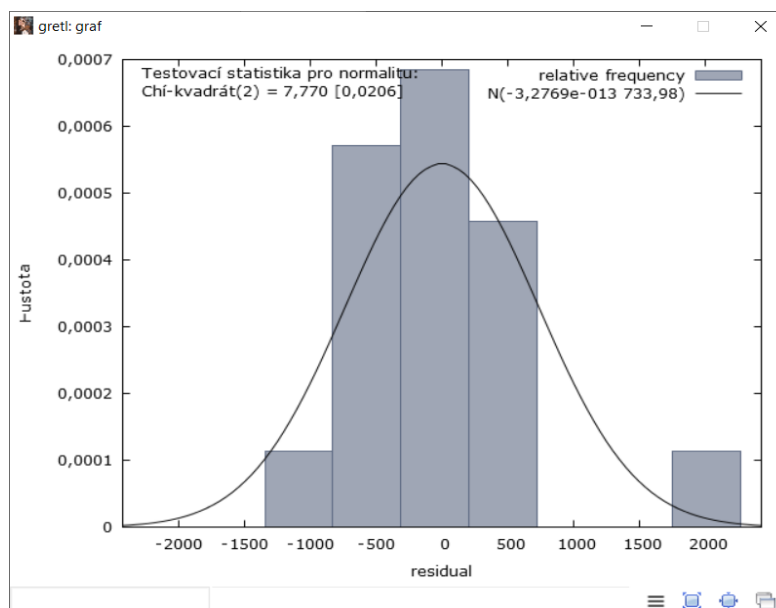
Neadjustovaný koeficient determinace = 0,493699

Testovací statistika: $TR^2 = 8,392881$,
s p-hodnotou = $P(\text{Chi-kvadrát}(5) > 8,392881) = 0,135871$

Obrázek 34: Test heteroskedasticity

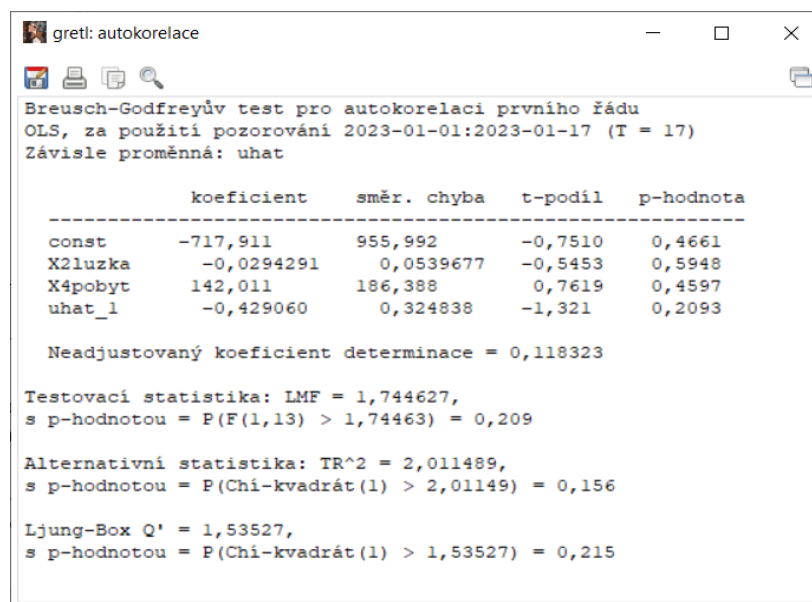
Pro testování normality vybereme ve výstupu modelu TESTY → Normalita reziduí. A dostaneme výsledek, který zachycuje Obrázek 35. Vyhodnocení provedeného testu normality je pravděpodobně nejsnazší odvodit z průběhu grafu předpokládaného normálního rozdělení v porovnání se skutečným rozdělením reziduí a analýzou p -hodnoty Chíkvadrát testu. Testuje se nulová hypotéza H_0 : Rezidua mají normální rozdělení, oproti

H_1 : Rezidua nemají normální rozdělení. P -hodnota = 0,02 je menší než zvolené $\alpha = 0,05$; proto H_0 zamítáme, a bylo tedy prokázáno, že rezidua nemají normální rozdělení.



Obrázek 35: Test normality

Pokud chceme pomocí programu GRETl testovat autokorelaci, musíme vstupní data uložit jako časovou řadu. Testuje se, zda je u_t závislé na u_{t-1} . Vybereme ve výstupu modelu záložku TESTY→Autokorelace. A dostaneme výsledek, který zachycuje Obrázek 36.



Obrázek 36: Test autokorelace

Testuje se nulová hypotéza H_0 : Rezidua nejsou autokorelována, oproti H_1 : Rezidua jsou autokorelována. P -hodnota = 0,209 je větší než zvolené $\alpha = 0,05$; proto H_0 nelze zamítnout, nebylo tedy prokázáno, že by rezidua byla autokorelována.



SAMOSTATNÉ ÚKOLY

6.1 V následující tabulce jsou uvedeny hodnoty obratu, výdajů na vědu a výzkum (VaV) a zisku za 18 průmyslových odvětví v USA v roce 2023. Vytvořte lineární regresní model závislosti zisku na obratu a výdajích na VaV. Zjistěte, zda je v modelu přítomna multikolinearita a heteroskedasticita. Použijte postupy, které jste se naučili v této kapitole.

Obrat	VaV	Zisk
6375,3	62,5	185,1
11626,4	92,9	1569,5
14655,1	178,3	276,8
21869,2	258,4	2828,1
26408,3	494,7	225,9
32405,6	1083,0	3751,9
35107,7	1620,6	2884,1
40295,4	421,7	4645,7
70761,6	509,2	5036,4
80552,8	6620,1	13869,9
95294,0	3918,6	4487,8
101314,1	1595,3	10278,9
116141,3	6107,5	8787,3
122315,7	4454,1	16438,8
141649,9	3163,8	9761,4
175025,8	13210,7	19774,5
230614,5	1703,8	22626,6
293543,0	9528,2	18415,4



ODPOVĚDI

6.1 $Y = 791,54 + 0,069 \cdot x_1 + 0,369 \cdot x_2$

x_1 ...obrat; x_2 ...výdaje na VaV; koeficient $b_2 = 0,369$ není statisticky významný

Korelační koeficient = 0,9 je statisticky významný na hladině významnosti 0,01. V modelu je přítomna multikolinearita.

Závislost zisku na obratu:

$$Y = 862,85 + 0,08 \cdot x_1$$

Koeficient 0,08 je statisticky významný.

Dvouvýběrový F-test pro rozptyl

	Soubor 1	Soubor 2
Stř. hodnota	-809,8808	809,8807591
Rozptyl	1219536	20761396,39
Pozorování	9	9
Rozdíl	8	8
F	0,058741	
P(F<=f) (1)	0,000289	
F krit (1)	0,290858	

Nulovou hypotézu: rozptyly obou souborů jsou stejné můžeme zamítnout, a lze tvrdit, že rozptyl náhodné složky není konstantní neboli heteroskedasticita je v modelu přítomna.

Závislost zisku na VaV je dána vztahem

$Y = 3817,11 + 1,4 \cdot x_2$ a koeficient 1,4 je statisticky významný.

Dvouvýběrový F-test pro rozptyl

	<i>Soubor 1</i>	<i>Soubor 2</i>
Stř. hodnota	-1348,771	1348,770762
Rozptyl	7292620	43919891,06
Pozorování	9	9
Rozdíl	8	8
F	0,166044	
P(F<=f) (1)	0,010033	
F krit (1)	0,290858	

Nulovou hypotézu: rozptyly obou souborů jsou stejné nelze zamítnout, a proto rozptyl náhodné složky je konstantní neboli heteroskedasticita není v modelu přítomna.

SHRNUTÍ KAPITOLY



Tato kapitola se věnovala identifikaci a analýze problémů, které způsobuje nesplnění hlavních předpokladů klasického vícerozměrného lineárního regresního modelu. Jednalo se o multikolinearitu, heteroskedasticitu a autokorelaci. Multikolinearitou rozumíme vzájemnou statistickou závislost, tj. korelaci, mezi vysvětlujícími proměnnými ve vícenásobném lineárním regresním modelu. Další důležitou vlastností klasického lineárního regresního modelu je homoskedasticita, která spočívá v tom, že rozptyl poruchy ε_i v populačním lineárním regresním modelu je konstantní. Autokorelace je korelace mezi pozorováními uspořádanými v čase, (data jsou časové řady) nebo v prostoru (data jsou průřezová, tj. v jednom časovém okamžiku/intervalu).

7 ZÁKLADY ANALÝZY ČASOVÝCH ŘAD



RYCHLÝ NÁHLED KAPITOLY

Důležitým nástrojem ke zkoumání dynamiky ekonomických procesů je analýza časových řad. Časovou řadou přitom rozumíme věcně a prostorově srovnatelná pozorování uspořádaná v čase směrem od minulosti přes přítomnost k budoucnosti. Obsahem této kapitoly je objasnit typizaci ekonomických časových řad, vysvětlit elementární charakteristiky časových řad, uvést základní modely časových řad a popsat jejich složky. Analýza časových řad je vedena snahou po vysvětlení minulosti a předvídání budoucnosti, v ekonomické oblasti se jedná o vývojové trendy ukazatelů hospodářské činnosti. Analýza časových řad jako soubor metod a postupů nabízí širokou škálu nástrojů a technik. Ke klasickým analytickým postupům založeným na regresi z předchozích kapitol a syntetickým přístupům založeným na technikách vyrovnání časových řad, přistupuje moderní, výpočetně náročnější Box-Jenkinsova metodologie.



CÍLE KAPITOLY

Po prostudování této kapitoly budete umět:

- uvést typy ekonomických časových řad,
- vypočítat hodnoty očištěné časové řady,
- vypočítat základní charakteristiky časových řad.



ČAS POTŘEBNÝ KE STUDIU

K prostudování této kapitoly budete potřebovat asi 60 minut.



KLÍČOVÁ SLOVA KAPITOLY

Časová řada, difference časové řady, koeficient růstu, očištěná hodnota časové řady.

7.1 Typy ekonomických časových řad

Důležitým nástrojem ke zkoumání dynamiky ekonomických procesů je analýza časových řad. Časovou řadou přitom rozumíme věcně a prostorově srovnatelná pozorování uspořádaná v čase směrem od minulosti přes přítomnost k budoucnosti. Časové řady členíme následujícím způsobem:

- podle charakteru časové řady na *intervalové časové řady* a *okamžikové časové řady*,
- podle periodicity, s jakou jsou sledovány, na *krátkodobé* časové řady (méně než roční periodičita), *střednědobé* časové řady (roční periodičita) a *dlouhodobé* časové řady (delší, než roční periodičita),
- podle druhu sledovaných ukazatelů (údajů) na časové řady *absolutních* ukazatelů a časové řady *odvozených* ukazatelů.

Intervalovou časovou řadou se rozumí časová řada intervalového ukazatele y_t , tj. ukazatele, jehož velikost (hodnota) závisí na délce intervalu, za který je sledován. Pro ukazatele tohoto typu je možné tvořit součty, z jejich povahy však vyplývá, že se vztahují ke stejně dlouhým časovým intervalům, jinak by byly hodnoty vzájemně nesrovnatelné. Není např. správné srovnávat výrobu za leden a únor, neboť únor je z hlediska počtu pracovních dní kratší. Abychom zajistili srovnatelnost, přepočítáváme všechna sledovaná období na stejný časový interval. Tato operace se nazývá očišťování časových řad od kalendářních variací. Údaje očištěné časové řady $y_t^{(0)}$ dostaneme z hodnoty očišťovaného ukazatele y_t takto:

$$y_t^{(0)} = y_t \frac{\bar{k}_t}{k_t}, \quad (7.1)$$

kde \bar{k}_t je průměrný počet dnů v příslušném dílčím období, k_t je skutečný počet dnů v příslušném dílčím období t .

Okamžikovou časovou řadou rozumíme časovou řadu ukazatelů, které se vztahují k určitému okamžiku, např. počátku nebo konci určitého časového intervalu (období). Protože součet za několik za sebou jdoucích okamžikových hodnot obvykle nemá reálný smysl, shrnují se řady tohoto typu pomocí chronologického průměru.

Pro dané *ekvidistantní* (stejně vzdálené) časové okamžiky t_1, t_2, \dots, t_n , ke kterým přísluší hodnoty okamžikových ukazatelů y_1, y_2, \dots, y_n je *prostý chronologický průměr* definován jako aritmetický průměr z aritmetických průměrů vždy dvou po sobě jdoucích hodnot, tedy:

$$\bar{y}_{ch} = \frac{\frac{y_1+y_2}{2} + \frac{y_2+y_3}{2} + \dots + \frac{y_{n-1}+y_n}{2}}{n-1} \quad (7.2)$$

Není-li délka mezi jednotlivými časovými okamžiky stejná, definujeme *vážený chronologický průměr*, kde vahami jsou délky jednotlivých časových intervalů $d_k = t_{k+1} - t_k$, $k = 1, 2, \dots, n-1$:

$$\bar{y}_{ch} = \frac{\frac{y_1 + y_2}{2}d_1 + \frac{y_2 + y_3}{2}d_2 + \dots + \frac{y_{n-1} + y_n}{2}d_{n-1}}{d_1 + d_2 + \dots + d_{n-1}} \quad (7.3)$$

Časový rozdíl mezi časovými okamžiky, tedy délka časového intervalu v okamžikové časové řadě, se nazývá *periodicita* časové řady. Je-li periodicita ekonomických časových řad kratší než jeden rok, hovoříme o *krátkodobých časových řadách*. Nejčastější periodicitou je *měsíční* periodicita. Je-li periodicita *roční*, hovoříme často o *střednědobých časových řadách*, při delší periodicitě, např. pětileté, hovoříme o *dlouhodobých časových řadách*.

Časovou řadou absolutních hodnot se obvykle rozumí časová řada přímo zjištěných údajů (v naturálních jednotkách) očištěná od kalendářních variací. Odvozené údaje a z nich vytvořené časové řady získáme obvykle matematickými operacemi z absolutních údajů. Většinu důležitých ekonomických časových řad tvoří časové řady ukazatelů vyjádřených v peněžní formě. Vzhledem ke změnám cenové hladiny, které jsou v tržní ekonomice přirozené, však v delší časové řadě často dostáváme posloupnost údajů, které nejsou vždy zcela souměřitelné. Proto důležitým problémem v analýze časových řad je srovnatelnost údajů, konkrétně cenová srovnatelnost. Při sestavování delší časové řady je možno v zásadě postupovat dvojím způsobem: použít *běžné ceny* a vyjádřit z nich absolutní objem určitého ukazatele, resp. tempa růstu, nebo vycházet ze *stálých cen*, tj. cen fixovaných k určitému datu. Používání stálých cen v ekonomice vede ke zmírnění negativních tendencí v účinnosti základních fondů vyplývajících z vlivu technického rozvoje na výrobu, dále vede ke zrealnění výsledků hospodářského vývoje vzhledem k mezinárodnímu srovnání.

Vývoj základních ekonomických ukazatelů v České republice je možné sledovat na webových stránkách Českého statistického úřadu. Pro potřeby vrcholového řízení ve firmách a podnicích slouží především údaje o vývoji základních ukazatelů podle měsíců, neboť jde o informace s určitým vztahem k okamžité odezvě v chování ekonomických subjektů, ať už výrobců, nebo spotřebitelů. Jsou to zejména informace o inflaci (index spotřebitelských cen a indexy životních nákladů), dále informace o peněžních příjmech a výdajích obyvatelstva, o celkovém prodeji v maloobchodě, průmyslové, zemědělské a stavební výrobě a též údaje o nezaměstnanosti.

Zdrojem informací a dat jsou webové stránky Českého statistického úřadu (ČSÚ), www.czso.cz případně Statistického úřadu Evropské komise EUROSTAT: <http://epp.eurostat.ec.europa.eu>, Česká národní banka <https://www.cnb.cz/cs/>, The world bank <https://databank.worldbank.org/source/world-development-indicators#>.

7.2 Elementární charakteristiky časových řad

Mezi elementární metody analýzy časových řad patří *vizuální analýza* chování ukazatele využívající grafů spolu s určováním elementárních statistických charakteristik, ke kterým patří *absolutní diference* různého řádu a *koeficient růstu* časové řady.

Označíme-li y_t hodnoty určitého ukazatele v čase $t = 1, 2, \dots, n$ (např. v jednotlivých měsících), potom *absolutní diferencí prvního řádu* rozumíme rozdíl:

$$\Delta^{(1)}y_t = y_t - y_{t-1}, t = 2, 3, \dots, n. \quad (7.4)$$

Obdobně lze definovat *absolutní diference vyšších řádů*:

$$\Delta^{(2)}y_t = \Delta^{(1)}y_t - \Delta^{(1)}y_{t-1} = y_t - 2y_{t-1} + y_{t-2}, t = 3, 4, \dots, n,$$

$$\Delta^{(3)}y_t = \Delta^{(2)}y_t - \Delta^{(2)}y_{t-1} = y_t - 3y_{t-1} + 3y_{t-2} - y_{t-3}, t = 4, 5, \dots, n$$

Další používanou elementární charakteristikou je *koeficient růstu*, který udává, o kolik procent vzrostla hodnota časové řady v daném časovém okamžiku oproti období v předchozím časovém okamžiku:

$$k_t = \frac{y_t}{y_{t-1}}, t = 2, 3, \dots, n. \quad (7.5)$$

Při hodnocení vývoje za celou analyzovanou řadu zjišťujeme souhrnné charakteristiky – *průměrný absolutní přírůstek*:

$$\bar{\Delta} = \frac{1}{n-1} \sum_{t=2}^n \Delta^{(1)}y = \frac{y_n - y_1}{n-1} \quad (7.6)$$

a průměrný koeficient růstu:

$$\bar{k} = \sqrt[n-1]{k_2 k_3 \dots k_n} = \sqrt[n-1]{\frac{y_n}{y_1}}. \quad (7.7)$$

Jak průměrný absolutní přírůstek, tak průměrný koeficient růstu závisí pouze na první a poslední hodnotě časové řady. Průměrný absolutní přírůstek ukazuje, o kolik by se měl ukazatel pravidelně měnit (v absolutních jednotkách), aby se hodnota ukazatele změnila z původní první hodnoty y_1 na poslední hodnotu y_n . Naproti tomu průměrný koeficient růstu poskytuje informaci, o kolik procent by se měla hodnota ukazatele měnit, tj. jaká by měla být rychlost růstu (poklesu), aby se hodnota ukazatele změnila z původní první hodnoty y_1 na poslední hodnotu y_n .

7.3 Modely ekonomických časových řad

Modelový přístup k analýze časových řad bude vycházet z předpokladu, že *jediným faktorem dynamiky ukazatele v časové řadě je čas*. Ostatní faktory působící na hodnotu ukazatele budeme většinou zanedbávat. Model časové řady tohoto typu můžeme zapsat ve formě:

$$y_t = f(t, \varepsilon_t), \quad (7.8)$$

kde y_t je hodnota analyzovaného ukazatele v čase t , f je určitá funkce (typ závislosti), t je časová proměnná, ε_t je hodnota náhodné složky. Modely časových řad založené na výše uvedeném principu se nazývají *jednorozměrné modely*.

