

Univerzita Karlova v Praze
Matematicko-fyzikální fakulta

BAKALÁŘSKÁ PRÁCE



Jan Vávra

Testy homoskedasticity v lineárním modelu

Katedra pravděpodobnosti a matematické statistiky

Vedoucí bakalářské práce: doc. RNDr. Arnošt Komárek, Ph.D.

Studijní program: Matematika

Studijní obor: Obecná matematika

Praha 2016

Poděkování

Rád bych na tomto místě věnoval poděkování vedoucímu bakalářské práce doc. RNDr. Arnoštu Komárkovi, Ph.D. za to, že si vždy našel chvíli, aby mi poskytnul nesčetné množství odborných rad a připomínek, které dopomohly k finální podobě této práce.

Prohlašuji, že jsem tuto bakalářskou práci vypracoval samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona v platném znění, zejména skutečnost, že Univerzita Karlova v Praze má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle §60 odst. 1 autorského zákona.

V dne

Podpis autora

Název práce: Testy homoskedasticity v lineárním modelu

Autor: Jan Vávra

Katedra: Katedra pravděpodobnosti a matematické statistiky

Vedoucí bakalářské práce: doc. RNDr. Arnošt Komárek, Ph.D., Katedra pravděpodobnosti a matematické statistiky

Abstrakt: Tato práce se zabývá testováním předpokladu homoskedasticity v lineárním modelu, neboli předpokladu o konstantním rozptylu chyb tohoto modelu. Takových testů existuje celá řada, ale ne všechny se dají aplikovat na konkrétním modelu a ne všechny dosahují uspokojivých výsledků za různých okolností. Práce se zaměří na testy, které lze odvodit na základě asymptotické teorie maximální věrohodnosti, zvláště pak teorie testů s rušivými parametry. Odvozeny jsou dva základní testy, první v situaci modelu analýzy rozptylu jednoduchého třídění a druhý v situaci, kdy je připuštěna závislost rozptylu na doprovodných veličinách. V následných numerických studiích jsou prověřeny vlastnosti odvozených testových statistik.

Klíčová slova: homoskedasticita, heteroskedastický lineární model

Title: Homoscedasticity Tests in a Linear Model

Author: Jan Vávra

Department: Department of Probability and Mathematical Statistics

Supervisor: doc. RNDr. Arnošt Komárek, Ph.D., Department of Probability and Mathematical Statistics

Abstract: This thesis deals with testing the assumption of homoscedasticity in linear model, that is the assumption of constant variance of this model. There is plenty of such tests, but not all of them can be applied to specific model and not all of them reach satisfactory results under various circumstances. Thesis focuses on tests which can be derived on the basis of the asymptotic theory for maximum likelihood estimation, particularly the test theory with nuisance parameters. There are derived two basic tests, the first one in the situation of analysis of variance model and the second one in the situation when we allow the dependence of variance to concomitant quantities. In subsequent numerical studies there are examined characteristics of derived test statistics.

Keywords: homoscedasticity, heteroscedastic linear model

Obsah

Používané značení	2
Úvod	4
1 Lineární model	5
1.1 Klasický lineární model	5
1.2 Metoda nejmenších čtverců	7
1.3 Heteroskedastické modely	9
2 Maximální věrohodnost	11
2.1 Metoda maximální věrohodnosti	11
2.2 Testy s rušivými parametry	14
3 Bartlettův test	17
3.1 Test založený na věrohodnostním poměru	18
3.2 Porovnání testových statistik	24
4 Breusch-Paganův test	26
5 Numerické studie	33
5.1 Experimentální porovnání statistik B a LR	33
5.2 Breusch-Paganův test	36
Závěr	41
Literatura	42
Seznam obrázků	43
Seznam tabulek	44
Přílohy	45

Používané značení

V této práci se budeme držet následujících konvencí:

- ▷ Malými římskými a řeckými písmeny budeme označovat jednorozměrné konstanty, parametry či funkce. Písmena i, j, k, l, m, n budou označovat přirozená čísla.
- ▷ Velkými římskými písmeny budeme (až na výjimky - ε) označovat náhodné veličiny (X, Y, Z), jejich realizace potom příslušnými malými písmeny (x, y, z). Dále také v několika případech použijeme velkých římských a řeckých písmen pro označení množin (M, Ω, Θ, \dots).
- ▷ Tučnými písmeny budeme značit vektory. Každý vektor bude chápán jako sloupcový. Pro náhodné vektory budeme (až na výjimky - \mathbf{u}) užívat velká římská písmena ($\mathbf{X}, \mathbf{Y}, \mathbf{Z}$) a jejich realizace příslušnými malými písmeny ($\mathbf{x}, \mathbf{y}, \mathbf{z}$).
- ▷ Pro matice budeme používat zdvojených velkých římských písmen (\mathbb{A}, \mathbb{X}) nebo Σ . V textu budeme používat značení \mathbb{X} jak pro náhodnou matici, tak i její realizaci, je tedy třeba význam symbolu určit z kontextu. Symboly \mathbb{N} a \mathbb{R} máme vyhrazeny pro označení množiny přirozených a reálných čísel.
- ▷ Pro míry máme vyhrazený speciální font řeckých písmen (μ, ν, λ).
- ▷ Bezpatkové písmo je používáno pro značení pravděpodobnostních rozdělání ($\mathbb{N}, \text{Mult}, \text{LM}, \dots$) nebo určitých funkcionalů či operátorů ($\mathbb{E}, \text{var}, \text{Tr}, \dots$).
- ▷ Strojové písmo budeme používat pro pojmy spojené s výpočetním softwarem (`bartlett.test`, `bptest`, `...`).
- ▷ Použití stříšky či vlnky nad písmenem značí, že se jedná o odhad tohoto parametru, často je následován indexem n označující rozsah náhodného výběru. Například $\tilde{\beta}_n, \hat{\sigma}_n^2, \dots$

Speciální symboly:

- $(\Omega, \mathcal{A}, \mathbb{P})$ pravděpodobnostní prostor s množinou jevů Ω ,
 σ -algebrou \mathcal{A} a pravděpodobnostní mírou \mathbb{P}
s.j. skoro jistě vzhledem k pravděpodobnostní míře \mathbb{P}
- $[\mu]$ s.v. skoro všude/všechna vzhledem k míře μ
- λ, λ^n Lebesgueova (n -rozměrná) míra
- \mathcal{B}_0 borelovská σ -algebra na \mathbb{R}

$Y \mathbf{X}$	podmíněné rozdělení náhodné veličiny Y při \mathbf{X}
$E[Y \mathbf{X}]$	podmíněná střední hodnota náhodné veličiny Y při \mathbf{X}
$\text{var}[Y \mathbf{X}]$	podmíněný rozptyl náhodné veličiny Y při \mathbf{X}
$\text{diag}(\dots)$	diagonální matice s danými prvky na diagonále
$\text{Tr}(\mathbb{A})$	stopa matice \mathbb{A}
$\det(\mathbb{A})$	determinant matice \mathbb{A}
$\mathbf{x}^\top, \mathbb{X}^\top$	operátor transpozice vektoru \mathbf{x} a matice \mathbb{X}
\mathbb{A}^{-1}	inverzní matice k matici \mathbb{A}
$\ \mathbf{x}\ $	eukleidovská norma vektoru $\mathbf{x} = (x_1, \dots, x_n)^\top$, $\ \mathbf{x}\ = \sqrt{x_1^2 + \dots + x_n^2}$
$\mathbf{1}_n$	n -složkový sloupcový vektor, jehož prvky jsou jen 1
\mathbb{I}_n	jednotková matice řádu n , neboli $\text{diag}(1, \dots, 1)$
$\mathbf{0}^n, \mathbf{0}^{n \times 1}$	n -složkový sloupcový nulový vektor
$\mathbf{0}^{1 \times n}$	n -složkový řádkový nulový vektor
$\mathbb{O}^{n \times m}$	nulová matice řádu $n \times m$
δ_{ij}	Kroneckerovo delta, $\delta_{ij} = 1$ jedině tehdy, když $i = j$, jinak $\delta_{ij} = 0$
ℓ	logaritmická věrohodnost
\mathcal{F}_m	systém hustot generovaný parametrickým prostorem o dimenzi m
H_0, H_1	nulová hypotéza a alternativa
$\xrightarrow{\text{D}}$	konvergence v distribuci
χ_q^2	chí kvadrát rozdělení o $q \in \mathbb{N}$ stupních volnosti
$\chi_q^2(\alpha)$	α -kvantil rozdělení chí kvadrát o q stupních volnosti
R	symbol pro označení statistického výpočetního softwaru R