Každá časová řada může obsahovat 4 složky, které vyjadřují různé druhy pohybu analyzovaného ukazatele:

- trendovou složku (trend) T_t ,
- sezónní složku S_t ,
- cyklickou složku C_t ,
- náhodnou složku ε_t .
-

Trendová, sezónní a cyklická složka tvoří společně *systematickou (deterministickou) složku*, kterou značíme Y_t , tj. $Y_t = T_t + S_t + C_t$. Zpravidla se uvažuje, že složky Y_t jsou v aditivním vztahu, takže model časové řady můžeme zapsat ve tvaru:

$$y_t = T_t + S_t + C_t + \varepsilon_t. \quad (7.9)$$

V tom případě mluvíme o *aditivním modelu* časové řady. V ekonomických časových řadách se nejčastěji setkáváme se dvěma speciálními případy modelu (7.9). U střednědobých modelů (s roční periodicitou) se obvykle předpokládá $S_t = C_t = 0$, pak model časové řady (7.9) má tvar:

$$y_t = T_t + \varepsilon_t. \quad (7.10)$$

U krátkodobých modelů časových řad (s čtvrtletní nebo měsíční periodicitou) se předpokládá, že $C_t = 0$, a tedy model (7.9) má tvar:

$$y_t = T_t + S_t + \varepsilon_t, \quad (7.11)$$

mluvíme pak o *časové řadě se sezónní složkou*.

Vedle aditivního modelu (8.9) je *multiplikativní model* založen na předpokladu, že vzájemný vztah jednotlivých složek obsažených v modelu je dán vzájemným násobením:

$$y_t = T_t \cdot S_t \cdot C_t \cdot \varepsilon_t. \quad (7.12)$$

Popis a kvantifikace jednotlivých složek modelu časové řady patří k hlavním úkolům analýzy časových řad.

ŘEŠENÁ ÚLOHA 7.1

V tabulce jsou uvedeny průměrné měsíční výdaje na vzdělávání zaměstnanců ve firmě A+B v letech 2015-2023. Pro tuto časovou řadu vypočítejte:

- absolutní přírůstky a průměrný absolutní přírůstek,
- koeficienty růstu a průměrný koeficient růstu.

Roky	2015	2016	2017	2018	2019	2020	2021	2022	2023
Mzda	2980	3110	4500	5650	7460	8930	10670	12820	13250

Řešení:

- Absolutní přírůstky vypočítáme podle vztahu (7.4):

$$\Delta^{(1)}y_2 = y_2 - y_1 = 3110 - 2980 = 130, \text{ atd.}$$

Výsledek říká, že průměrné měsíční výdaje na vzdělávání zaměstnanců ve firmě A+B stouply v letech 2015-2016 o 130 Kč.

Všechny absolutní přírůstky jsou uvedeny v následující tabulce.

Průměrný absolutní přírůstek je podle (7.6):

$$\bar{\Delta} = \frac{y_n - y_1}{n-1} = \frac{13250 - 2980}{8} = 1283,75.$$

- Koeficienty růstu vypočítáme podle vztahu (7.5). Např.: $k_2 = \frac{y_2}{y_1} = \frac{3110}{2980} = 1,0436.$

Průměrné měsíční výdaje na vzdělávání zaměstnanců ve firmě A+B vzrostly v letech 2015-2016 o 4,36%.

Hodnoty ostatních koeficientů růstu jsou uvedeny v následující tabulce.

Průměrný koeficient růstu vypočítáme podle (7.7):

$$\bar{k} = \sqrt[n-1]{\frac{y_n}{y_1}} = \sqrt[8]{\frac{13250}{2980}} = 1,205.$$

Výsledek ukazuje, že měsíční výdaje na vzdělávání zaměstnanců ve firmě A+B rostly ročně v průměru o 20,5%.

Roky	2015	2016	2017	2018	2019	2020	2021	2022	2023
Mzda	2980	3110	4500	5650	7460	8930	10670	12820	13250
$\Delta^{(1)}y$.	130	1390	1150	1810	1470	1740	2150	430
k	.	1,04	1,45	1,26	1,32	1,20	1,19	1,20	1,03

SAMOSTATNÉ ÚKOLY

7.1 V tabulce jsou uvedeny počty prodaných automobilů v autocentru A+A v letech 2016 až 2023. Pro tuto časovou řadu vypočítejte:

- absolutní přírůstky a průměrný absolutní přírůstek
- koeficienty růstu a průměrný koeficient růstu.

Rok	2016	2017	2018	2019	2020	2021	2022	2023
Počet	120	159	167	175	197	172	199	240

7.2 Uvedené údaje v tabulce zachycují zisk firmy v tis. Kč v letech 2017-2023. Pro tuto časovou řadu vypočítejte:

- absolutní přírůstky a průměrný absolutní přírůstek
- koeficienty růstu a průměrný koeficient růstu.

Rok	2017	2018	2019	2020	2021	2022	2023
Počet	1303,6	1381,1	1447,7	1432,8	1401,3	1390,6	1433,8



ODPOVĚDI

7.1

Rok	Počet	Abs.přírůstky	Koeficienty růstu
2016	120	xxx	xxx
2017	159	39	1,325
2018	167	8	1,050
2019	175	8	1,048
2020	197	22	1,126
2021	172	-25	0,873
2022	199	27	1,157
2023	240	41	1,206

Průměrný absolutní přírůstek je podle (7.6): $\bar{\Delta} = 17,14$.

Průměrný koeficient růstu vypočítáme podle (7.7): $\bar{k} = 1,104$.

Počet prodaných automobilů rostl ročně v průměru o 10,4%.

7.2

Rok	Počet	Abs.přírůstky	Koeficienty růstu
2017	1303,6	xxx	xxx
2018	1381,1	77,5	1,059
2019	1447,7	66,6	1,048
2020	1432,8	-14,9	0,990
2021	1401,3	-31,5	0,978
2022	1390,6	-10,7	0,992
2023	1433,8	43,2	1,031

Průměrný absolutní přírůstek je podle (7.6): $\bar{\Delta} = 21,7$ tis. Kč.

Průměrný koeficient růstu vypočítáme podle (7.7): $\bar{k} = 1,016$.

Zisk firmy rostl ročně v průměru o 1,6%.

SHRNUTÍ KAPITOLY



Obsahem této kapitoly bylo objasnit typizaci ekonomických časových řad, vysvětlit elementární charakteristiky časových řad, uvést základní modely časových řad a popsat jejich složky. Časová řada se dá rozložit na čtyři složky. Jedná se o složku trendovou, sezónní, cyklickou a náhodnou. Cyklickou složku v ekonomických časových řadách zanedbáváme, protože popisuje jevy, které se opakují za období delší než 1 rok. V případě, že se jednotlivé složky sčítají, tak se jedná o aditivní model, v případě násobení jednotlivých složek mluvíme o multiplikačním modelu. Analýza časových řad je vedena snahou po vysvětlení minulosti a předvídání budoucnosti, v ekonomické oblasti se jedná o vývojové trendy ukazatelů hospodářské činnosti.

8 ANALÝZA TRENDU ČASOVÝCH ŘAD



RYCHLÝ NÁHLED KAPITOLY

V této kapitole se budete zabývat trendovou složkou časové řady, která představuje nejdůležitější komponentu analyzované časové řady. Proto popis trendu je jedním z nejdůležitějších úkolů analýzy časových řad. Vycházíme přitom z předpokladu, že jediným faktorem vývoje dynamiky analyzovaného ukazatele je čas. Trendová složka totiž poskytuje rozhodující informaci pro prognózování hodnot časové řady do budoucna. K určení trendové složky používáme dva obecné přístupy: analytický a syntetický. Analytický přístup stanovení trendu vychází z předem známých typů trendových funkcí vyznačujících se přítomností parametrů, které je třeba stanovit co nejlépe s ohledem na skutečné hodnoty ukazatele časové řady. Z velkého množství používaných trendových funkcí se zaměříme na několik z nich, které mají význam především v ekonomických aplikacích. Jsou to: lineární trend, parabolický trend, exponenciální trend, logistický trend a Gompertzův trend. Syntetický přístup stanovení trendu spočívá ve vyrovnání odchylek daného ukazatele v časové řadě tak, že získané vyrovnané hodnoty vyjadřují trendový faktor obsažený pouze v časové řadě, nikoliv faktor vložený z vnějšku. Nemusíte proto znát předem typ trendové funkce, což je přednost syntetického přístupu oproti přístupu analytickému. Jeho nevýhodou je naopak obtížnější využití pro prognózování hodnot časové řady. Z existujících metod syntetického přístupu uvedeme metody klouzavého průměru a exponenciální vyrovnání.



CÍLE KAPITOLY

Po prostudování této kapitoly budete umět:

- uvést přístupy používané k určení trendové složky,
 - napsat lineární, kvadratickou, exponenciální a logaritmickou trendovou funkci,
 - vztahy pro výpočet odhadů parametrů lineární trendové funkce,
 - vypočítat koeficient determinace,
 - vyrovnat časovou řadu klouzavými průměry,
 - použít pro vyrovnání časové řady exponenciální vyrovnání.
-

ČAS POTŘEBNÝ KE STUDIU

K prostudování této kapitoly budete potřebovat asi 120 minut.

KLÍČOVÁ SLOVA KAPITOLY

Trendová složka, lineární trendová funkce, koeficient determinace, klouzavé průměry, koeficient korelace.

8.1 Trendová složka časových řad

Jak již bylo v průvodci studiem řečeno, v této kapitole vycházíme z předpokladu, že jediným faktorem vývoje dynamiky analyzovaného ukazatele je čas t . Jednoduchý způsob volby časové proměnné spočívá v jejím zavedení tak, že časová řada začíná v okamžiku 1, ke kterému se vztahuje první člen analyzované časové řady y_1 . Další časové okamžiky označujeme po řadě přirozenými čísly 2, 3, ..., n . Symbol n označuje poslední uvažovaný časový okamžik a zároveň i počet uvažovaných časových okamžiků.

Jiný jednoduchý a výhodný způsob označení časové proměnné spočívá v zavedení nové časové proměnné t' následujícím způsobem:

$$t' = (t - \bar{t}), \quad (8.1)$$

je-li počet členů časové řady n lichý, pak $\bar{t} = \frac{n+1}{2}$, jak ukazuje Tabulka 12,

$$\text{nebo } t' = 2(t - \bar{t}), \quad (8.2)$$

je-li počet členů n sudý, jak ukazuje Tabulka 13. Nová časová proměnná splňuje důležitý požadavek: $\sum_{t=1}^n t' = 0$. (8.3)

Tabulka 12: Transformovaná proměnná při lichém časová n

Rok	2017	2018	2019	2020	2021	2022	2023
t	1	2	3	4	5	6	7
t'	-3	-2	-1	0	1	2	3

Tabulka 13: Transformovaná časová proměnná při sudém n

Rok	2018	2019	2020	2021	2022	2023
t	1	2	3	4	5	6
t'	-5	-3	-1	1	3	5

Dále uvedené vztahy pro výpočet odhadů teoretických hodnot parametrů jsou uváděny po zavedení transformací v Tabulkách 12 a 13.

Trendová složka představuje nejdůležitější komponentu analyzované časové řady, a proto popis trendu je jedním z nejdůležitějších úkolů analýzy časových řad. Trendová složka totiž poskytuje rozhodující informaci pro prognózování hodnot časové řady do budoucna. K určení trendové složky používáme dva obecné přístupy: analytický a syntetický.

Analytický přístup stanovení trendu vychází z předem známých typů trendových funkcí vyznačujících se přítomností parametrů, které je třeba stanovit co nejlépe s ohledem na skutečné hodnoty ukazatele časové řady.

Syntetický přístup stanovení trendu spočívá ve vyrovnání odchylek daného ukazatele v časové řadě (tzv. vyrovnání) tak, že získané vyrovnané hodnoty vyjadřují trendový faktor obsažený pouze v časové řadě, nikoliv faktor vložený z vnějšku. Nemusíme proto znát předem typ trendové funkce, což je přednost syntetického přístupu oproti přístupu analytickému. Jeho nevýhodou je naopak obtížnější využití pro prognózování hodnot časové řady. Z existujících metod syntetického přístupu uvedeme metody klouzavého průměru a exponenciální vyrovnání.

8.2 Trendové funkce

Z velkého množství používaných trendových funkcí se zaměříme na několik z nich, které mají význam především v ekonomických aplikacích. Jsou to: lineární trend, parabolický trend, exponenciální trend, logistický trend a Gompertzův trend. Výhodou těchto trendových funkcí je to, že je lze snadno použít pro účely prognózování. Nevýhodou je fakt, že typ trendové funkce musíme stanovit předem na základě externích, mnohdy subjektivních předpokladů a informací. Nejužívanější metodou odhadu neznámých parametrů trendové funkce je *metoda nejmenších čtverců (MNC)*, s níž jsme se setkali již v kapitole 3. Zde tuto metodu aplikujeme na speciální typ jednoduché regrese pro data ve formě ekonomické časové řady, tedy případ, kdy nezávisle proměnnou je čas a závisle proměnnou tvoří sledovaný ekonomický ukazatel. Kromě metody nejmenších čtverců pro nelineární trendové funkce uvedeme alternativní *metodu vybraných bodů (MVB)*.

8.2.1 LINEÁRNÍ TREND

Nejčastěji používanou trendovou funkcí je lineární trendová funkce:

$$T_t = \beta_0 + \beta_1 t, \quad (8.4)$$

kde β_0, β_1 jsou neznámé parametry a $t = 1, 2, \dots, n$ je časová proměnná. Odhady neznámých parametrů, které označujeme b_0, b_1 , získáme metodou nejmenších čtverců, která dává nejlepší nestranné odhady. V souladu s postupem z kapitoly 3 je zapotřebí vyřešit 2 normální rovnice (3.12), kde x_i nahradíme t :

$$\sum y_t = b_0 n + b_1 \sum t, \quad (8.5)$$

$$\sum ty_t = b_0 \sum t + b_1 \sum t^2. \quad (8.6)$$

Použijeme-li nyní časové transformace (8.1), (8.2) a s využitím vztahu (8.3) dostaneme jednoduché řešení normálních rovnic (8.5), (8.6):

$$b_0 = \frac{\sum y_t}{n}, \quad b_1 = \frac{\sum t'y_t}{\sum (t')^2}. \quad (8.7)$$

Parametr b_0 interpretujeme jako aritmetický průměr hodnot časové řady, parametr b_1 udává, jaký přírůstek hodnoty T_t odpovídá jednotkovému přírůstku proměnné t .

ŘEŠENÁ ÚLOHA 8.1



V následující tabulce jsou uvedeny počty prodaných automobilů v autocentru A+A v letech 2016 až 2023. Pro tuto časovou řadu vypočítejte:

Rok	2016	2017	2018	2019	2020	2021	2022	2023
Počet	120	159	167	175	197	172	199	240

- Trend v prodeji automobilů popište lineární trendovou funkcí.
- Jaký počet prodaných automobilů lze očekávat v roce 2024 s 95 % pravděpodobností? (Stanovte bodový odhad a 95 %-ní interval spolehlivosti prognózy.)
- Stanovte koeficient determinace a na jeho základě určete přílehlavost dat k trendové funkci.

Řešení:

- Podle vztahu (8.2) zavedeme novou časovou proměnnou t' (viz následující tabulka).

Rok	t'	y_t	t'^2	$y_t t'$	\hat{T}	$(y - \hat{T})^2$	$(y - \bar{y})^2$
2016	-7	120	49	-840	133,818	190,937	3436,891
2017	-5	159	25	-795	146,620	153,264	385,141
2018	-3	167	9	-501	159,422	57,426	135,141
2019	-1	175	1	-175	172,224	7,706	13,141
2020	1	197	1	197	185,026	143,377	337,641
2021	3	172	9	516	197,828	667,086	43,891
2022	5	199	25	995	210,630	135,257	415,141
2023	7	240	49	1680	223,432	274,499	3766,891
Součet	0	1429	168	1077		1629,552	8533,875

Odhady b_0, b_1 parametrů β_0, β_1 trendové funkce:

$$T_t = \beta_0 + \beta_1 t', t' = -7, -5, -3, \dots$$

vypočítáme podle vztahů:

$$b_0 = \frac{\sum y_t}{n} = \frac{1429}{8} = 178,625, \quad b_1 = \frac{\sum t' y_t}{\sum t'^2} = \frac{1077}{168} = 6,410.$$

Odhadnutá trendová funkce má tvar:

$$\hat{T} = 178,625 + 6,41 t', t' = -7, -5, -3, \dots$$

b. Očekávaný prodej v roce 2024 vypočítáme dosazením t' , které odpovídá roku 2024, do rovnice trendu:

$$\hat{T} = 178,625 + 6,401 \cdot 9 \cong 236,32.$$

Intervalovou předpověď obdržíme dosazením potřebných hodnot do vztahu (4.8). Ve speciálním případě časové řady, kdy $t_i = x_i$, obdržíme po úpravách následující vztah pro interval spolehlivosti predikce na i časových okamžiků předem:

$$[y(n+i) - t_{1-\alpha/2}(n-2) s_R \sqrt{Q_n(i)}, y(n+i) + t_{1-\alpha/2}(n-2) s_R \sqrt{Q_n(i)}],$$

kde $y(n+i) = \hat{T} = 236,32$

$$t_{1-\alpha/2}(n-2) = 2,45$$

$$s_R = \sqrt{\frac{S_R}{n-p}} \quad Q_n(i) = \sqrt{(1-R^2) \frac{n(n^2-1)+12i^2}{(n^2-1)(n-2)}}, i = 1.$$

Z tabulky obdržíte $S_R = 1629,552$. Potom směrodatná chyba odhadu s_R je

$$s_R = \sqrt{\frac{1629,552}{8-2}} = 16,48.$$

K výpočtu $Q_n(i)$ je zapotřebí znát hodnotu koeficientu determinace R^2

$$R^2 = 1 - \frac{S_R}{S_y} = 1 - \frac{1629,552}{8533,875} = 0,809.$$

Výpočet součtu S_y je uveden v tabulce. Potom

$$Q_n(i) = \sqrt{(1-0,809) \frac{8(64-1)+12}{(64-1)(8-2)}} = \sqrt{0,191 \cdot \frac{516}{378}} = 0,51.$$

Dosazením výše vypočítaných hodnot do obecného vztahu obdržíte levou (L) a pravou (P) mez intervalové předpovědi.

$$L = 236,315 - 2,447 \cdot 16,48 \cdot \sqrt{0,51} = 207,52.$$

$$P = 236,315 + 2,447 \cdot 16,48 \cdot \sqrt{0,51} = 265,11.$$

Bodový odhad prodeje v roce 2018 je 236 automobilů. S 95 % pravděpodobností by se mělo v roce 2024 prodat mezi 208 a 265 automobily.

c. Koeficient determinace byl vypočten v **b**: $R^2 = 0,809$. Tato hodnota říká, že přiléhavost dat k trendové funkci je „vysoká“.

8.2.2 KVADRATICKÝ TREND

Rozšířením lineárního trendu o kvadratický člen dostaneme *parabolickou trendovou funkci*:

$$T_t = \beta_0 + \beta_1 t + \beta_2 t^2, \quad (8.8)$$

kde $\beta_0, \beta_1, \beta_2$ jsou neznámé parametry a $t = 1, 2, \dots, n$ je časová proměnná. Odhady neznámých parametrů, které označujeme b_0, b_1, b_2 , získáme metodou nejmenších čtverců řešením soustavy 3 lineárních rovnic o 3 neznámých:

$$\begin{aligned} \sum y_t &= b_0 n + b_1 \sum t' + b_2 \sum (t')^2, \\ \sum t' y_t &= b_0 \sum t' + b_1 \sum (t')^2 + b_2 \sum (t')^3, \\ \sum (t')^2 y_t &= b_0 \sum (t')^2 + b_1 \sum (t')^3 + b_2 \sum (t')^4. \end{aligned} \quad (8.9)$$

Z podmínky (8.3) dostaneme z rovnice (8.9) ihned řešení:

$$b_1 = \frac{\sum t' y_t}{\sum (t')^2}. \quad (8.10)$$

Dosazením (8.10) do zbývajících dvou normálních rovnic obdržíme ještě řešení b_0, b_2 :

$$b_0 = \frac{\sum y_t \sum (t')^4 - \sum (t')^2 \sum y_t (t')^2}{n \sum (t')^4 - (\sum (t')^2)^2}, \quad (8.11)$$

$$b_2 = \frac{n \sum y_t (t')^2 - \sum y_t \sum (t')^2}{n \sum (t')^4 - (\sum (t')^2)^2}. \quad (8.12)$$

8.2.3 MOCNINNÝ TREND

Mocninná trendová funkce má tvar:

$$T_t = \beta_0 t^{\beta_1}, \quad (8.13)$$

avšak namísto něj uvažujeme model, jenž vznikne logaritmováním obou stran (8.13):

$$\ln T_t = \ln \beta_0 + \beta_1 \ln t,$$

kde \ln je přirozený logaritmus o základu $e = 2,718\dots$ Použijeme analogický postup jako v případě jednoduché lineární regrese v kapitole 2.2.6. Jestliže nyní použijeme substituce

$$T'_t = \ln T_t, \quad t'' = \ln t, \quad (8.14)$$

$$\beta'_0 = \ln \beta_0, \quad \beta'_1 = \beta_1, \quad (8.15)$$

obdržíme „čárkovaný“ lineární trend:

$$T'_t = \beta'_0 + \beta'_1 t'' \quad (8.16)$$

jehož parametry β'_0, β'_1 (regresní koeficienty) odhadneme metodou nejmenších čtverců a obdržíme tak jejich odhady b'_0, b'_1 . Ze vztahů (8.15) vypočteme zpětně odhady b_0, b_1 :

$$b_0 = e^{b'_0}, b_1 = b'_1.$$

8.2.4 EXPONENCIÁLNÍ TREND

Exponenciální trendová funkce má tvar:

$$T_t = \beta_0 \beta_1^t \quad (8.17)$$

který substitucemi:

$$T'_t = \ln T_t, t'' = t, \quad (8.18)$$

$$\beta'_0 = \ln \beta_0, \beta'_1 = \ln \beta_1, \quad (8.19)$$

lze rovněž transformovat na „čárkovaný“ lineární trend, jehož parametry β'_0, β'_1 odhadneme metodou nejmenších čtverců, a obdržíme tak odhady b'_0, b'_1 . Ze vztahů (8.19) vypočteme odhady b_0, b_1 původního nelineárního regresního modelu (8.17):

$$b_0 = e^{b'_0}, b_1 = e^{b'_1}.$$

Použití exponenciálního trendu je uvedeno v následující řešené úloze.



ŘEŠENÁ ÚLOHA 8.2

V tabulce jsou uvedeny údaje o počtu vyrobených myček nádobí v letech 2017-2023.

- Trend ve výrobě tohoto výrobku popište exponenciální trendovou funkcí.
- Vypočítejte bodovou prognózu výroby na rok 2024, dále zjistěte koeficient determinace a na jeho základě zhodnoťte „přiléhavost“ dat k trendové funkci.

Rok	2007	2018	2019	2020	2021	2020	2021	2022	2023
Myčky nádobí (tis. ks)	8	9	17	20	38	40	70	101	180

Řešení:

Nejprve vypočítáte odhady b_0, b_1 parametrů exponenciální trendové funkce

$$T_t = \beta_0 \beta_1^t.$$

Logaritmováním této rovnice obdržíte vztah

$$\ln T_t = \ln \beta_0 + t \ln \beta_1.$$

Zavedením substituce

$$T'_t = \ln T_t, \quad t' = t,$$

$$\beta'_0 = \ln \beta_0, \quad \beta'_1 = \ln \beta_1$$

se původní rovnice exponenciálního trendu transformuje na rovnici lineárního trendu.

Zavedete novou časovou proměnnou t'' viz (8.1) a vypočítáte koeficienty b'_0, b'_1

$$b'_0 = \frac{\sum y'_t}{n} = \frac{31,4886}{9} = 3,4987, \quad b'_1 = \frac{\sum t'' y'_t}{\sum t''^2} = \frac{23,2315}{60} = 0,3872.$$

Potom

$$b_0 = e^{b'_0} = e^{3,4987} = 33,07, \quad b_1 = e^{b'_1} = e^{0,3872} = 1,47.$$

Rok	t''	y	$y' = \ln y$	t''^2	$t'' y'$	T	$(y - T)^2$	$(y - \bar{y})^2$
2004	-4	8	2,0794	16	-8,3178	7,0285	0,8425	2085,7489
2005	-3	9	2,1972	9	-6,5917	10,3519	1,9904	1995,4089
2006	-2	17	2,8332	4	-5,6664	15,2466	2,8771	1344,6889
2007	-1	20	2,9957	1	-2,9957	22,4558	6,2330	1133,6689
2008	0	38	3,6376	0	0	33,0737	24,3049	245,5489
2009	1	40	3,6889	1	3,6889	48,7122	74,1821	186,8689
2010	2	70	4,2485	4	8,4970	71,7452	2,1345	266,6689
2011	3	101	4,6151	9	13,8453	105,6690	16,3831	2240,1289
2012	4	180	5,1930	16	20,7718	155,6333	654,3364	15959,2689
Součet	0	490	31,4886	60	23,2315		783,2839	25458,0001

Hledaná trendová funkce má tvar

$$\hat{T}_{t''} = 33,07 \cdot 1,47^{t''}, \quad t'' = -4, -3, -2, \dots$$

K bodovému odhadu využijeme nalezenou trendovou funkci, kam dosadíme $t'' = 5$, což je hodnota, která odpovídá netransformované časové hodnotě $t = 2024$.

Koeficient determinace vyžaduje znát hodnotu celkového součtu S_y a reziduálního součtu S_R (viz poslední dva sloupce v tabulce).

Pro výpočet reziduálního součtu čtverců je dále třeba znát odhady teoretické hodnoty $\hat{T}_{t''}$, které obdržíme postupným dosazováním za t'' do rovnice trendu, tedy např. pro $t'' = -4$:

$$\hat{T} = 33,07 \cdot 1,47^{-4} = 7,08.$$

Všechny hodnoty \hat{T} i součtů S_y , S_R najdete v tabulce. Pro koeficient determinace platí:

$$R^2 = 1 - \frac{S_R}{S_y} = 1 - \frac{783,2839}{25458,0001} = 0,969.$$

Hodnota 0,969 říká, že přiléhavost dat k trendové křivce je vysoká.

8.2.5 LOGISTICKÝ TREND

Logistická trendová funkce patří k nelineárním trendům, které se vyznačují horní asymptotou, tj. hranicí, k níž se hodnoty ukazatele přibližují pro neomezeně rostoucí hodnoty času, a jedním inflexním bodem, v němž graf logistické funkce přechází z konvexního do konkávního tvaru. Pro tvar podobný písmenu S se takovými křivkám říká *S-křivky*. V ekonomické oblasti, speciálně v marketingu, se tato funkce používá při modelování poptávky po zboží dlouhodobé spotřeby, ale také při modelování vývoje výroby a prodeje některých druhů výrobků.

Na rozdíl od předchozích trendových funkcí, které byly definovány jednoznačně, logistická funkce bývá vyjadřována v několika různých variantách, uvedeme zde nejpoužívanější tvar:

$$T_t = \frac{\kappa}{1 + \beta_0 \beta_1^t}, \quad (8.20)$$

kde β_0, β_1, κ jsou neznámé parametry a $t = 1, 2, \dots, n$ je časová proměnná, přitom se kvůli zachování tvaru S-křivky předpokládá, že $0 < \kappa, 0 < \beta_0, 0 < \beta_1 < 1$. Odhady neznámých parametrů, označujeme je b_0, b_1, k lze opět získat metodou nejmenších čtverců, která dává nejlepší výsledky, i když vede na řešení soustavy nelineárních rovnic vyžadující použití složitějších výpočetních metod, např. *iteračních metod*. Proto zde ukážeme jinou metodu výpočtu neznámých parametrů, která sice nevede z teoretického pohledu k nejlepším odhadům, avšak její výhoda spočívá ve výpočetní nenáročnosti umožňující „ruční“ výpočet. Tato metoda se nazývá *metoda vybraných bodů* a spočívá v tom, že z daných údajů časové řady vybereme 3 charakteristické hodnoty (body), kterými necháme logistickou trendovou křivku procházet, jinými slovy, položíme empirické hodnoty rovny hodnotám teoretickým. Jestliže charakteristické hodnoty $T_{t_1}, T_{t_2}, T_{t_3}$ odpovídají časovým okamžikům t_1, t_2, t_3 , kde $t_1 < t_2 < t_3$, pak ze vztahu (4.33) obdržíme soustavu 3 rovnic o 3 neznámých β_0, β_1, κ :

$$T_{t_1} = \frac{\kappa}{1 + \beta_0 \beta_1^{t_1}}, T_{t_2} = \frac{\kappa}{1 + \beta_0 \beta_1^{t_2}}, T_{t_3} = \frac{\kappa}{1 + \beta_0 \beta_1^{t_3}}, \quad (8.21)$$

jejichž řešením získáme odhady neznámých parametrů b_0, b_1, k . Výpočty v metodě vybraných bodů můžeme usnadnit, když charakteristické body zvolíme ekvidistantně:

$$t_1 = 0, t_2 = \Delta, t_3 = 2\Delta,$$

kde Δ je určitý časový interval. Za tohoto předpokladu je řešení soustavy následující:

$$b_0 = \frac{k - T_{t_1}}{T_{t_1}},$$

$$b_1 = \left(\frac{T_{t_1} (k - T_{t_2})}{T_{t_2} (k - T_{t_1})} \right)^{\frac{1}{t_2}}, \quad (8.22)$$

$$k = \frac{2T_{t_1} T_{t_2} T_{t_3} - T_{t_2}^2 (T_{t_1} + T_{t_3})}{T_{t_1} T_{t_3} - T_{t_2}^2}. \quad (8.23)$$

Z výše uvedeného vztahu (8.23) lze přímo vypočítat parametr k , jeho dosazením do vztahu (8.22) vypočítáme parametry b_0, b_1 . Jak se snadno zjistí, hodnota asymptoty logistické křivky je $\frac{k}{1+\beta_0}$, což představuje horní mez, k níž se limitně přibližuje hodnota trendové funkce při velkých hodnotách času t .

ŘEŠENÁ ÚLOHA 8.3



V tabulce jsou uvedeny údaje o počtu výrobků určitého typu (v tis. ks) v letech 2015 - 2023. Nalezněte logistickou trendovou funkci, která charakterizuje trend dané časové řady. Prognózuje výrobu pomocí bodového odhadu na rok 2024.

Čas	2015	2016	2017	2018	2019	2020	2021	2020	2021	2022	2023
Zjištěné hodnoty	5	6	9	16	22	25	32	34	41	44	45

Řešení:

Hledáme odhady parametrů trendové funkce ve tvaru (8.20)

$$T_t = \frac{\kappa}{1 + \beta_0 \beta_1^t}.$$

Tyto odhady stanovíte metodou vybraných bodů. Abyste mohli k výpočtu použít vztahy (8.21), (8.22), (8.23), zvolíte opět novou časovou proměnnou t' , viz následující tabulka. Ze všech údajů v časové řadě vyberete tři časové okamžiky, např. na počátku, uprostřed a na konci časové osy: $t'_1 = 0, t'_2 = 5, t'_3 = 10$. V těchto okamžicích (jsou vyznačeny tučně) položíte empirické hodnoty rovny hodnotám teoretickým, tedy $T_{t'_1} = 5, T_{t'_2} = 25, T_{t'_3} = 45$.

t	2015	2016	2017	2018	2019	2020	2021	2020	2021	2022	2023
t'	0	1	2	3	4	5	6	7	8	9	10
Zjištěné hodnoty	5	6	9	16	22	25	32	34	41	44	45

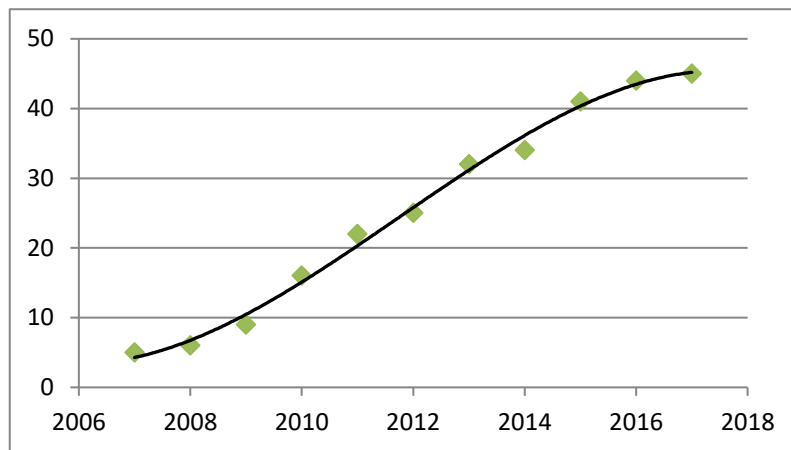
Potom ze vztahů (8.22), (8.23) postupně vypočítáte:

$$k = \frac{2T_{t'_1}T_{t'_2}T_{t'_3} - T_{t'_2}^2(T_{t'_1} + T_{t'_3})}{T_{t'_1}T_{t'_3} - T_{t'_2}^2} = \frac{2 \cdot 5 \cdot 25 \cdot 45 - 25^2(5 + 45)}{5 \cdot 45 - 25^2} = 50,$$

$$b_0 = \frac{k - T_{t'_1}}{T_{t'_1}} = \frac{50 - 5}{5} = 9, \quad b_1 = \left(\frac{T_{t'_1}(k - T_{t'_2})}{T_{t'_2}(k - T_{t'_1})} \right)^{\frac{1}{t'_2}} = \left(\frac{5(50 - 25)}{25(50 - 5)} \right)^{\frac{1}{5}} = 0,644.$$

Odhadovaný logistický trend má tvar

$$\hat{T}_{t'} = \frac{50}{1 + 9 \cdot 0,644^{t'}}.$$



Obrázek

Logistický trend

25:

Rok 2024 odpovídá v transformované časové ose hodnotě $t' = 13$. Dosazením do rovnice zjištěné trendové funkce obdržíte

$$\hat{T}_{2008} = \frac{50}{1 + 9 \cdot 0,644^{13}} = 48,57 \cong 49,$$

tj. prognózovaná výroba daného výrobku v roce 2024 je 49 tis. ks.

8.2.6 GOMPERTZŮV TREND

Ve srovnání s předchozí logistickou trendovou funkcí je *Gompertzův trend* jiným typem S-křivky:

$$T_t = k\beta_0\beta_1^t, \tag{8.24}$$

kde opět β_0, β_1, k jsou neznámé parametry a $t = 1, 2, \dots, n$ je časová proměnná, přitom se kvůli zachování tvaru S-křivky předpokládá, že $0 < k, 0 < \beta_0, 0 < \beta_1 < 1$. Odhady b_0, b_1 ,

těchto k parametrů získáme opět metodou nelineární regrese (metodou nejmenších čtverců), eventuálně metodou vybraných bodů, jako v předchozím odstavci. Asymptota Gompertzovy křivky je rovnoběžná s osou t ve vzdálenosti k , přičemž inflexní bod křivky není na rozdíl od logistického trendu (8.20) umístěn uprostřed mezi časovou osou a asymptotou.

8.3 Volba vhodného modelu trendu

Závažným problémem analýzy časových řad je problém stanovení konkrétního typu trendové funkce. Základem pro rozhodnutí o vhodném typu funkce by měla být věcně-ekonomická kritéria, tedy trendová funkce by měla být volena na základě věcné analýzy zkoumaného ekonomického jevu. Během věcného rozboru lze obvykle posoudit, zda jde o funkci rostoucí (nebo klesající), s trendem růstu nade všechny meze, či k určité konečné hodnotě (asymptotě).

Grafické znázornění časové řady umožní v hrubých rysech odhalit základní tendence ve vývoji analyzovaného ukazatele. Nebezpečí volby na základě vizuálního výběru spočívá však v jeho subjektivitě. Různí analytici mohou danou situaci posoudit různě a zvolit rozdílné typy trendové funkce. Nebezpečí tu plyne i z toho, že tvar grafu je do značné míry závislý na volbě použitého měřítka.

Přiléhavost dat k trendové (regresní) křivce jsme v kapitole 3 měřili koeficientem determinace R^2 , viz (3.18):

$$R^2 = \frac{S_T}{S_y} = 1 - \frac{S_R}{S_y}. \quad (8.25)$$

Tento koeficient můžeme k porovnání vhodnosti různých modelů trendu použít i nyní. V zásadě lze přijmout hodnocení, v němž nejvhodnější model trendu dává nejvyšší hodnotu koeficientu determinace R^2 . Vzhledem k tomu, že hodnota S_y je dána, závisí velikost R^2 na velikosti reziduálního součtu čtverců S_R ; čím je jeho hodnota menší, tím je hodnota R^2 větší (blíže k jedné). Taková metoda hodnocení trendu časové řady však upřednostňuje modely s větším počtem parametrů. Protože se zejména u ekonomických časových řad snažíme o nalezení jednoduchého tvaru trendu, je lepší k hodnocení vhodnosti modelu použít *reziduální rozptyl*:

$$s_R^2 = \frac{S_R}{n - p}, \quad (8.26)$$

kde $S_R = \sum_{i=1}^n (y_i - Y_i)^2$ je reziduální součet čtverců, n je počet datových bodů a p je počet parametrů v modelu. Z tvaru (8.26) je zřejmé, že hodnota reziduálního rozptylu roste s ros-

toucím počtem parametrů, což odpovídá výše uvedenému požadavku po co nejmenším počtu parametru v trendové funkci. Vhodný model trendu bude tedy „kompromisem“ mezi velikostmi hodnot R^2 a p .

Volbu vhodné trendové funkce lze podpořit také *testy hypotéz*. Z celé řady různých testů uvedeme známý F -test, který slouží pro rozhodování, zda má smysl dávat přednost složitějšímu modelu (s větším počtem parametrů) před jednodušším modelem (s menším počtem parametrů). Testujeme nulovou hypotézu, že totiž pokud jde o přiléhavost dat ke zvoleným trendovým funkcím, není mezi modely statisticky významný rozdíl. Tento test je založen na statistice:

$$F = \frac{\frac{S_T^{(2)} - S_T^{(1)}}{p_1 - p_2}}{\frac{S_R^{(1)}}{n - p_1}}, \quad (8.27)$$

kde hodnoty $S_T^{(1)}, S_R^{(1)}, p_1$ přísluší ke složitějšímu modelu, hodnoty $S_T^{(2)}, p_2$ přísluší k jednoduššímu modelu, tj. $p_1 > p_2, S_T^{(2)} > S_T^{(1)}$. Statistika (8.27) má přibližně Fisherovo rozdělení F s $p_1 - p_2$ a $n - p_1$ stupni volnosti. V případě, že vypočítaná hodnota statistiky padne do kritického oboru, lze na zvolené hladině významnosti α usuzovat, že model s větším počtem parametrů přináší výrazné zlepšení oproti jednoduššímu modelu.