Úvod

Úkolem této práce bude nejdříve představit klasický lineární model, kde se předpokládá konstantní rozptyl chyb, a následně ho zobecnit na heteroskedastický lineární model, který už bude obecně předpokládat proměnlivost tohoto rozptylu. Jednodušší definice lineárního modelu předpokládá, že uvažované regresory jsou předem známé konstanty. V tomto textu ovšem budeme předpokládat, že regresory budou obecně náhodné vektory. Toto zobecnění nás tedy nutí pracovat s podmíněným rozdělením, podmíněnou střední hodnotou a podmíněným rozptylem. Dále jelikož teorie maximální věrohodnosti a z ní vycházející teorie testů s rušivými parametry tvoří základy, ze kterých budeme při odvozování testů vycházet, shrneme ve druhé kapitole stěžejní definice a poznatky těchto přístupů, na které se budeme odvolávat. Ve třetí kapitole si představíme model analýzy rozptylu jednoduchého třídění, kde budeme obecně předpokládat nejen rozdílné střední hodnoty, ale i rozptyly v jednotlivých skupinách. Tento model si zobecníme do pojmů definovaných v první kapitole a následně podrobně odvodíme test homoskedasticity poměrem věrohodnosti. Seznámíme se také s velmi podobnou statistikou, kterou již v roce 1937 navrhl anglický statistik Maurice Stevenson Bartlett a dáme ji do souvislosti s odvozeným testem poměrem věrohodností. Ve čtvrté kapitole se podíváme na heteroskedastický lineární model, kde připouštíme závislost rozptylu na doprovodných regresorech. Důkladně v tomto modelu odvodíme takzvaný skórový test, který navrhli v roce 1979 australští statistici Trevor Breusch a Adrian Pagan. V závěrečné kapitole za pomoci výpočetního prostředí **R** (R Core Team, 2016) prostudujeme vlastnosti odvozených testových statistik. Zejména nás bude zajímat, zda dodržují předepsanou hladinu a jak přesná je použitá asymptotika při malém rozsahu výběru. Dále se zaměříme na studium síly těchto testů a chování příslušných statistik při nesplnění předpokladu normality.

Kapitola 1

Lineární model

V této kapitole se seznámíme s klasickým a heteroskedastickým lineárním modelem a základními pojmy, se kterými budeme dále pracovat. Uvedeme si také základní tvrzení a věty, které z těchto definic vyplývají.

1.1 Klasický lineární model

Nejprve si zavedeme klasický lineární model. Uvažme reálný náhodný vektor (Y, \mathbf{X}) na pravděpodobnostním prostoru $(\Omega, \mathcal{A}, \mathbb{P})$ se sdruženou hustotou $f(y, \mathbf{x})$ vzhledem k σ -konečné míře $\nu = \nu_y \times \nu_x$, kde Y je náhodná veličina a $\mathbf{X} = (X_0, X_1, \dots, X_k)^\top$ je náhodný vektor dimenze $k + 1$ pro $k \in \mathbb{N}_0$. Dále budeme předpokládat, že ν_y je Lebesgueova míra λ na $(\mathbb{R}, \mathcal{B}_0)$.

Definice 1. Řekneme, že (Y, \mathbf{X}) se řídí lineárním modelem, pokud

$$\begin{aligned} \mathbb{E}[Y|\mathbf{X}] &= \mathbf{X}^\top \boldsymbol{\beta} & \text{s.j.}, \\ \text{var}[Y|\mathbf{X}] &= \sigma^2 & \text{s.j.}, \end{aligned} \tag{1.1}$$

kde $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$ a $\sigma^2 > 0$ jsou neznámé parametry, a pokud rozdělení náhodného vektoru \mathbf{X} na těchto parametrech nezávisí.

Nadále budeme psát $(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2)$.

Poznámka. Několik doplňujících poznámek k definici 1.