8.4 Klouzavé průměry

Podstata vyrovnání časové řady pomocí klouzavých průměrů spočívá v tom, že posloupnost hodnot časové řady nahradíme novou řadou průměrů vypočítaných s kratších úseky časové řady, přičemž tyto kratší úseky postupně posouváme (kloužeme) směrem od začátku ke konci časové řady, a současně vypočítáváme dílčí průměry, tzv. *klouzavé průměry*. Vzniká důležitý problém, který je nutno předem řešit: jaký má být počet členů klouzavé části průměru. Klouzavou částí průměru budeme tedy rozumět časový interval určité délky, který se posunuje po časové ose vždy o jednotku. Volba rozsahu klouzavé části závisí na věcném (ekonomickém) charakteru časové řady a nelze ji obvykle stanovit na podkladě exaktních statistických metod. V praxi jsou u ekonomických neperiodických časových řad voleny většinou klouzavé části menší liché délky, např. 3, 5 nebo 7 časových jednotek, což souvisí se snadnější interpretací výsledků, neboť pak můžeme hodnotu klouzavého průměru přiřadit prostřednímu časovému okamžiku klouzavé části. U periodických časových řad se volí délka klouzavých částí totožná s délkou periody (sezóny, cyklu).

Uvažujme časovou řadu $y_1, y_2, y_3, \dots, y_n$. *Prosté klouzavé průměry* získáme tak, že úseky časové řady o délce $m = 2p + 1$, přičemž $m < n, p \geq 1$, celé číslo, vyrovnáme lineárním trendem s využitím metody nejmenších čtverců. Výsledkem je vzorec pro hodnoty vyrovnané časové řady ve formě aritmetického průměru:

$$\bar{y}_t = \frac{1}{2p+1} \sum_{i=-p}^p y_{t+i} = \frac{y_{t-p} + y_{t-p+1} + \dots + y_{t+p-1} + y_{t+p}}{2p+1}, \quad (8.28)$$

kde $t = p + 1, p + 2, \dots, n - p$. Přitom p hodnot na začátku a p hodnot na konci časové řady zůstává nevyrovnáno.

Kromě prostých klouzavých průměrů se někdy používají složitější *vážené klouzavé průměry*, případně *centrované klouzavé průměry*. Ty získáme tak, že namísto lineárního trendu v každém úseku použijeme polynomický trend vyššího řádu, tj. kvadratickou parabolu, kubickou parabolu apod. Metodou nejmenších čtverců obdržíme poměrně složité vzorce pro výpočet vyrovnaných hodnot. Vzhledem k poměrně řídkému použití těchto složitějších klouzavých průměrů se jimi zde nebudeme dále zabývat. Zájemce odkazujeme na literaturu, např. Seger (1998).

8.5 Exponenciální vyrovnání

Další metodou vyhlazování časové řady, tedy syntetického stanovení trendu, je *exponenciální vyrovnání*. Při něm se nová vyrovnaná hodnota stanoví na základě exponenciálně váženého průměru současné hodnoty a všech předchozích hodnot časové řady. Přitom se používá systém koeficientů, které nazýváme váhy, kdy novější hodnota má vždy větší váhu (tj. důležitost), než hodnota starší.

Nechť y_t značí pozorovanou hodnotu v časovém okamžiku t , w je váha přiřazená současné hodnotě, přičemž $0 < w < 1$, \hat{y}_t je vyrovnaná hodnota v čase t . Metoda exponenciálního vyrovnání začíná tím, že první vyrovnanou hodnotu časové řady \hat{y}_1 (v čase 1) položíme rovnu pozorované hodnotě y_1 , tedy $\hat{y}_1 = y_1$.

Následující vyrovnané hodnoty definujeme rekurentním vztahem:

$$\hat{y}_t = wy_t + (1 - w)\hat{y}_{t-1}, \quad t = 2, 3, \dots, n, \quad (8.29)$$

který umožňuje postupně vypočítat všechny vyrovnané hodnoty dané časové řady. Ze vztahu (8.28) lze snadno odvodit vztah:

$$\hat{y}_t = wy_t + w(1 - w)y_{t-1} + w(1 - w)^2y_{t-2} + \dots + w(1 - w)^{t-1}y_1.$$

Z posledního vztahu je vidět, že vyrovnaná hodnota časové řady v čase t závisí na všech předchozích nevyrovnaných hodnotách s tím, že do celkového součtu vstupují starší hodnoty s menší vahou

$$w_{t-i} = w(1 - w)^i, \quad (8.30)$$

kde $i = 0, 1, \dots, t-2$. Vzhledem k tomu, že platí $0 < w < 1$, je zřejmé, že se hodnota w_{t-i} exponenciálně zmenšuje s rostoucím i , tj. rostoucím „stářím“ dat. Váhu w nazýváme *koeficient exponenciálního zapomínání*. Ze vztahu (8.30) vyplývá, že čím vyšší je koeficient zapomínání, tím menší je hodnota $(1 - w)$, a tedy také $(1 - w)^i$, což znamená, že váha, tedy význam starších dat klesá, starší data se rychleji zapomínají. Je-li např. $w = 0,9$, tedy koeficient zapomínání je 90%, potom za jednotku času se vliv hodnoty y_{t-i} zmenší na $(1 - w)y_{t-i} = 0,1y_{t-i}$, což znamená, že se „zapomene“ 90% hodnoty. V praxi se používají obvykle váhy z intervalu 0,7 až 1,0. Pro výpočet exponenciálně vyrovnaných hodnot časové řady je ovšem výhodnější rekurentní vztah (8.29).

Kromě výše uvedené metody se v praxi využívají i složitější postupy exponenciálního vyrovnaní, které se zařazují do skupiny metod, kterým se říká *adaptivní metody*. Zájemce odkazujeme např. na práce Seger (1998), Cipra (1986).



ŘEŠENÁ ÚLOHA 8.4

V následující tabulce jsou uvedeny údaje o spotřebě pitné vody v jednotlivých dnech tří po sobě jdoucích týdnů.

- Stanovte odpovídající interval klouzavého průměru a vyrovnejte tuto řadu prostými klouzavými průměry.
- Vyrovnejte časovou řadu pomocí metody exponenciálního vyrovnaní, použijte koeficient zapomínání $w = 0,7$.

Po	0,64	0,75	0,54
Út	0,78	0,63	0,61
St	0,93	0,82	0,7
Čt	0,66	0,63	0,56
Pá	0,99	1,3	0,79
So	1,22	0,65	1,3
Ne	1,05	1,3	1,24

Řešení:

- Z charakteru dat vyplývá, že pro analyzovanou časovou řadu budou vhodné klouzavé průměry o délce $m = 7$ pozorování, tj. v rámci týdne. Použijete proto prosté 7-členné klouzavé průměry, které vypočítáte podle vztahu (8.28):

$$\bar{y}_1 = \frac{y_1 + y_2 + \dots + y_7}{7} = \frac{0,64 + 0,78 + 0,93 + 0,66 + 0,99 + 1,22 + 1,05}{7} = 0,896.$$

Tuto hodnotu přiřadíte prostřednímu časovému okamžiku klouzavé části, tj. ke čtvrté hodnotě dané časové řady.

Druhý klouzavý průměr vypočítáte analogicky posunutím o jeden den a přiřadíte jej k páté hodnotě původní časové řady:

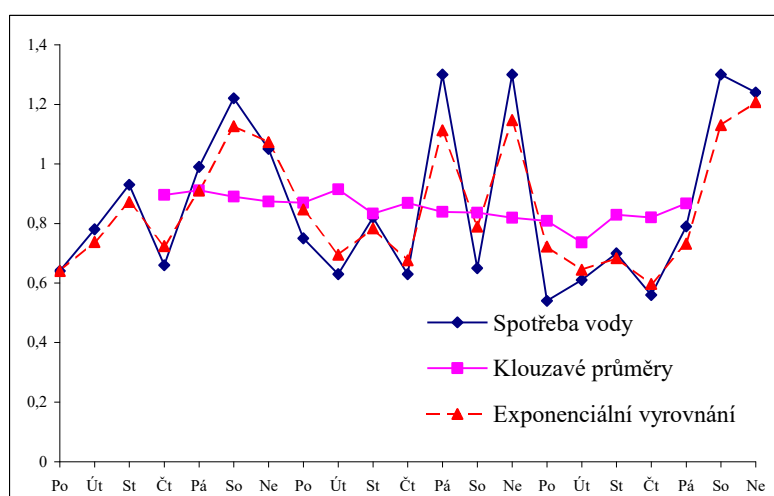
$$\bar{y}_2 = \frac{y_2 + y_3 + \dots + y_8}{7} = \frac{0,78 + 0,93 + 0,66 + 0,99 + 1,22 + 1,05 + 0,75}{7} = 0,911.$$

Ostatní klouzavé průměry vypočítáte obdobně postupným klouzáním směrem ke konci časové řady. Empirické hodnoty jakož i klouzavé průměry ukazuje Obrázek 26.

b. Exponenciální vyrovnání se provede podle (8.29):

$$\hat{y}_1 = y_1,$$

$$\hat{y}_t = wy_t + (1 - w)\hat{y}_{t-1}, t = 2, 3, \dots, n, \quad \text{kde } w = 0,7.$$



Obrázek 37: Klouzavé průměry a exponenciální vyrovnání

Proto:

$$\hat{y}_1 = 0,64,$$

$$\hat{y}_2 = 0,7y_2 + (1 - 0,7) \cdot \hat{y}_1 = 0,7 \cdot 0,78 + 0,3 \cdot 0,64 = 0,738.$$

Další hodnoty \hat{y}_t vypočítáme rekurentně, viz následující tabulka.

Den	Spotřeba vody (m ³ /os.)	Klouzavé průměry	Exponenciální vyrovnání
Po	0,64		0,640
Út	0,78		0,738
St	0,93		0,872
Čt	0,66	0,896	0,724
Pá	0,99	0,911	0,910
So	1,22	0,890	1,127
Ne	1,05	0,874	1,073
Po	0,75	0,870	0,847
Út	0,63	0,914	0,695
St	0,82	0,833	0,783
Čt	0,63	0,869	0,676
Pá	1,30	0,839	1,113
So	0,65	0,836	0,789
Ne	1,30	0,819	1,147
Po	0,54	0,809	0,722
Út	0,61	0,736	0,644
St	0,70	0,829	0,683
Čt	0,56	0,820	0,597
Pá	0,79	0,867	0,732
So	1,30		1,130
Ne	1,24		1,207

Je zřejmé, že koeficient zapomínání $w = 0,7$ ještě nevyhlazuje původní data dostatečně, k většímu vyhlazení by byla zapotřebí menší hodnota koeficientu zapomínání.



SAMOSTATNÉ ÚKOLY

8.1 V tabulce jsou údaje o počtu vyrobených kuchyňských robotů v letech 2015 až 2023.

Rok	2015	2016	2017	2018	2019	2020	2021	2020	2021	2022	2023
Kuchyňské roboty (tis. ks)	5	4	8	16	35	32	40	56	100	120	195

- Trend ve výrobě tohoto výrobku popište exponenciální trendovou funkcí.
- Jaké množství vyrobených kuchyňských robotů lze očekávat v roce 2024?
- Znaménkovým testem (bude vysvětlen v následující kapitole) ověřte na hladině významnosti $\alpha = 0,05$ náhodnost reziduí.

8.2 Časová řada představuje počet vyrobených pneumatik Barum v letech 2014 až 2023.

Rok	2014	2015	2016	2017	2018	2019	2020	2021	2020	2021	2022	2023
Pneumatiky (mil. ks)	0,8	1,6	1,5	2,4	5	3,88	4,47	3,88	6,89	7,69	5,83	8,25

- Nalezněte lineární trend časové řady.
- Jaké množství vyrobených pneumatik lze očekávat v roce 2024? Stanovte bodový i intervalový odhad na hladině významnosti $\alpha = 0,05$.

ODPOVĚDI



- 8.1** a) $\hat{T} = 29,55 \cdot 1,47^t, t = -5, -4, -3, \dots$
b) v roce 2024, tzn. $t = 8$; $\hat{T} = 644,31$.
c) $S = 5$; testové kritérium $U = 0$; obor přijetí $A = (-1,96; 1,96)$; přijímáme nulovou hypotézu o náhodném uspořádání reziduí
- 8.2** a) $\hat{T} = 4,35 + 0,32 \cdot t, t = -11, -9, -7, \dots$
b) v roce 2024, tzn. $t = 13$; $\hat{T} = 8,5$ mil. ks; 95 % intervalový odhad (5,97; 11,05)
-

SHRNUTÍ KAPITOLY



Zopakujme si získané poznatky této kapitoly: trendová složka poskytuje rozhodující informaci pro prognózování hodnot časové řady do budoucna. K určení trendové složky používáme dva obecné přístupy: analytický a syntetický. Analytický přístup vychází z předem známých typů trendových funkcí vyznačujících se přítomností parametrů. V této kapitole jsme se zabývali lineárním trendem, parabolickým trendem, exponenciálním trendem, logistickým trendem a Gompertzovým trendem. Z metod syntetického přístupu byly uvedeny metody klouzavého průměru a exponenciální vyrovnání.

9 SEZÓNÍ SLOŽKA, NÁHODNÁ SLOŽKA



RYCHLÝ NÁHLED KAPITOLY

Při analýze ekonomických časových řad se setkáváme téměř vždy s existencí sezónních vlivů, reprezentovaných v modelu časové řady sezónní složkou. Sezónními vlivy rozumíme soubor příčin, které se pravidelně opakují v důsledku koloběhu přírody. Pokud se u časových řad vyskytují podobné vlivy v delším časovém horizontu, hovoříme o cyklické složce časové řady, v kratším časovém horizontu, hovoříme o sezónní složce časové řady. Souhrnně se sezónní a cyklické složky označují jako periodické složky časové řady. Úkolem modelování periodické složky časové řady je nalézt její vhodné vyjádření, které by umožnilo periodickou (nejčastěji sezónní) složku vhodně identifikovat a následně použít k predikci chování časové řady v budoucnu. Naučíte se aplikovat metody konstantní sezónnosti se schodovitým a lineárním trendem a metodu proporcionální sezónnosti. V závěru se budete věnovat analýze náhodné složky.



CÍLE KAPITOLY

Po prostudování této kapitoly budete umět:

- popsat sezónní a náhodnou složku,
 - použít metodu konstantní sezónnosti se schodovitým trendem,
 - použít metodu konstantní sezónnosti s lineárním trendem,
 - testovat vlastnosti náhodné složky.
-



ČAS POTŘEBNÝ KE STUDIU

K prostudování této kapitoly budete potřebovat asi 90 minut.

KLÍČOVÁ SLOVA KAPITOLY

Sezónní složka, náhodná složka, model konstantní sezónnosti se schodovitým trendem, model konstantní sezónnosti s lineárním trendem, znaménkový test, Durbin-Watsonův test.

9.1 Model konstantní sezónnosti se schodovitým trendem

Označení časové proměnné $t = 1, 2, \dots, n$, budeme používat pro označení časových intervalů (např. roků), které se člení na dalších r dílčích časových obdobích, které nazýváme *sezóny* (např. měsíce nebo čtvrtletí) a označujeme $j = 1, 2, \dots, r$ (např. v případě, že sezóny jsou měsíce je $r = 12$, v případě že sezóny představují kvartály, platí $r = 4$). Model časové řady lze zapsat ve tvaru:

$$y_{tj} = T_{tj} + P_{tj} + \varepsilon_{tj}, t = 1, 2, \dots, n, j = 1, 2, \dots, r. \quad (9.1)$$

U modelu *konstantní sezónnosti* se vychází z předpokladu, že:

$$P_{tj} = \gamma_j \text{ pro sezónu } j \text{ v letech } t = 1, 2, \dots, n, \quad (9.2)$$

kde γ_j jsou neznámé sezónní parametry, o nichž dále předpokládáme, že splňují rovnost:

$$\sum_{j=1}^r \gamma_j = 0. \quad (9.3)$$

Předpoklady (9.2) a (9.3) vycházejí z představy, že v důsledku pravidelného (ročního) koloběhu sezónních vlivů se v j -té sezóně opakují sezónní výkyvy γ_j , které se mezi léty neliší, to je podmínka (9.2). Dále se tyto vlivy během roku (r sezón) vykompenzují, takže jejich roční součet je nulový, což odpovídá podmínce (9.3).

Nejprve budeme předpokládat, že trendová složka T_{tj} nabývá ve všech sezónách hodnotu roku t hodnotu α_t , takže posloupnost těchto hodnot v letech $t = 1, 2, \dots, n$ představuje *schodovitý trend*. Model (9.1) pak bude mít tvar:

$$y_{tj} = \alpha_t + \gamma_j + \varepsilon_{tj}, t = 1, 2, \dots, n, j = 1, 2, \dots, r. \quad (9.4)$$

Odhady a_t, c_j $n + r$ parametrů tohoto modelu získáme metodou nejmenších čtverců:

$$a_t = \frac{1}{r} \sum_{j=1}^r y_{tj} = \bar{y}_t, \quad c_j = \frac{1}{n} \sum_{t=1}^n y_{tj} - \frac{1}{rn} \sum_{t=1}^n \sum_{j=1}^r y_{tj}. \quad (9.5)$$

Všimněte si v prvním vzorci, že odhadem výšky schodu v roce t je průměr hodnot v roce t . Z druhého vzorce pak vyplývá, že hodnota sezónního vlivu c_j , tzv. j -tého sezónního koeficientu, je představována průměrnou hodnotou vypočítanou z j -tých sezón ve všech letech po odečtení celkového průměru ze všech hodnot v celé časové řadě. Například sezónní koeficient c_1 se vypočítá jako průměr ze všech lednových hodnot v časové řadě měsíčních údajů po odečtení celkového průměru ze všech hodnot v celé časové řadě. V tomto případě je měsíc leden uvažován jako první sezóna z 12 měsíčních sezón.

9.2 Model konstantní sezónnosti s lineárním trendem

Při popisu trendové složky v předchozím odstavci jsme používali posloupnost časové proměnné $t = 1, 2, \dots, n$, o trendové funkci jsme předpokládali, že je konstantní během všech sezón daného roku t , tj. $T_{tj} = \alpha_t$ pro $j = 1, 2, \dots, r$. Přitom hodnota α_t mohla být v každém roce jiná a tvořila výšku „schodu“ v roce t . Model časové řady bude opět aditivní, tedy

$$y_{tj} = T_t + \gamma_j + \varepsilon_{tj}, t = 1, 2, \dots, n, j = 1, 2, \dots, r, \quad (9.6)$$

kde stejně jako v modelu (9.1) jsou γ_j neznámé sezónní parametry, o nichž dále předpokládáme, že splňují podmínku $\sum_{j=1}^r \gamma_j = 0$.

Nyní budeme předpokládat, že trendová složka T_{tj} má lineární tvar, potom model (9.6) bude mít tvar:

$$y_{tj} = \alpha + \beta(t - \bar{t}) + \gamma_j + \varepsilon_{tj}, t = 1, 2, \dots, n, j = 1, 2, \dots, r. \quad (9.7)$$

Odhady a, b, c_j z $(r + 2)$ parametrů tohoto modelu získáme metodou nejmenších čtverců, řešení má komplikovaný tvar, který zde neuvádíme, zájemce odkazujeme na Segera (1998).

9.3 Model proporcionální sezónnosti

Nyní budeme používat $t = 1, 2, \dots, n$, k označení časových intervalů (např. roků), které se člení na dalších r dílčích časových obdobích, které nazýváme *sezóny* (např. měsíce nebo čtvrtletí) a označujeme $j = 1, 2, \dots, r$ (např. v případě, že sezóny jsou měsíce je $r = 12$, v případě že sezóny představují kvartály, platí $r = 4$). Regresní model lze s použitím uvedené symboliky zapsat ve tvaru:

$$y_{tj} = T_{tj} + P_{tj} + \varepsilon_{tj}, t = 1, 2, \dots, n, j = 1, 2, \dots, r. \quad (9.8)$$

U modelu *proporcionální sezónnosti* se vychází z předpokladu, že periodická složka je proporcionální (tj. přímo úměrná) velikosti trendové složky:

$$P_{tj} = C_j T_{tj} \text{ pro sezónu } j \text{ v letech } t = 1, 2, \dots, n, \quad (9.9)$$

tedy po dosazení (9.9) do (9.8) obdržíte

$$y_{tj} = (1 + C_j)T_{tj} + \varepsilon_{tj}. \quad (9.10)$$

Aplikací MNC obdržíme c_j odhad koeficientů C_j takto

$$1 + c_j = \frac{\sum_{i=1}^n y_{ij} \bar{y}_i}{\sum_{i=1}^n \bar{y}_i^2}, \quad j = 1, 2, \dots, r. \quad (9.11)$$

Dosazením do (9.10) obdržíte konečnou podobu modelu *proporcionální sezónnosti*

$$y_{tj} = \frac{\sum_{i=1}^n y_{ij} \bar{y}_i}{\sum_{i=1}^n \bar{y}_i^2} T_{tj} + \varepsilon_{tj}, \quad t = 1, 2, \dots, n, \quad j = 1, 2, \dots, r. \quad (9.12)$$

Přitom $\bar{y}_i = \frac{1}{r} \sum_{j=1}^r y_{ij}$ je aritmetický průměr y_{ij} přes j . V konkrétním případě můžeme uvažovat, že trendová složka má lineární tvar, tedy například

$$T_{tj} = \alpha + \beta(t - \bar{t}). \quad (9.13)$$

9.4 Analýza náhodné složky

Náhodnou složku ε_t lze v modelu (9.8) vyjádřit v tvaru:

$$\varepsilon_t = y_t - Y_t, \quad t = 1, 2, \dots, n, \quad (9.14)$$

kde $Y_t = T_t + P_t$. Jedná se zde o vyjádření blíže nespecifikovaných náhodných vlivů. Zdrojem této složky jsou obvykle nepodchycené drobné vzájemně nezávislé náhodné vlivy. Chceme-li zajistit spolehlivé předpovědi na základě modelu časové řady, potom je třeba mít zajištěny některé předpoklady o náhodné složce. Konkrétně je výhodné, když jsou splněny předpoklady klasického lineárního regresního modelu, které jsme uvedli v kapitole 3.5. Byly to předpoklady 1. až 3., které pro přehlednost zopakujeme, avšak při současném označení, kdy nezávisle proměnná x je nyní čas t . Jedná se tedy o tyto předpoklady:

1. Hodnoty vysvětlující proměnné t se volí předem, obvykle $t = 1, 2, \dots, n$.
2. Náhodné složky ε_t mají *normální rozdělení* pravděpodobnosti se střední hodnotou 0 a (neznámým) rozptylem σ^2 .

Konstantnost rozptylu nazýváme *homoskedasticita*.

3. Náhodné složky jsou *nekorelované*, tj.
 $Cov(\varepsilon_t, \varepsilon_{t'}) = 0$ pro každé $t \neq t'$, $t, t' = 1, 2, \dots, n$.

Jak již bylo řečeno v kapitole 3.5, v praxi jsou podmínky klasického modelu často splněny. Nejsme-li si však jejich platností jisti, můžeme provést testy hypotéz jak o normalitě rozdělení náhodné složky (např. Chi-kvadrát test dobré shody), tak i testy homoskedasticity (Bartleyův test). Při ověřování těchto předpokladů zjišťujeme, zda jsou všechny systematické složky z časové řady eliminovány. Jakákoliv nenáhodnost u reziduí naznačuje nevhodnost zvoleného modelu časové řady.

Jednoduchým nástrojem, kterým lze ověřit náhodnost reziduí, je *znaménkový test*. Při tomto testu vyčíslíme počet případů, kdy rozdíl sousedních reziduí $e_t - e_{t-1}$ je kladný, jejich počet označíme S . Přitom je:

$$e_t = y_t - Y_t, \quad (9.15)$$

kde $Y_t = T_t + P_t$ je odhad teoretické hodnoty časové řady, T_t je odhad trendu (s regresními koeficienty získanými např. metodou nejmenších čtverců), P_t je odhad periodické složky, např. (9.11), kde parametry α_j, β_j jsou rovněž odhadnuty metodou nejmenších čtverců. Náhodné složky ε_t , které jsou dány (9.14), jsou tedy náhodné veličiny, zatímco rezidua e_t , (9.15), jsou realizacemi, jsou to odhady těchto náhodných veličin. Je-li posloupnost reziduí e_t náhodně uspořádána, potom pro střední hodnotu S platí: $E(S) = \frac{n-1}{2}$. Testujeme proto nulovou hypotézu: $H_0 : E(S) = \frac{n-1}{2}$, proti alternativní hypotéze $H_1 : E(S) \neq \frac{n-1}{2}$. Použijeme testové kritérium:

$$U = \frac{\sqrt{12} \left(S - \frac{1}{2}(n-1) \right)}{\sqrt{n+1}}, \quad (9.16)$$

které má již pro $n \geq 13$ přibližně normované normální rozdělení. Pro stanovení kritických hodnot tedy použijeme *kvantily* normovaného normálního rozdělení $u_{1-\alpha/2}$.

Vlastnost časových řad, která často způsobuje porušení předpokladů 1. až 3. se nazývá *autoregrese* náhodných složek, viz též kapitola 6.5, která znamená, že mezi náhodnými složkami platí následující vztah:

$$\varepsilon_t = \rho \varepsilon_{t-1} + u_t, \quad (9.17)$$

kde $0 < \rho < 1$ je *autokorelační koeficient* a u_t splňuje předpoklady 1. až 3. Nulovou hypotézu: $H_0 : \rho = 0$ (což je totéž, jako $\varepsilon_t = u_t$) testujeme proti alternativní hypotéze $H_1 : \rho \neq 0$ pomocí testového kritéria:

$$D = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2}. \quad (9.18)$$

Funkce D , nazývaná *Durbin-Watsonova statistika*, bývá tabelována pro různé hladiny významnosti α , viz např. Gujarati (2003). Test založený na této statistice nazýváme *Durbin-Watsonův test autokorelace*.

ŘEŠENÁ ÚLOHA 9.1



Data v tabulce představují objem přepravy po vodních tocích ČR v jednotlivých čtvrtletích pěti po sobě jdoucích let.

tj	Čtvrtletí					
Roky	1	2	3	4	Součet	Průměr
1	120	138	132	114	504	126,00
2	118	138	150	119	525	131,25
3	149	161	155	145	610	152,50
4	150	173	181	159	663	165,75
5	178	195	198	183	754	188,50
Součet	715	805	816	720	3056	
Průměr	143	161	163,2	144		152,80

- Naleznete pro tuto časovou řadu model konstantní sezónnosti se schodovitým trendem.
- Na hladině významnosti $\alpha = 0,05$ ověřte náhodnost reziduí.

Řešení:

- Úkolem je nalézt odhady parametrů α_t, γ_j modelu

$$y_{tj} = \alpha_t + \gamma_j + \varepsilon_{tj}, \quad t = 1, 2, \dots, n, \quad j = 1, 2, \dots, r,$$

kde α_t je trendová složka,

γ_j je sezónní složka.

Odhady a_t, c_j $n + r$ parametrů tohoto modelu vypočítáme ze vztahů (9.5):

$$a_t = \frac{1}{r} \sum_{j=1}^r y_{tj} = \bar{y}_t, \quad c_j = \frac{1}{n} \sum_{t=1}^n y_{tj} - \frac{1}{rn} \sum_{t=1}^n \sum_{j=1}^r y_{tj}.$$

Všechny potřebné součty a průměry jsou uvedeny v tabulce, jejich dosazením do daných vztahů obdržíte:

$$\text{trendová složka: } a_1 = 126 \quad a_2 = 131,25 \quad a_3 = 152,5 \quad a_4 = 165,75 \quad a_5 = 188,5$$

$$\text{sezónní složka: } c_1 = 143 - 152,8 = -9,8 \quad c_2 = 161 - 152,8 = 8,2$$

$$c_3 = 163,2 - 152,8 = 10,4 \quad c_4 = 144 - 152,8 = -8,8.$$

Výsledky ukazují, že působení sezónních vlivů klesl v prvním čtvrtletí objem přepravy o 9,8 tun a ve čtvrtém čtvrtletí o 8,8 tuny. Tento pokles je vykompenzován růstem přepravy ve zbylých dvou čtvrtletích o 8,2 a 10,4 tun, tj. ve čtvrtletích pro říční přepravu klimaticky příznivějších. Z vývoje ročních průměrů a_t je zřejmé, že se průměrný roční objem přepravy neustále zvyšoval.

- Nejdříve vypočítáte odhady teoretických hodnot \hat{Y} dané časové řady tak, že odhadnete trendovou i sezónní složku. Např.:

sezónní složka, náhodná složka

$$\hat{Y}_{1,1} = a_1 + c_1 = 126 + (-9,8) = 116,2$$

$$\hat{Y}_{1,2} = a_1 + c_2 = 126 + 8,2 = 134,2$$

Všechny hodnoty $\hat{Y}_{t,j}$ jsou uvedeny v následující tabulce.

t/j	1	2	3	4
1	116,20	134,2	136,4	117,2
2	121,45	139,5	141,7	122,5
3	142,70	160,7	162,9	143,7
4	155,95	174,0	176,2	157,0
5	178,70	196,7	198,9	179,7

Dále vypočítáme hodnoty reziduí. Např.:

$$e_{1,1} = y_{1,1} - \hat{Y}_{1,1} = 120 - 116,2 = 3,8, \quad e_{1,2} = y_{1,2} - \hat{Y}_{1,2} = 138 - 134,2 = 3,8.$$

Hodnoty všech reziduí jsou uvedeny v následující tabulce:

t/j	1	2	3	4
1	3,80	3,80	-4,40	-3,20
2	-3,45	-1,45	8,35	-3,45
3	6,30	0,30	-7,90	1,30
4	-5,95	-0,95	4,85	2,05
5	-0,70	-1,70	-0,90	3,30

K testu náhodnosti reziduí použijeme znaménkový test. Je proto třeba určit počet případů S , kdy je rozdíl sousedních reziduí $e_t - e_{t-1}$ kladný. Např.:

$$e_{1,2} - e_{1,1} = 3,8 - 3,8 = 0,$$

$$e_{1,3} - e_{1,2} = -4,4 - 3,8 = -8,2.$$

V následující tabulce jsou případy, kdy $e_t - e_{t-1} > 0$, označeny „+“, ostatní „-“.

t/j	1	2	3	4
1	.	-	-	+
2	-	+	+	-
3	+	-	-	+
4	-	+	+	-
5	-	-	+	+

Z tabulky vidíme, že $S = 9$. Hodnotu testového kritéria vypočítáme podle (9.16):

$$U = \frac{\sqrt{12} \left(S - \frac{1}{2}(n-1) \right)}{\sqrt{n+1}} = \frac{\sqrt{12} \left(9 - \frac{1}{2}(20-1) \right)}{\sqrt{20+1}} = -0,378.$$

V tabulce normovaného normálního rozdělení nalezneme $u_{1-\alpha/2}$, tj.: $u_{0,975} = 1,96$.

Protože hodnota testového kritéria $-0,378$ leží v oboru přijetí $A = (-1,96; 1,96)$, lze na zvolené hladině významnosti přijmout nulovou hypotézu, tj. hypotézu o náhodném uspořádání reziduí.

SAMOSTATNÉ ÚKOLY



9.1 V následující tabulce jsou uvedeny měsíční tržby jedné obchodní organizace za posledních 60 měsíců od ledna 2019 až do prosince 2023.

- a. Nalezněte model konstantní sezónnosti se schodovým trendem.
- b. Pro rok 2024 uvažujte s růstem 5% (tj. výška schodu). Prognózuje tržby na rok 2024.

1	2	3	4	5	6	7	8	9	10	11	12
6489	5971	6272	6944	7217	7448	7259	7602	7651	8064	7952	8498
13	14	15	16	17	18	19	20	21	22	23	24
6930	6391	6979	7315	7798	7861	7994	7798	8022	8155	8694	8764
25	26	27	28	29	30	31	32	33	34	35	36
7560	7182	7077	7847	8603	8659	8827	8855	8337	8379	8834	9709
37	38	39	40	41	42	43	44	45	46	47	48
7833	7406	7791	8190	8869	8988	8736	9254	9240	9380	9422	9954
49	50	51	52	53	54	55	56	57	58	59	60
8442	7987	8673	8925	9534	9534	9331	9877	9695	9730	10192	10661

9.2 Použijte data z řešené úlohy 9.1. Nalezněte pro tuto časovou řadu model konstantní sezónnosti s lineárním trendem.

9.3 Je dána reziduální složka, která obsahuje tyto hodnoty:

0,652 0,767 -1,667 2,579 -0,254 0,963 0,188 -0,936 0,572 -2,863.

- Proveďte: a) znaménkový test náhodnosti reziduí,
b) Durbin – Watsonův test autokorelace.

ODPOVĚDI



9.1 a) $a_1 = 7280,6$; $a_2 = 7630,6$; $a_3 = 8322,4$; $a_4 = 8755,3$; $a_5 = 9381,7$; $a_6 = 9850,8$

$c_1 = -823,3$; $c_2 = -1286,7$; $c_3 = -915,7$; $c_4 = -429,9$; $c_5 = 130,1$; $c_6 = 223,9$;
 $c_7 = 155,3$; $c_8 = 403,1$; $c_9 = 314,9$; $c_{10} = 467,5$; $c_{11} = 744,7$; $c_{12} = 1243,1$

b)

leden 2024	9027,51	červenec 2024	10006,1
únor 2024	8564,11	srpen 2024	10253,9
březen 2024	8935,11	září 2024	10165,7

duben 2024	9420,91	říjen 2024	10318,3
květen 2024	9980,91	listopad 2024	10595,5
červen 2024	10074,7	prosinec 2024	11093,9

$$9.2 \quad Y_t = 6782,2 + 49,536.t + c_j$$

$$c_1 = -569,8; c_2 = -1082,7; c_3 = -761,3; c_4 = -325; c_5 = 185,4; c_6 = 229,7; \\ c_7 = 111,6; c_8 = 309,8; c_9 = 172,1; c_{10} = 275,2; c_{11} = 502,8; c_{12} = 951,7$$

leden 2024	9234,1	červenec 2024	10212,7
únor 2024	8770,7	srpen 2024	10460,5
březen 2024	9141,7	září 2024	10372,3
duben 2024	9627,5	říjen 2024	10524,9
květen 2024	10187,5	listopad 2024	10802,1
červen 2024	10281,3	prosinec 2024	11300,5

9.3 a) Počet kladných hodnot $S = 4$; $U = -0,522$. Protože hodnota $-0,522$ leží v oboru přijetí $A = (-1,96; 1,96)$, lze na zvolené hladině významnosti přijmout nulovou hypotézu, tj. hypotézu o náhodném uspořádání reziduí.

b) Hodnota Durbin – Watsonova koeficientu $D = 2,368$. Protože $k = 1$ a $n = 10$ najdeme pro $\alpha = 0,05$ v tabulkách $d_L = 0,879$; $d_U = 1,32$. Nelze zamítnout nulovou hypotézu, což znamená, že v modelu nebyla prokázána statisticky významná autokorelace.



SHRNUTÍ KAPITOLY

V této kapitole jste se zabývali časovými řadami, jejichž hodnoty se periodicky opakují, tzv. sezónními časovými řadami. Nejprve jste si objasnili význam sezónní složky časové řady. Poté jste se naučili aplikovat jednoduché metody konstantní sezónnosti se schodovitým a lineárním trendem a rovněž metodu proporcionální sezónnosti. Dále zde byly uvedeny metody testování náhodné složky (znaménkový test, Durbin-Watsonův test autokorelace).

10 MODELY TYPU ARIMA A PREDIKCE ČASOVÝCH ŘAD

RYCHLÝ NÁHLED KAPITOLY



Nejprve se budete zabývat časovými řadami typu ARIMA. Box-Jenkinsova metodologie, která se modely analýzy časových řad typu ARIMA zabývá, klade důraz nikoliv na konstrukci jedno-rovnicového nebo vícerovnicového modelu, jak je tomu např. v regresní analýze, nýbrž na analýzu vlastních stochastických vlastností ekonomických ČŘ. Postupně se seznámíte s vlastnostmi autoregresivních procesů AR, procesů pohyblivých průměrů MA, integračních procesů I, jakož i procesů vzniklých jejich kombinací: ARIMA. Dále lze tyto procesy rozšířit též na sezónní procesy. Úkolem pak je pro časovou řadu nalézt vhodný model typu ARIMA a nalezený model použít pro účely prognózy (predikce, extrapolace) hodnot dané časové řady. Celý postup tvorby prognózy ČŘ autoři metody ARIMA formulovali ve 4 krocích, které nazýváme Box-Jenkinsova metodologie prognózování ČŘ. Jednotlivé kroky jsou (1) Identifikace modelu, (2) Odhad modelu, (3) Verifikace modelu a (4) Prognóza pomocí modelu, a budou ilustrovány na příkladu časové řady čtvrtletního HDP České republiky s pomocí statistického programu GRETL.

CÍLE KAPITOLY



Po prostudování této kapitoly budete umět:

nalézt vhodný model typu ARIMA

použít model pro účely predikce časové řady,

formulovat 4 kroky Box-Jenkinsovy metodologie,

použít pro výpočet ARIMA modelu časové řady program GRETL.

ČAS POTŘEBNÝ KE STUDIU



K prostudování této kapitoly budete potřebovat asi 120 minut.



KLÍČOVÁ SLOVA KAPITOLY

ARIMA model, Box-Jenkinsova metodologie, identifikace modelu, odhad modelu, verifikace modelu, predikce pomocí modelu, program GRETl.

10.1 Program GRETl

GRETl (Gnu Regression, Econometrics and Time-series Library) je open-source statistický software navržený pro analýzu dat, ekonometrii, regresní analýzu a analýzu časových řad. Software je vyvinut v rámci projektu GNU a je k dispozici zdarma pro všechny uživatele.