- Náhodné veličině Y se říká *vysvětlovaná proměnná* (nebo také závisle proměnná) a náhodnému vektoru \mathbf{X} říkáme *regresory* (nebo také nezávisle proměnné). Dále říkáme, že se jedná o *lineární model*, protože podmíněná střední hodnota Y při \mathbf{X} závisí na $\boldsymbol{\beta}$ lineárně.
- Častým předpokladem v lineárních modelech bývá, že náhodná veličina X_0 má degenerované (nenáhodné) rozdělení $X_0 = 1$ skoro jistě. Náhodný vektor \mathbf{X} jsme tedy zavedli jako $(k + 1)$ -složkový, abychom tento případ mohli odlišit ve zvláštní (nulté) složce.
- V textu budeme dále pracovat s podmíněným rozdělením $Y|\mathbf{X}$, jehož hustotu budeme značit $f_{Y|\mathbf{X}}$. Tato hustota závisí na parametrech $\boldsymbol{\beta}$ a σ^2 . Označme $f_{\mathbf{X}}$ hustotu náhodného vektoru \mathbf{X} , tato už na $\boldsymbol{\beta}$ a σ^2 nezávisí. Ze známých vlastností hustot z teorie pravděpodobnosti platí

$$f(y, \mathbf{x}; \boldsymbol{\beta}, \sigma^2) = f_{Y|\mathbf{X}}(y|\mathbf{x}; \boldsymbol{\beta}, \sigma^2) f_{\mathbf{X}}(\mathbf{x}) \quad [\nu] \text{ s.v.} \tag{1.2}$$

Definice 2. V lineárním modelu $(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2)$ zavedeme náhodnou veličinu ε jako

$$\varepsilon := Y - \mathbf{X}^\top \boldsymbol{\beta}$$

a budeme ji nazývat chybový člen modelu.

Poznámka. ε z definice 2 odpovídá náhodné veličině popisující odchylku (či chybu - „error“) vysvětlované proměnné od její podmíněné střední hodnoty. Jakožto náhodnou veličinu bychom ji tedy měli dle naší konvence značit velkým římským písmenem. Ovšem použití písmena E by mohlo snadno vést k záměně s operátorem střední hodnoty E , proto si zde dovolíme od naší konvence upustit a značit tuto veličinu právě písmenem ε .

Lemma 1. Necht' $(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2)$, potom pro ε z definice 2 platí

$$E \varepsilon = 0, \quad \text{var } \varepsilon = \sigma^2.$$

Důkaz. Oba vzorce dostaneme použitím základních vlastností podmíněné střední hodnoty a vzorce (1.1) z definice 1 lineárního modelu.

Pro střední hodnotu náhodné veličiny ε platí

$$E \varepsilon = E(E[\varepsilon | \mathbf{X}]) = E(E[Y | \mathbf{X}] - E[\mathbf{X}^\top \boldsymbol{\beta} | \mathbf{X}]) = E(\mathbf{X}^\top \boldsymbol{\beta} - \mathbf{X}^\top \boldsymbol{\beta}) = 0.$$

Pro rozptyl ε si nejprve připomeneme vzorec (1.3) pro \mathbf{U}, \mathbf{V} obecné náhodné vektory. Platí, že

$$\text{var } \mathbf{U} = E(\text{var}[\mathbf{U} | \mathbf{V}]) + \text{var}(E[\mathbf{U} | \mathbf{V}]). \quad (1.3)$$

Tento vzorec aplikujeme na $U = \varepsilon$ a $\mathbf{V} = \mathbf{X}$. Už v minulém výpočtu jsme došli ke zjištění, že

$$E[\varepsilon | \mathbf{X}] = 0 \quad \text{s.j.,}$$

proto druhý člen ve vzorci (1.3) bude také nulový. Zaměřme se proto na výpočet $\text{var}[\varepsilon | \mathbf{X}]$, pro který platí

$$\begin{aligned} \text{var}[\varepsilon | \mathbf{X}] &= \text{var}[Y - \mathbf{X}^\top \boldsymbol{\beta} | \mathbf{X}] \\ &= E\left[(Y - \mathbf{X}^\top \boldsymbol{\beta} - E[Y - \mathbf{X}^\top \boldsymbol{\beta} | \mathbf{X}])^2 | \mathbf{X}\right] \\ &= E\left[(Y - E[Y | \mathbf{X}])^2 | \mathbf{X}\right] \\ &= \text{var}[Y | \mathbf{X}] = \sigma^2 \quad \text{s.j.} \end{aligned}$$

Odtud již plyne, že $\text{var } \varepsilon = E(\text{var}[\varepsilon | \mathbf{X}]) = E\sigma^2 = \sigma^2$.