Klíčové rysy a funkce softwaru GRETl jsou tyto:

- **Ekonometrické funkce:** GRETl poskytuje širokou škálu funkcí pro ekonometrickou analýzu, včetně lineární regrese, logit a probit modelů, ARIMA modelů pro analýzu časových řad, panelových datových modelů a dalších.
- **Uživatelské rozhraní:** Software nabízí grafické uživatelské rozhraní (GUI), což umožňuje uživatelům provádět analýzu dat bez nutnosti psaní kódu. Nicméně je také možné používat skriptovací jazyk GRETl pro vytvoření vlastních analýz a scénářů.
- **Import a export dat:** GRETl umožňuje importovat data z různých formátů, včetně textových souborů, Excelu a databází. Také umožňuje exportovat výsledky analýz do různých formátů.
- **Grafy:** Software obsahuje nástroje pro vytváření grafů a vizualizaci dat. Uživatelé mohou vytvářet grafy pro zobrazení vztahů mezi proměnnými a vizualizaci výsledků analýzy.
- **Dokumentace a podpora:** GRETl je doprovázeno podrobnou dokumentací a nápovědou, která uživatelům pomáhá pochopit jeho funkce a použití. Také existuje uživatelská komunita a fóra, kde lze získat pomoc a diskutovat o různých aspektech softwaru.
- **Rozšiřitelnost a platformní nezávislost:** Software je navržený tak, aby byl rozšiřitelný pomocí doplňků a skriptů. Uživatelé mohou vytvářet vlastní funkce, modely a rozšíření, což umožňuje přizpůsobit GRETl specifickým potřebám. GRETl je dostupný pro různé operační systémy.

Celkově lze GRETL považovat za užitečný nástroj pro analýzu dat, zejména v oblasti ekonometrie a analýzy časových řad. Díky kombinaci grafického rozhraní a možnosti skriptování je vhodný jak pro začátečníky, tak i pro pokročilé uživatele se znalostmi statistiky a ekonometrie.

Prognózování (předvídání, předpovídání) je důležitou součástí ekonomických (ekonomických) analýz, dá se říci, že z určitého pohledu nejdůležitější. Jak prognózovat budoucí hodnoty ekonomických veličin, jako jsou HDP, inflace, kurzy měn, ceny akcií, míra nezaměstnanosti, počet nově nakažených osob a dalších? Jednu klasikou metodu již znáte: lineární, (resp. nelineární) regresní analýza, s níž jste se seznámili již v kapitolách 3 a 4. V této kapitole se dozvíte o nové metodě, která se stala v posledních letech velmi populární: tzv. modely autoregresivních a integrovaných procesů a klouzavých průměrů ARIMA (Auto Regresive Integrated Moving Average), která je známa také pod názvem Box-Jenkinsova metodologie (podle autorů metody G.P.E. Boxe a G.M. Jenkinse).

Téma ekonomického prognózování je velmi široké a existuje k němu množství specializovaných knih a dalších publikací. My zde chceme podat pouze stručný vhled do problematiky. Naštěstí k problematice prognózování ekonomických ČŘ existuje nejen vhodná literatura, její přehled lze nalézt např. u Arlta (1999), u Gujaratho (2003) aj., ale též příslušný specializovaný SW v podobě programových balíků, jakými jsou GRETL, SPSS, EViews, STATISTICA, SAS a další. V této kapitole budeme využívat konkrétně program GRETL.

Jak jsme již dříve zmínili, k analýze časových řad existuje řada různých metod a přístupů. Kromě již zmíněné (1) jednoduché regresní analýzy a (2) metody ARIMA, které jsou předmětem tohoto textu, je zapotřebí ještě jmenovat (3) metody *exponenciálního vyrovnání* (Holtova-Wintersova metoda a jejich varianty), (4) metody *simultánních rovnic* a (5) *vektorové autoregresivní metody* VAR, (6) metody ARCH a GARCH a další. S nimi se zájemci mohou blíže seznámit např. v Seger (1998).

10.2 Modelování časových řad pomocí ARIMA modelu

Podle svých autorů známa jako Box-Jenkinsova metodologie, avšak technicky nazývaná ARIMA metodologie klade důraz nikoliv na konstrukci jednorovnicového nebo vícerovnicového modelu, jak je tomu např. v regresní analýze, nýbrž na analýzu vlastních stochastických vlastností ekonomických ČŘ podle filosofie „ať data hovoří sama za sebe“. V regresních modelech je závisle proměnná Y vysvětlována několika vysvětlujícími proměnnými – regresory, zatímco v ARIMA metodách je závisle proměnná Y v čase t vysvětlována hodnotami téže Y v minulých časových okamžicích a zároveň chybovými členy v sou-

časných anebo minulých okamžicích. Na rozdíl od regresních modelů a modelů simultánních rovnic, které jsou založeny na ekonomické teorii, nejsou modely ARIMA na teorii přímo závislé. Teoretické závislosti jsou u nich vyjádřeny zprostředkovaně skrze sledované hodnoty v minulých časových okamžicích.

10.2.1 AUTOREGRESIVNÍ PROCES (AR)

Budeme předpokládat, že Y_t se chová podle vztahu

$$(Y_t - \mu) = \varphi_1(Y_{t-1} - \mu) + u_t, \quad (10.1)$$

kde μ je střední hodnota Y_t a u_t je bílý šum, φ_1 je konstanta. V tom případě říkáme, že ČŘ Y_t je *autoregresivní proces 1. řádu* neboli AR (1). Podle modelu (10.1) je prognóza $Y_t - \mu$ v čase t je přímo úměrná $Y_{t-1} - \mu$ v čase $(t-1)$ prostřednictvím koeficientu úměry φ_1 plus/mínus náhodná chyba (bílý šum). Pokud pro konstantu v modelu (10.1) platí $-1 < \varphi_1 < 1$, pak se dá ukázat, že proces AR (1) je stacionární. Dále si všimněte, že speciálně při $\varphi_1 = 0$ je z (10.1) proces AR (1) bílý šum a při $\varphi_1 = 1$ je z (10.1) proces AR (1) náhodná procházka. Také pro $\varphi_1 \geq 1$ nebo $\varphi_1 < -1$ je proces AR (1) nestacionární (Arlt, 1999).

Podobně *autoregresivní proces 2. řádu* neboli AR (2) má tvar

$$(Y_t - \mu) = \varphi_1(Y_{t-1} - \mu) + \varphi_2(Y_{t-2} - \mu) + u_t. \quad (10.2)$$

Analogicky *autoregresivní proces p-tého řádu*, neboli AR(p) má tvar

$$(Y_t - \mu) = \varphi_1(Y_{t-1} - \mu) + \varphi_2(Y_{t-2} - \mu) + \dots + \varphi_p(Y_{t-p} - \mu) + u_t. \quad (10.3)$$

Otázka stacionarity procesů AR(p) pro $p > 1$ je složitější problém, kterým se zde zabývat nebudeme. Eventuální zájemce odkazujeme na literaturu, např. knihu Arlt (1999). Všimněte si, že kromě hodnot Y v různých časových okamžicích se ve výše uvedených modelech nevyskytují jiné regresory. V tomto smyslu říkáme, že „data hovoří sama za sebe“.

10.2.2 PROCES KLOUZAVÝCH PRŮMĚRŮ (MA)

Výše uvedený AR proces není jediný, kterým lze generovat hodnoty Y . Nyní budeme předpokládat, že Y_t se chová podle vztahu

$$(Y_t - \mu) = u_t - \theta_1 u_{t-1}, \quad (10.4)$$

kde μ je střední hodnota Y_t a u_t je bílý šum. V tom případě říkáme, že ČŘ Y_t je *proces klouzavých průměrů 1. řádu* neboli MA (1). Podle modelu (10.4) je prognóza $Y_t - \mu$ v čase

t je přímo úměrná náhodné chybě v čase $(t-1)$ prostřednictvím koeficientu úměry $-\theta_1$ plus/mínus náhodná chyba (bílý šum).

Podobně *proces klouzavých průměrů 2. řádu* neboli MA (2) má tvar

$$(Y_t - \mu) = u_t - \theta_1 u_{t-1} - \theta_2 u_{t-2}, \quad (10.5)$$

Analogicky *proces klouzavých průměrů q -tého řádu* neboli MA(q) má tvar

$$(Y_t - \mu) = u_t - \theta_1 u_{t-1} - \theta_2 u_{t-2} - \dots - \theta_q u_{t-q}. \quad (10.6)$$

Jednoduše řečeno, proces klouzavých průměrů je lineární kombinací minulých náhodných chyb bílého šumu. Na rozdíl od AR procesů jsou procesy MA (q) pro všechna $q \geq 1$ stacionární nezávisle na hodnotách koeficientů θ_i .

10.2.3 AUTOREGRESIVNÍ PROCES KLOUZAVÝCH PRŮMĚRŮ (ARMA)

Časová řada, která má charakteristiky jak AR, tak MA procesů, je ARMA proces. Konkrétně ARMA proces 1. řádu, tj. ARMA (1,1) má tvar

$$Y_t = \delta + \varphi_1 Y_{t-1} + u_t - \theta_1 u_{t-1}, \quad (10.7)$$

kde δ je konstantní člen. Analogicky můžete uvažovat procesy ARMA (p, q), které mají p autoregresivních a q klouzavých členů. Vzhledem ke stacionaritě procesu MA(q) je podmínka stacionarity procesu ARMA (p, q) totožná s podmínkou stacionarity procesu AR(p). Jinak řečeno, proces ARMA (p, q) je stacionární, právě když je stacionární proces AR(p).

10.2.4 AUTOREGRESIVNÍ A INTEGROVANÝ PROCES KLOUZAVÝCH PRŮMĚRŮ (ARIMA)

Časové procesy, které jste doposud poznali, byly vesměs za určitých podmínek stacionární. Dobře však víte, že mnohé ekonomické časové řady jsou nestacionární. Říkáme, že časová řada Y_t , tj. *stochastický proces Y_t je integrovaný 1. řádu* neboli je to I (1) proces, jestliže 1. diference této časové řady je stacionární. Jinak řečeno, ČŘ Y_t je integrovaná 1. řádu, jestliže $\Delta Y_t = Y_t - Y_{t-1}$ je stacionární ČŘ. Analogicky lze zavést pojem integrované časové řady d -tého řádu, jestliže d -tá diference této ČŘ je stacionární neboli $\Delta^d Y_t = \Delta^{d-1} Y_t - \Delta^{d-1} Y_{t-1}$ je stacionární, přitom $\Delta^1 = \Delta$. Stacionární proces se této symbolice označuje jako I(0) proces.

Proto když nejprve proces d -krát diferencujeme a poté obdržíme ARMA (p, q) proces, nazývá se původní proces ARIMA (p, d, q). V tomto symbolickém vyjádření znamenají např. ARIMA ($p, 0, q$) a ARMA (p, q) stejný proces, stejně tak ARIMA (0, 0, q) = MA (q), ARIMA ($p, 0, 0$) = AR (p), ARMA ($p, 0$) = AR (p), apod.

10.3 Box – Jenkinsova metodologie prognózování časových řad

Představte si, že máte analyzovat nějakou časovou řadu, jako třeba čtvrtletní HDP ČR. Jak zjistíte, o který typ procesu se jedná? Jde o realizaci AR procesu, nebo snad MA procesu, či jejich kombinaci ARMA? Může být konkrétní časová řada realizací více různých typů procesu, např. jak AR (1), tak současně MA (1)? V této souvislosti hledáme model časové řady a hned je třeba říci, že konkrétní časová řada může mít několik „správných“ modelů.

Box-Jenkinsova metodologie, známá také jako ARIMA (AutoRegressive Integrated Moving Average) modelování, je přístup používaný pro analýzu, modelování a prognózování časových řad. Tato metodologie byla vyvinuta v 60. letech Georgeem E.P. Boxem a Gwilymem M. Jenkinsem. Je široce využívána v oblasti statistiky a ekonometrie pro práci s neperiodickými časovými řadami.

ARIMA modelování kombinuje tři základní komponenty:

AutoRegressive (AR) složka: Tato složka zahrnuje autoregresní členy, což znamená, že hodnota časové řady v daném okamžiku závisí na předchozích hodnotách řady. Autoregresní modely zachycují korelaci mezi aktuální hodnotou a jejími minulými hodnotami.

Integrated (I) složka: Integrovaná složka zahrnuje diferencování dat, což může pomoci přeměnit nestacionární časovou řadu na stacionární. Diference odstraňují trend a sezónní složky, čímž zjednodušují analýzu.

Moving Average (MA) složka: Tato složka zahrnuje klouzavý průměr reziduí, což jsou odchylky mezi aktuální hodnotou a hodnotou předpovězenou autoregresní částí modelu.

ARIMA model se tedy označuje jako ARIMA (p, d, q), kde:

- p značí řád autoregresní složky (počet předchozích hodnot zahrnutých do modelu),
- d značí stupeň diferencování potřebný k dosažení stacionarity,
- q značí řád klouzavého průměru složky (počet reziduí zahrnutých do modelu).

Samotný proces Box-Jenkinsovy metodologie zahrnuje několik kroků:

1. **Identifikace modelu:** Na základě analýzy časové řady se pokoušíme identifikovat potenciální hodnoty p , d a q pro ARIMA model. Využívá se tvarů funkcí ACF a PACF.
2. **Odhad parametrů:** Následuje odhad parametrů modelu na základě historických dat.
3. **Kontrola reziduí:** Provádíme analýzu reziduí modelu, abychom ověřili, zda jsou náhodně rozložená. Pokud rezidua nejsou náhodně rozložená, může být třeba provést další úpravy modelu.

4. **Prognóza:** Nakonec použijeme model k prognózování budoucích hodnot časové řady.

Je důležité mít na paměti, že Box-Jenkinsova metodologie vyžaduje určitý stupeň odbornosti a zkušeností v oblasti časových řad a statistiky, ačkoli existují i softwary, které mohou asistovat v procesu identifikace a odhadu modelu.

Aplikaci jednotlivých kroků s využitím programu GRETL si ukážeme na konkrétním příkladu v závěru této kapitoly. Ještě předtím se seznámíte s dalšími nástroji a metodami, které se využívají v prvním kroku při identifikaci modelu ČŘ.

Významným nástrojem ke stanovení typu modelu (AR, MA, I, ARMA, ARIMA) je autokorelační funkce ρ_k , $k = 1, 2, \dots$, (ACF) a korelogram, resp. výběrová autokorelační funkce $\hat{\rho}_k$, $k = 1, 2, \dots$, a výběrový korelogram. Korelace mezi 2 náhodnými veličinami je často způsobena tím, že obě tyto veličiny jsou korelovány s veličinou třetí. Velká část korelace mezi veličinami Y_t a Y_{t-k} může být zapříčiněna jejich korelací s mezilehlými veličinami $Y_{t-1}, Y_{t-2}, \dots, Y_{t-k+1}$. Pojem parciální autokorelace zachycuje korelaci mezi veličinami Y_t a Y_{t-k} očištěnou o vliv veličin mezi nimi. *Parciální autokorelační koeficient* ρ_{kk} , $k = 0, 1, 2, \dots$, (2 indexy kk) je analogií k pojmu parciální regresní koeficient. Uvažujte k -násobnou lineární regresi Y_t s regresory $Y_{t-1}, Y_{t-2}, \dots, Y_{t-k}$

$$Y_t = \rho_{k1} Y_{t-1} + \rho_{k2} Y_{t-2} + \dots + \rho_{kk} Y_{t-k} + e_t. \quad (10.8)$$

Regresní koeficient ρ_{kk} je ve (10.8) právě parciální autokorelační koeficient. Vztahu (10.8) se využívá k výpočtu *výběrového parciálního autokorelačního koeficientu* $\hat{\rho}_{kk}$, viz Arlt (1999).

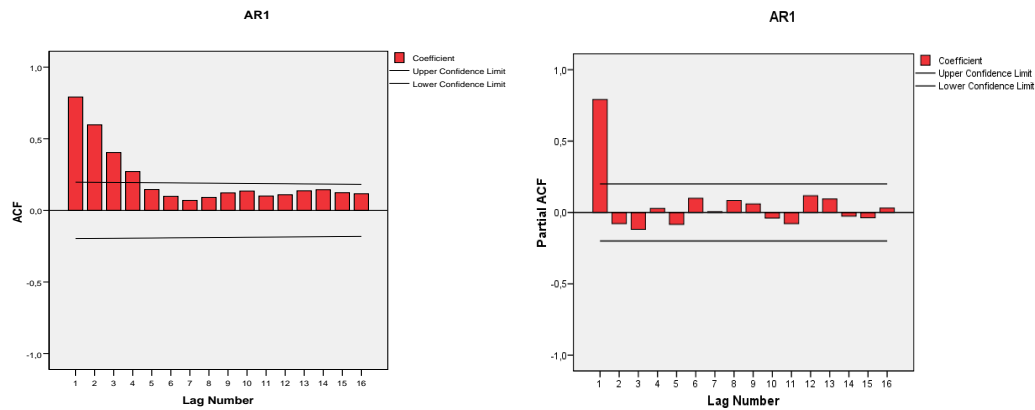
Důležitou roli hraje tzv. *parciální autokorelační funkce* (PACF) stochastického procesu ρ_{kk} pro $k = 0, 1, 2, \dots$, PACF má následující vlastnosti:

$$\begin{aligned} \rho_{00} &= 1, -1 \leq \rho_{kk} \leq 1 \text{ pro } k = 1, 2, \dots \\ \rho_{kk} &= \rho_{-k,-k} \text{ pro } k = 1, 2, \dots, \text{ tj. PACF je symetrická kolem } k = 0. \end{aligned}$$

Grafickým znázorněním PACF je *parciální korelogram*. Vzhledem k uvedeným vlastnostem stačí, aby parciální korelogram zobrazoval hodnoty pro posuvy $k > 0$.

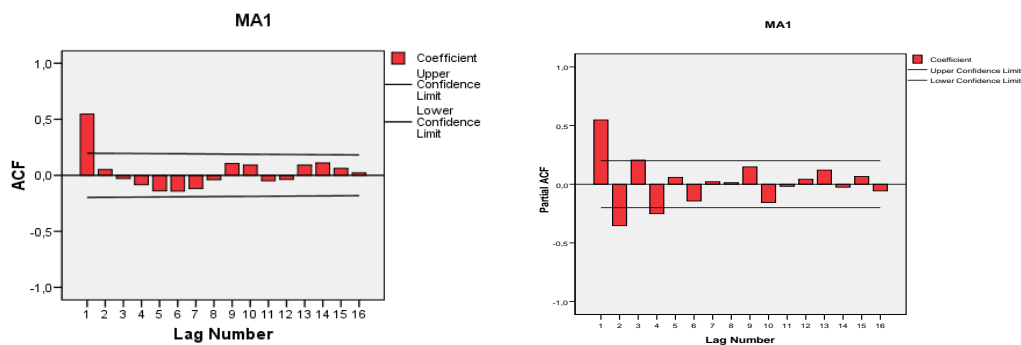
Při identifikaci typu procesu ARIMA a jeho řádů využíváme charakteristických tvarů ACF a PACF. Různé typy procesů ARIMA mají charakteristické tvary korelogramů a parciálních korelogramů. V programu GRETL využíváme nabídku: Proměnná \rightarrow Korelogram. Jednotlivé typy procesů mají následující charakteristiky:

- a. **Proces AR(p):** Prvních p hodnot PACF je „velkých“, další = 0 a „rychlý“ pokles (v absolutních hodnotách) ACF.



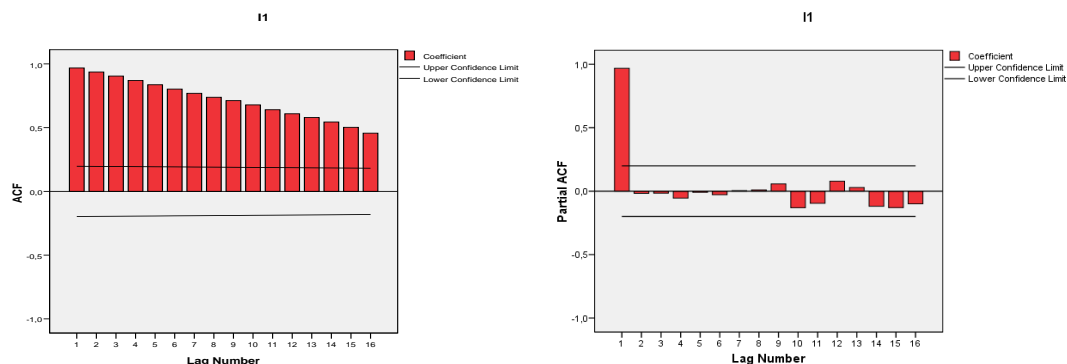
Obrázek 38: Příklady korelogramů AR (1)

- b. Proces $MA(q)$: Prvních q hodnot ACF je „velkých“, další = 0 a „rychlý“ pokles (v absolutních hodnotách) PACF.



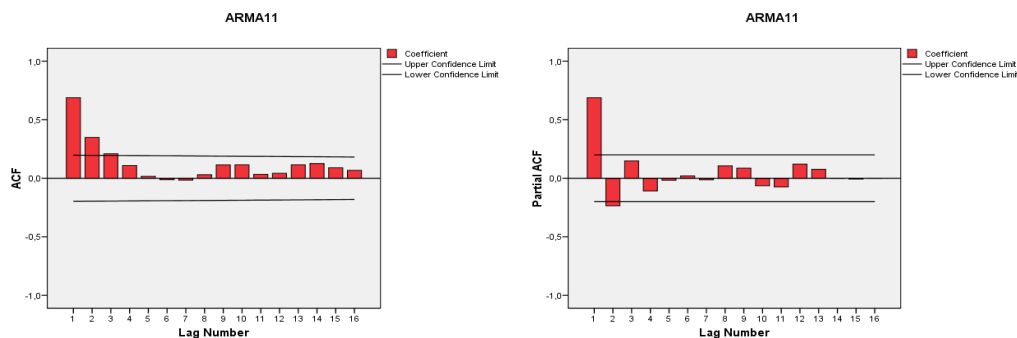
Obrázek 39: Příklady korelogramů MA (1)

- c. Proces $I(d)$: „Pomalý“ pokles ACF, prvních d hodnot PACF je „velkých“, další = 0.



Obrázek 40: Příklady korelogramů I (1): „Náhodná procházka“

- d. Proces $ARMA(p, q)$: Prvních q hodnot ACF je „velkých“, další = 0 a prvních p hodnot PACF je „velkých“, další = 0.



Obrázek 41: Příklady korelogramů ARMA (1,1)

ŘEŠENÁ ÚLOHA 10.1



Uvažujte časovou řadu „Čtvrtletní HDP České republiky“ v mil. Kč (zdroj Český statistický úřad). Hodnoty časové řady jsou uvedeny v následující Tabulce 14 a zobrazeny v grafu na Obrázku 42.

Tabulka 14: HDP ČR v mil. Kč v letech 2005–2023

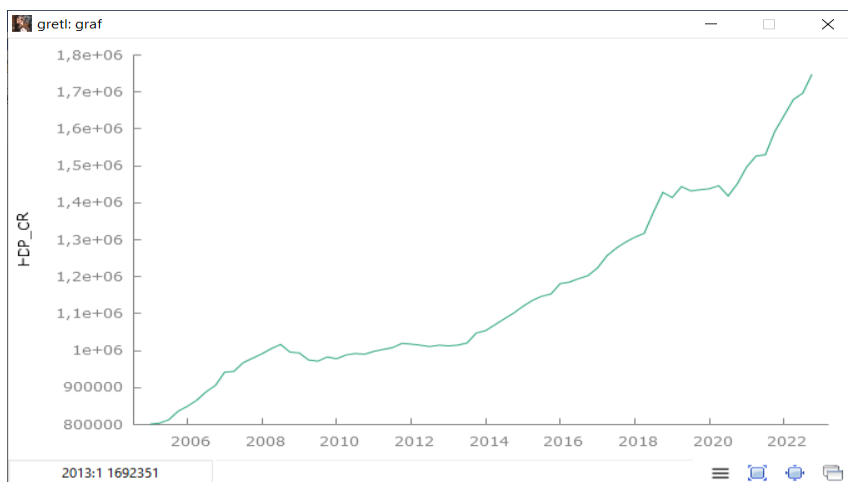
1Q/2005	801 486	1Q/2010	978 222	1Q/2015	1 119 947	1Q/2020	1 438 206
2Q/2005	803 815	2Q/2010	988 348	2Q/2015	1 135 893	2Q/2020	1 446 467
3Q/2005	813 013	3Q/2010	992 387	3Q/2015	1 147 260	3Q/2020	1 418 529
4Q/2005	836 035	4Q/2010	990 361	4Q/2015	1 153 639	4Q/2020	1 451 908
1Q/2006	849 444	1Q/2011	998 308	1Q/2016	1 181 683	1Q/2021	1 497 225
2Q/2006	865 904	2Q/2011	1 003 298	2Q/2016	1 185 584	2Q/2021	1 526 632
3Q/2006	888 564	3Q/2011	1 008 390	3Q/2016	1 195 161	3Q/2021	1 530 072
4Q/2006	905 632	4Q/2011	1 019 675	4Q/2016	1 203 410	4Q/2021	1 592 640
1Q/2007	942 051	1Q/2012	1 017 859	1Q/2017	1 224 225	1Q/2022	1 635 908
2Q/2007	944 126	2Q/2012	1 014 883	2Q/2017	1 256 656	2Q/2022	1 679 421
3Q/2007	967 813	3Q/2012	1 011 265	3Q/2017	1 277 280	3Q/2022	1 696 463
4Q/2007	979 829	4Q/2012	1 014 942	4Q/2017	1 293 581	4Q/2022	1 748 070
1Q/2008	991 805	1Q/2013	1 013 154	1Q/2018	1 306 933		
2Q/2008	1 005 646	2Q/2013	1 015 088	2Q/2018	1 317 350		
3Q/2008	1 017 012	3Q/2013	1 020 879	3Q/2018	1 374 997		
4Q/2008	996 313	4Q/2013	1 047 879	4Q/2018	1 428 474		
1Q/2009	993 972	1Q/2014	1 054 375	1Q/2019	1 414 353		
2Q/2009	974 788	2Q/2014	1 070 196	2Q/2019	1 443 994		
3Q/2009	971 644	3Q/2014	1 086 133	3Q/2019	1 432 287		
4Q/2009	983 089	4Q/2014	1 101 830	4Q/2019	1 435 590		

Najděte vhodný ARIMA model této časové řady a pomocí něj prognózuje čtvrtletní hodnoty HDP až do konce roku 2024.

Řešení:

K řešení využijeme Box-Jenkinsovu metodologii prognózování ČŘ formulovanou ve 4 krocích popsanych v subkapitole 10.3. Použijeme k tomu statistický program GRET. Nejprve sestrojíme graf časové řady HDP. Označíme proměnnou HDP a v nabídce ZOBRAZIT → Vyskreslit zadané proměnné → Vykreslit časové řady... a dostaneme graf zobrazený na Obrázku 42.

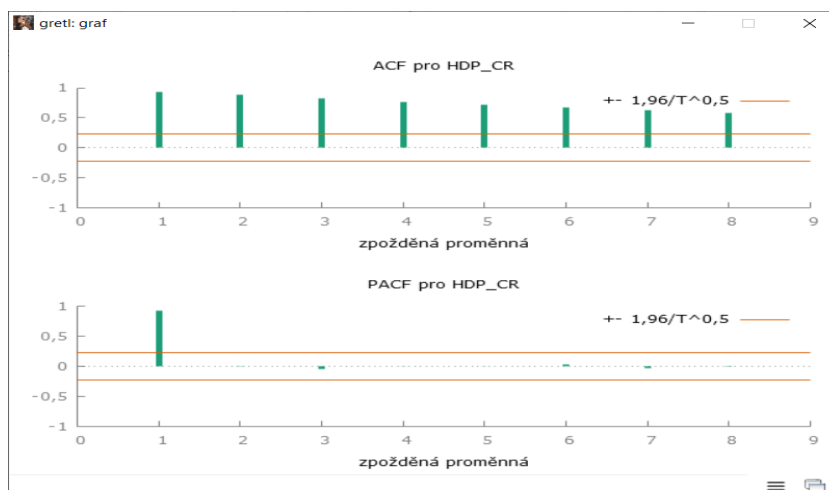
Z prostého pohledu na spojnicový graf na Obrázku 42 lze usoudit, že se jedná o nestacionární časovou řadu. Tento předpoklad potvrdíme analýzou korelogramů ACF a PACF.



Obrázek 42: HDP ČR v mil. Kč v letech 2005–2023

Krok 1: Identifikace modelu procesu ARIMA.

V menu: PROMĚNNÁ → KORELOGRAM → maximální počet zpoždění = 8, a ve výstupu obdržíme korelogramy, které zachycuje Obrázek 43.

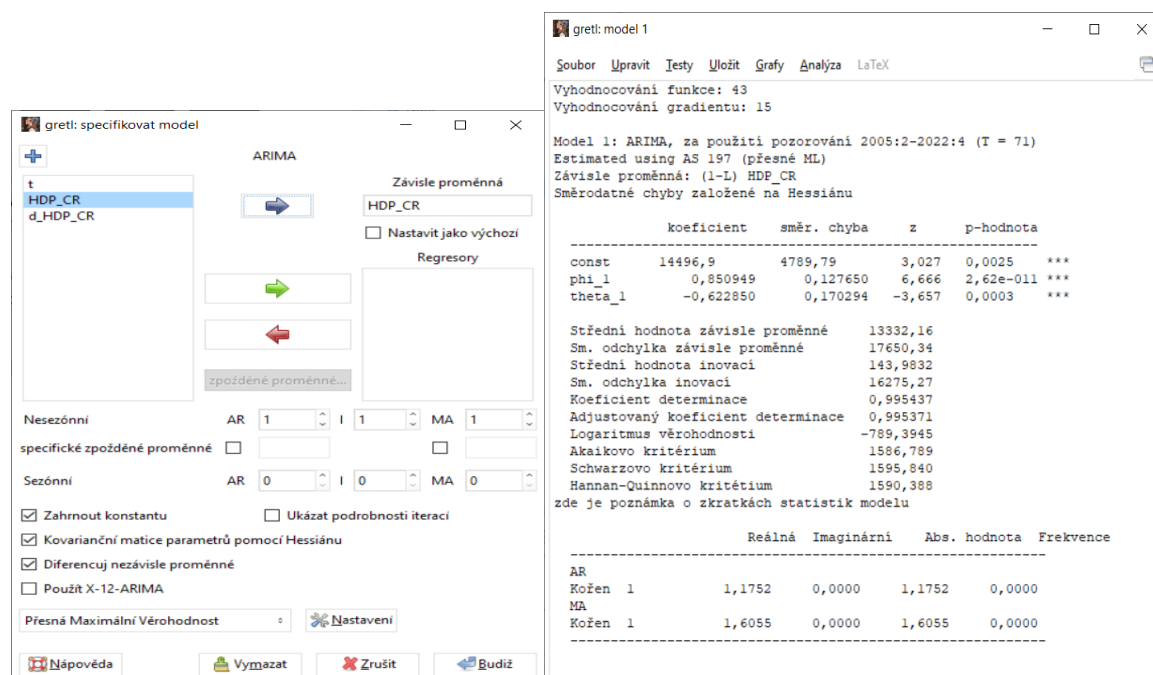


Obrázek 43: Korelogramy HDP ČR

V korelogramu hodnoty ACF pomalu klesají, v PACF je „velká“ první hodnota. Z toho vyvozujeme, že se jedná o nestacionaritu 1. řádu, tj. typu I (1). Stacionarizujeme proto časovou řadu jedním diferencováním. Dále vypočteme diferencované hodnoty proměnné

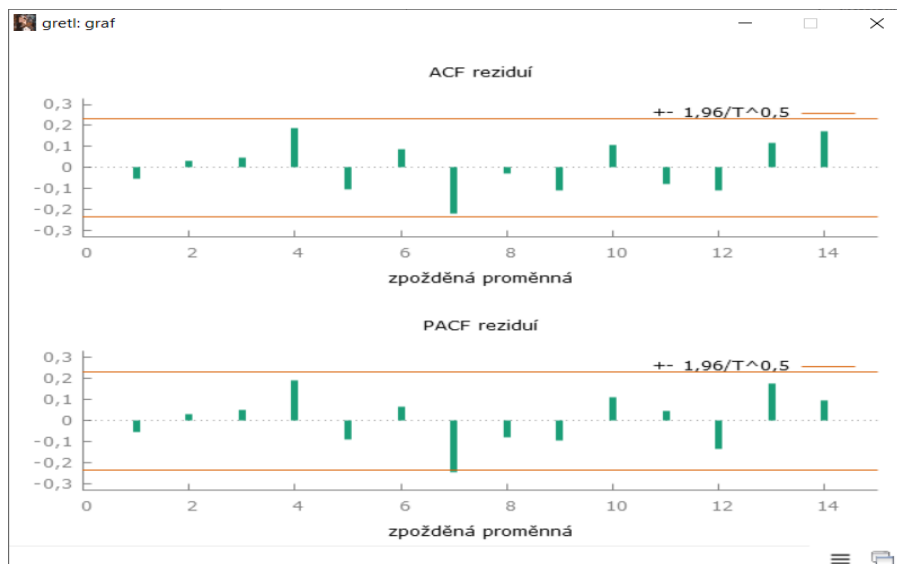
HDP, označíme proměnnou HDP a dále v menu vybereme PŘIDAT → První diference vybraných proměnných. Následně sestrojíme opět korelogramy této diferencované proměnné a na základě tvarů ACF a PACF vybereme model ARIMA (1, 1, 1).

Krok 2: Odhad parametru modelu – výpočet koeficientů provedeme v menu: MODEL → Univariate time series → ARIMA, Závisle proměnná: HDP viz. Obrázek 34, který ukazuje zadání modelu ARIMA (1, 1, 1) a jeho výstup, což je odhad koeficientů. Na Obrázku 34 vidíme, že koeficienty jsou statisticky významné, hodnota koeficientu determinace je 0,99.



Obrázek 34: Zadání modelu ARIMA modelu (1, 1, 1) a odhad parametrů

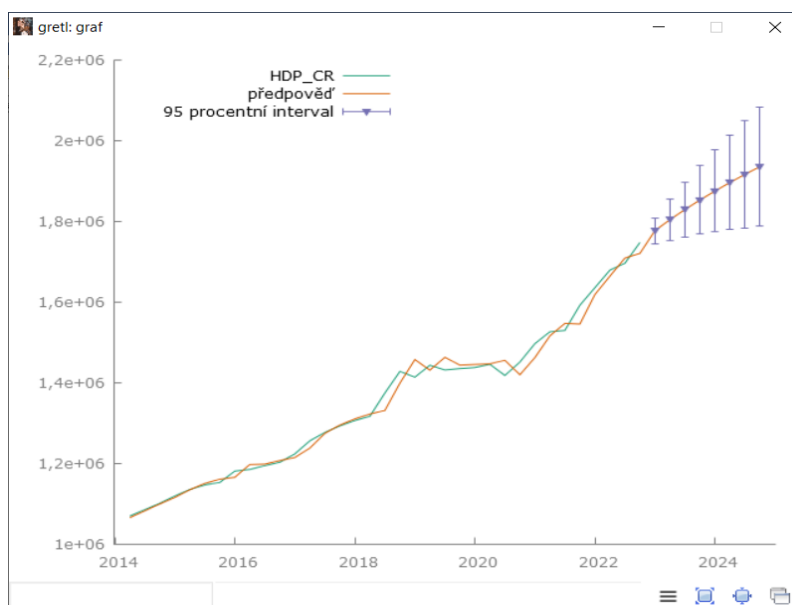
Krok 3: Verifikace modelu – spočívá v ověření předpokladu, že reziduum je bílým šumem. Ve výstupu modelu vybereme v menu GRAFY → Korelogram reziduí → potvrdíme a dostáváme Obrázek 35.



Obrázek 35: Korelogramy reziduí časové řady HDP ČR

Uvedené korelogramy potvrzují, že ACF i PACF jsou nulové,

Krok 4: Prognózu odhadneme do konce roku 2024. Výsledky ukazuje Graf 36, který zobrazuje hodnoty původní časové řady a hodnoty modelované časové řady. Graf sestrojíme tak, že v menu modelu vybereme ANALÝZA → Předpovědi.



Obrázek 36: Grafické zobrazení původní a odhadnuté časové řady



SAMOSTATNÉ ÚKOLY

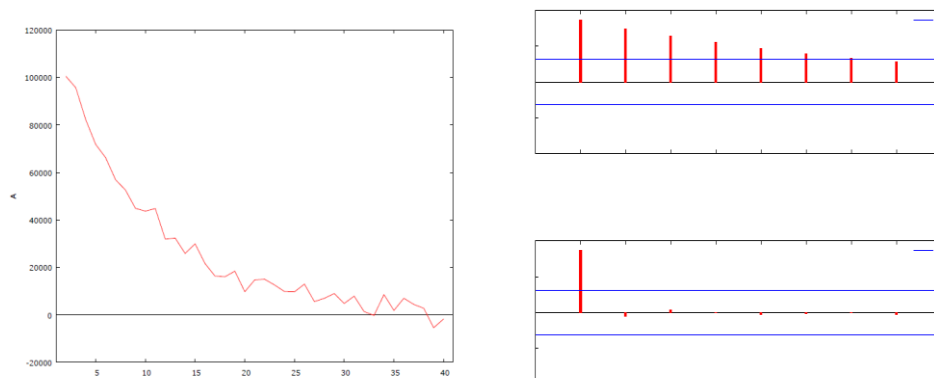
10.1 Uvažujte časovou řadu počtu vyrobených součástek v tis. ks v letech 2005-2022. Hodnoty časové řady jsou uvedeny v následující tabulce. Najděte vhodný ARIMA model této

časové řady a pomocí něj prognózuje čtvrtletní hodnoty až do konce roku 2023. Použijte přitom 4 kroky Box-Jenkinsovy metodologie.