□

Standardní situací bývá, že parametry $\boldsymbol{\beta}$ a σ^2 jsou neznámé a naším cílem je odhadnout tyto parametry na základě naměřených dat, které reprezentujeme náhodným výběrem z příslušného rozdělení. Dále si definujeme řadu užitečných charakteristik, které nám přiblíží chování našeho modelu.

Mějme náhodný výběr $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$ z rozdělení

$$(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2).$$

Zavedme si náhodný vektor $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ a náhodnou matici \mathbb{X} s řádky $\mathbf{X}_i^\top, i \in \{1, \dots, n\}$. Pak z nezávislosti a definice 1 plyne, že

$$\begin{aligned} \mathbf{E}[\mathbf{Y}|\mathbb{X}] &= \mathbb{X}\boldsymbol{\beta} \quad \text{s.j.}, \\ \text{var}[\mathbf{Y}|\mathbb{X}] &= \sigma^2\mathbb{I}_n \quad \text{s.j.}, \end{aligned}$$

kde \mathbb{I}_n značí jednotkovou matici řádu n .

Definice 3. Budeme říkat, že (\mathbf{Y}, \mathbb{X}) se řídí lineárním modelem, a zapisovat to budeme také jako

$$(\mathbf{Y}, \mathbb{X}) \sim \text{LM}(\mathbb{X}\boldsymbol{\beta}, \sigma^2\mathbb{I}_n).$$

Danou realizaci náhodného výběru budeme značit (y_i, \mathbf{x}_i) pro $i \in \{1, \dots, n\}$ a dále $\mathbf{y} = (y_1, \dots, y_n)^\top$. Symbolem \mathbb{X} budeme rozumět jak onu realizaci

$$\mathbb{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix},$$

tak i náhodnou matici \mathbb{X} . V celém textu budeme pro jednoduchost předpokládat, že $n > k + 1$ a že tato matice má skoro jistě plnou sloupcovou hodnotu $k + 1$, kterou budeme značit r , tj. $r = k + 1$.

Náhodný výběr z definice 3 pak má dle (1.2) sdruženou hustotu vzhledem k součinnové míře $\nu^n = \nu \times \dots \times \nu$ tvaru

$$\begin{aligned} f_n(\mathbf{y}, \mathbb{X}; \boldsymbol{\beta}, \sigma^2) &= \prod_{i=1}^n f(y_i, \mathbf{x}_i; \boldsymbol{\beta}, \sigma^2) = \prod_{i=1}^n f_{Y|X}(y_i|\mathbf{x}_i; \boldsymbol{\beta}, \sigma^2) f_X(\mathbf{x}_i) \\ &= f_{\mathbf{Y}|\mathbb{X}}(\mathbf{y}|\mathbb{X}; \boldsymbol{\beta}, \sigma^2) f_{\mathbb{X}}(\mathbb{X}) \quad [\nu^n] \text{ s.v.}, \end{aligned} \quad (1.4)$$

kde $f_{\mathbf{Y}|\mathbb{X}}$ značí hustotu sdruženého podmíněného rozdělení $\mathbf{Y}|\mathbb{X}$ a $f_{\mathbb{X}}$ značí sdruženou hustotu náhodné matice \mathbb{X} .

Obdobně si také zavedeme vektor chybových členů modelu $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$, kde $\varepsilon_i = Y_i - \mathbf{X}_i^\top \boldsymbol{\beta}$ pro každé $i \in \{1, \dots, n\}$, tedy $\boldsymbol{\varepsilon} = (\mathbf{Y} - \mathbb{X}\boldsymbol{\beta})$. Potom pro něj dle důkazu lemmatu 1 platí

$$\mathbf{E}[\boldsymbol{\varepsilon}|\mathbb{X}] = 0 \quad \text{s.j.}, \quad \text{var}[\boldsymbol{\varepsilon}|\mathbb{X}] = \sigma^2\mathbb{I}_n \quad \text{s.j.}$$