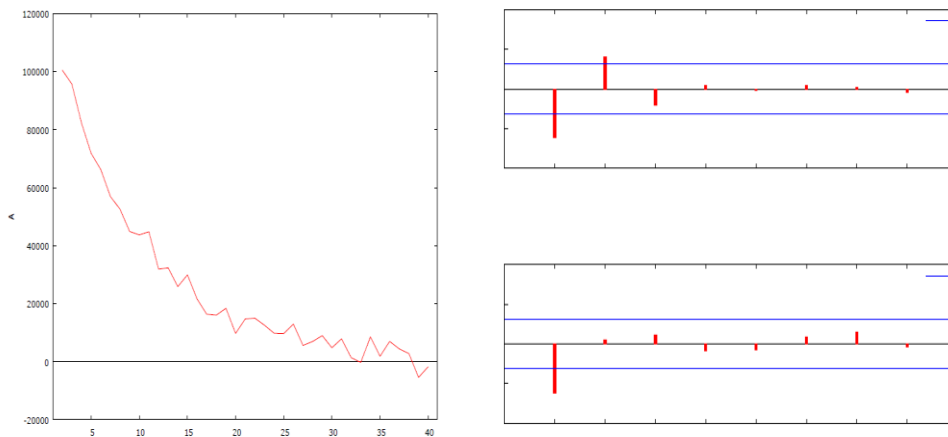
1Q/2005	2872,8	1Q/2010	3154	1Q/2015	3830,8	1Q/2020	4221,8
2Q/2005	2860,3	2Q/2010	3190,4	2Q/2015	3732,6	2Q/2020	4254,8
3Q/2005	2896,6	3Q/2010	3249,9	3Q/2015	3733,5	3Q/2020	4309
4Q/2005	2873,7	4Q/2010	3292,5	4Q/2015	3808,5	4Q/2020	4333,5
1Q/2006	2942,9	1Q/2011	3356,7	1Q/2016	3860,5	1Q/2021	4390,5
2Q/2006	2947,4	2Q/2011	3369,2	2Q/2016	3844,4	2Q/2021	4387,7
3Q/2006	2966	3Q/2011	3381	3Q/2016	3864,5	3Q/2021	4412,6
4Q/2006	2980,8	4Q/2011	3416,3	4Q/2016	3803,1	4Q/2021	4427,1
1Q/2007	2927,3	1Q/2012	3466,4	1Q/2017	3756,1	1Q/2022	4460
2Q/2007	3089,7	2Q/2012	3525	2Q/2017	3771,1	2Q/2022	4515,3
3Q/2007	3125,8	3Q/2012	3574,4	3Q/2017	3754,4	3Q/2022	4559,3
4Q/2007	3175,5	4Q/2012	3567,2	4Q/2017	3759,6	4Q/2022	4625,5
1Q/2008	3253,3	1Q/2013	3591,8	1Q/2018	3783,5		
2Q/2008	3267,6	2Q/2013	3707	2Q/2018	3886,5		
3Q/2008	3264,3	3Q/2013	3735,6	3Q/2018	3944,4		
4Q/2008	3289,1	4Q/2013	3779,6	4Q/2018	4012,1		
1Q/2009	3259,4	1Q/2014	3780,8	1Q/2019	4089,5		
2Q/2009	3267,6	2Q/2014	3784,3	2Q/2019	4144		
3Q/2009	3239,1	3Q/2014	3807,5	3Q/2019	4166,4		
4Q/2009	3226,4	4Q/2014	3814,6	4Q/2019	4194,2		

10.2 Z následujících grafů časových řad se pokuste určit stacionaritu a z korelogramů určete identifikační body, od kterých se hodnoty již statisticky významně neliší od nuly, a identifikujte řád autoregresního procesu AR (p) a řád procesu klouzavých průměrů MA (q).

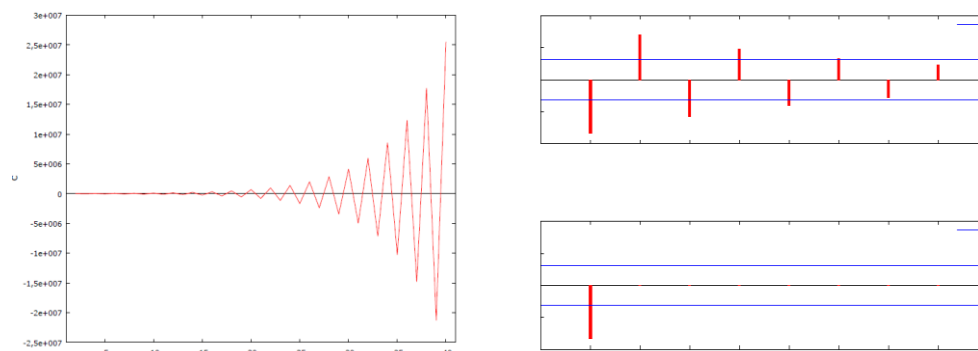
a)



b)

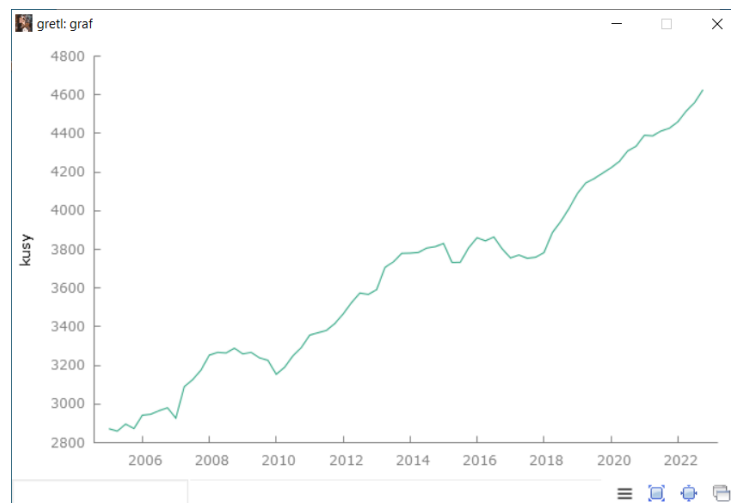


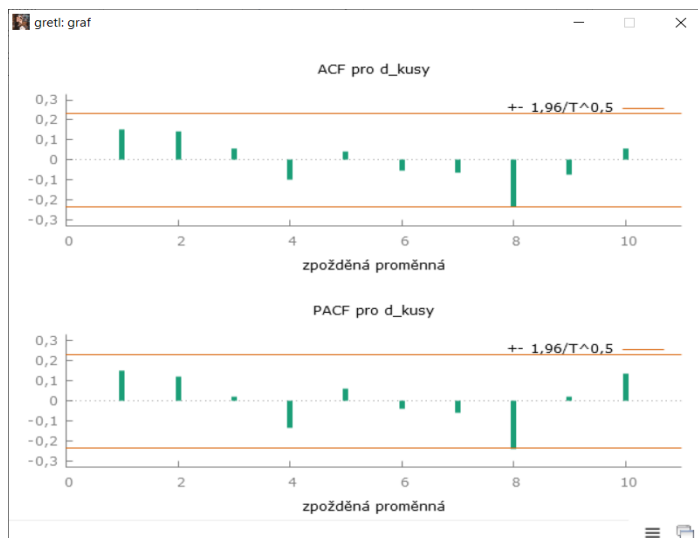
c)



ODPOVĚDI

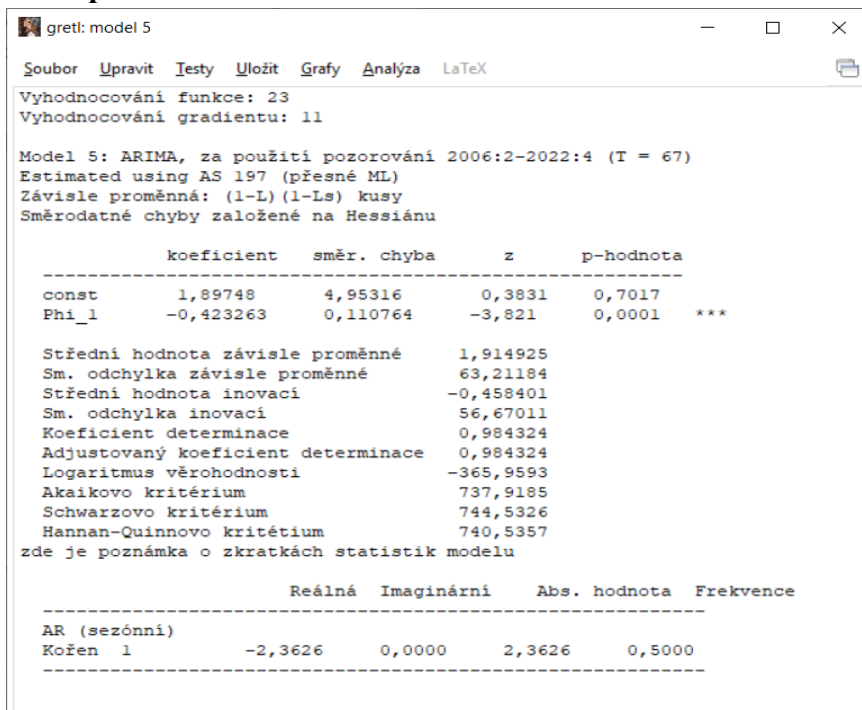
10.1 a) Identifikace modelu





Na základě tvaru korelačních funkcí diferencované časové řady vybíráme model ARIMA (0, 1, 0) x (1, 1, 0).

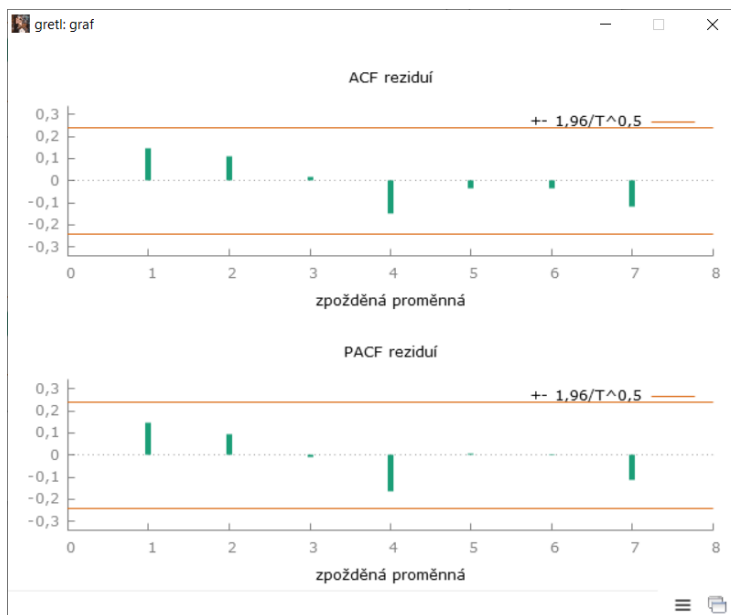
1. b) Odhad parametrů modelu



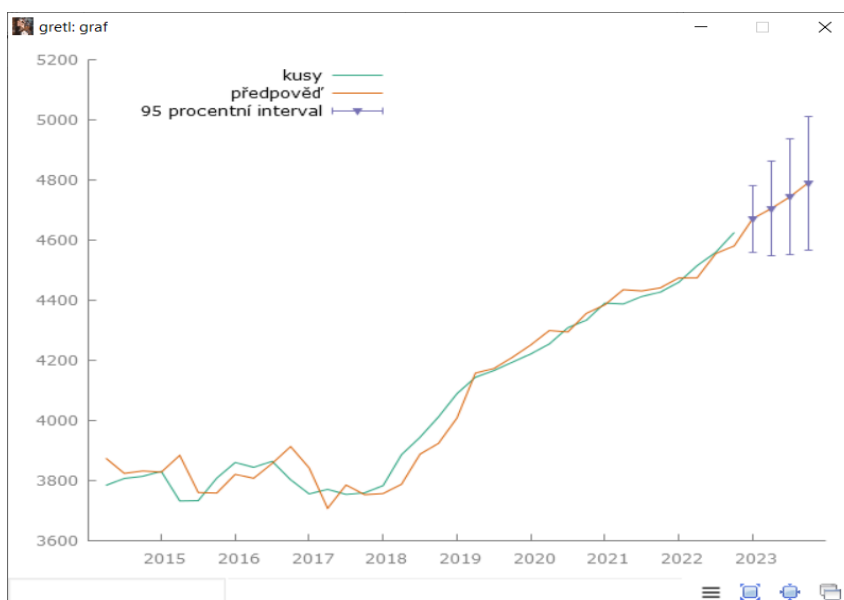
Koeficient SAR1 = -0,423 je statisticky významný na hladině významnosti 0,01 (protože hodnota signifikance = 0,000 je menší než 0,01).

c) Verifikace modelu

Korelogramy potvrzují, že ACF i PACF reziduální složky jsou nulové.



d) Predikce počtu vyrobených výrobků do 4. čtvrtletí 2023



Období	Bodový odhad	Intervalový odhad (95%)	
Q1/ 2023	4671,3	4560,2	4782,4
Q2/2023	4704,7	4547,6	4861,8
Q3/2023	4743,3	4550,9	4935,7
Q4/2023	4790,3	4568,2	5012,5

10.2

a) Časová řada nebude stacionární, musíme odstranit exponenciální trend. Počáteční model AR (1) – ACF exponenciálně klesá, PACF má významnou pouze první hodnotu. Model ARIMA (1, 1, 0).

b) Časová řada bude stacionární, má nulovou střední hodnotu, a rozptyl neroste. Počáteční modely AR (1) – ACF exponenciálně klesá, PACF má významnou pouze první hodnotu, MA (2) – PACF exponenciálně klesá, ACF má významné pouze první dvě hodnoty, model ARIMA (1, 0, 1).

c) Časová řada není stacionární. Počáteční model AR (1) – ACF exponenciálně klesá, PACF má významnou pouze první hodnotu. Model ARIMA (1, 1, 0).

SHRNUTÍ KAPITOLY



V této závěrečné kapitole jste se seznámili s časovými řadami typu ARIMA. Box-Jenkinsova metodologie, která se touto problematikou zabývá, klade důraz nikoliv na konstrukci jednorovnicového nebo víceroznicového modelu, jak tomu bylo např. v regresní analýze, nýbrž na analýzu vlastních stochastických vlastností ekonomických ČŘ. Postupně jste se seznámili s vlastnostmi autoregresivních procesů AR, procesů pohyblivých průměrů MA, integračních procesů I, jakož i procesů vzniklých jejich kombinací ARIMA. Dále byly tyto procesy rozšířeny též na sezónní procesy. Úkolem pak bylo pro konkrétní časovou řadu nalézt vhodný konkrétní model typu ARIMA a nalezený model použít pro účely prognózy (predikce, extrapolace) hodnot dané časové řady. Celý postup tvorby prognózy ČŘ autoři metody ARIMA formulovali ve 4 krocích, které nazýváme Box-Jenkinsova metodologie prognózování ČŘ. Jednotlivé kroky jsou (1) Identifikace modelu, (2) Odhad modelu, (3) Verifikace modelu a (4) Prognóza pomocí modelu. Jednotlivé kroky Box-Jenkinsovy metodologie byly ilustrovány na příkladu časové řady čtvrtletního HDP České republiky s pomocí statistického programu GRETL.

LITERATURA

- ANDĚL, Jiří, 2007. *Statistické metody*. 4. upr. vyd. Praha: Marfyzpress, 299 s. ISBN 80-7378-003-8.
- ARLT, Josef, 1999. *Moderní metody modelování ekonomických časových řad*. 1.vyd. Praha: Grada Publishing, 307 s. ISBN 80-716-9539-4.
- CIPRA, Tomáš, 1986. *Analýza časových řad s aplikacemi v ekonomii*. 1.vyd. Praha: Státní nakladatelství technické literatury, 246 s.
- GUJARATI, Damodar N, c2003. *Basic econometrics*. 4th ed. Boston: McGraw-Hill, xxix, 1002 s. ISBN 978-0-07-233542-2.
- HÁTLE, Jaroslav a LIKEŠ, Jiří, 1974. *Základy počtu pravděpodobnosti a matematické statistiky*. 2. vyd. Praha: SNTL. 463 s.
- HINDLS, Richard, SEGER, Jan a HRONOVÁ, Stanislava, 2002. *Statistika pro ekonomy*. 1. vyd. Praha: Professional Publishing, 415 s. ISBN 80-864-1926-6.
- KAŇKA, Miloš, 1998. *Vybrané partie z matematiky pro ekonomy*. 1.vyd. Praha: VŠE, 231 s. ISBN 80-707-9537-9.
- MAREK, Luboš a kol., 2007. *Statistika pro ekonomy: aplikace*. 2. vyd. Praha: Professional Publishing. 485 s. ISBN 978-80-86946-40-5.
- RAMÍK, Jaroslav a Šárka ČEMERKOVÁ, 2000. *Statistika A*. Vyd. 3., rozš. a upr. V Opavě: Slezská univerzita, Obchodně podnikatelská fakulta v Karviné, 162 s. ISBN 80-7248-097-9.
- RAMÍK, Jaroslav a Šárka ČEMERKOVÁ, 2000. *Statistika B*. Vyd. 2., rozš. a upr. V Opavě: Slezská univerzita, Obchodně podnikatelská fakulta v Karviné, 143 s. ISBN 80-724-8099-5.
- RAMÍK, Jaroslav a Šárka ČEMERKOVÁ, 2003. *Kvantitativní metody B: statistika*. Vyd. 1. Karviná: Slezská univerzita v Opavě, Obchodně podnikatelská fakulta v Karviné, 206 s. ISBN 80-724-8198-3.
- SEGER, Jan, HRONOVÁ, Stanislava a HINDLS, Richard, 1998. *Statistika v hospodářství*. 1.vyd. Praha: ETC Publishing, 636 s. ISBN 80-860-0656-5.

SHRnutí STUDIjNÍ OPORY

Tento text slouží jako pomocný materiál pro studium všech akreditovaných magisterských programů na Slezské univerzitě, konkrétně na Obchodně podnikatelské fakultě v Karviné. Předmět Statistické zpracování dat navazuje na bakalářský předmět Statistika, který se vyučuje na SU OPF, nebo na podobný předmět základů statistiky na bakalářské úrovni na jiných ekonomických fakultách v České republice. Tento text představuje inovaci oproti původnímu studijnímu materiálu. V rámci tohoto předmětu je klíčový důraz kladen na praktické využití statistických metod při zpracování ekonomických dat v oblastech aplikované ekonomie, zejména v oblastech marketingu a managementu.

Tato studijní opora umožňuje studentům plnohodnotnou a současně samostatnou studijní práci. Tento materiál je rozdělen do deseti tematických kapitol.

Vysokoškolské studium tohoto předmětu, Statistické zpracování dat, vyžaduje od studentů značné úsilí věnované pravidelnosti a trpělivosti při studiu a samostudiu, schopnost soustředění na téma, aktivní přístup, který zahrnuje samostatné řešení úloh. Tato studijní opora by měla studentům pomoci v těchto oblastech. Dalšími doplňkovými zdroji pro studium mohou být tradiční učebnice, skripta a doporučená literatura.

PŘEHLED DOSTUPNÝCH IKON



Čas potřebný ke studiu



Klíčová slova



Průvodce studiem



Rychlý náhled



Tutoriály



K zapamatování



Řešená úloha



Kontrolní otázka



Odpovědi



Samostatný úkol



Pro zájemce



Cíle kapitoly



Nezapomeňte na odpočinek



Průvodce textem



Shrnutí



Definice



Případová studie



Věta



Korespondenční úkol



Otázky



Další zdroje



Úkol k zamyšlení

Název: **Statistické zpracování dat**

Autor: **Prof. RNDr. Jaroslav Ramík, CSc., Mgr. Radmila Krkošková**

Vydavatel: Slezská univerzita v Opavě
Obchodně podnikatelská fakulta v Karviné

Určeno: studentům SU OPF Karviná

Počet stran: 18180

Tato publikace neprošla jazykovou úpravou.