1.2 Metoda nejmenších čtverců

Nyní můžeme přistoupit k odhadu vektorového parametru $\boldsymbol{\beta}$. Standardní metodou pro nalezení tohoto odhadu je *metoda nejmenších čtverců*. Součtem čtverců rozumíme funkci $SS(\boldsymbol{\beta})$ při daných (\mathbf{Y}, \mathbb{X}) definovanou jako

$$SS(\boldsymbol{\beta}) = \|\mathbf{Y} - \mathbb{X}\boldsymbol{\beta}\|^2 = (\mathbf{Y} - \mathbb{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbb{X}\boldsymbol{\beta}) = \sum_{i=1}^n (Y_i - \mathbf{X}_i^\top \boldsymbol{\beta})^2.$$

Odhadem \mathbf{b} vektorového parametru $\boldsymbol{\beta}$ metodou nejmenších čtverců potom rozumíme takovou hodnotu $\boldsymbol{\beta}$, která minimalizuje funkci $SS(\boldsymbol{\beta})$, to jest

$$\mathbf{b} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{k+1}} SS(\boldsymbol{\beta}).$$

Za našeho předpokladu o plné sloupcové hodnosti matice \mathbb{X} lze toto \mathbf{b} jednoznačně určit jako řešení tzv. *soustavy normálních rovnic*

$$\mathbb{X}^\top \mathbb{X} \mathbf{b} = \mathbb{X}^\top \mathbf{Y} \quad \text{ve tvaru} \quad \mathbf{b} = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y}.$$

Definice 4. V modelu $(\mathbf{Y}, \mathbb{X}) \sim \text{LM}(\mathbb{X}\boldsymbol{\beta}, \sigma^2 \mathbb{I}_n)$ definujeme

- vektor vyrovnaných hodnot $\hat{\mathbf{Y}} := \mathbb{X} \mathbf{b} = \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y}$,
- vektor reziduí $\mathbf{u} := \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbb{I}_n - \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top) \mathbf{Y}$,
- reziduální součet čtverců $RSS := SS(\mathbf{b}) = \|\mathbf{u}\|^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$,
- reziduální rozptyl $S^2 := \frac{RSS}{n-r}$, kde $r = k + 1$ je hodnost matice \mathbb{X} .

Poznámka. $\mathbf{Y} - \hat{\mathbf{Y}}$ je opět náhodný vektor, tedy bychom měli užívat pro označení reziduí velkého písmene \mathbf{U} . Zde se držíme značení zavedeného v učebnici Zvára (2008), ze které primárně vycházíme.

Poznámka. Vektor $\hat{\mathbf{Y}}$ je nestranným odhadem vektoru $\mathbb{X}\boldsymbol{\beta}$, neboť

$$\mathbb{E} [\hat{\mathbf{Y}} | \mathbb{X}] = \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbb{E} [\mathbf{Y} | \mathbb{X}] = \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbb{X} \boldsymbol{\beta} = \mathbb{X} \boldsymbol{\beta} \quad \text{s.j.}$$

Pro podmíněný rozptyl tohoto odhadu platí

$$\text{var} [\hat{\mathbf{Y}} | \mathbb{X}] = \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top (\sigma^2 \mathbb{I}_n) \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top = \sigma^2 \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \quad \text{s.j.}$$

Reziduální součet čtverců odpovídá dle této definice součtu čtverců vzdáleností odhadů \hat{Y}_i od Y_i , tedy nám vypovídá o tom, jak moc dobře nám $\hat{\mathbf{Y}}$ aproximuje \mathbf{Y} . Jeho vztah k podmíněnému rozptylu σ^2 je popsán v následujícím lemmatu.

Lemma 2. V modelu $\mathbf{Y} \sim \text{LM}(\mathbb{X}\boldsymbol{\beta}, \sigma^2 \mathbb{I}_n)$, kde daná matice \mathbb{X} má hodnost r , platí

$$\mathbb{E} [RSS | \mathbb{X}] = (n-r)\sigma^2 \quad \text{s.j.}, \quad \mathbb{E} [S^2 | \mathbb{X}] = \sigma^2 \quad \text{s.j.}$$

Důkaz. Důkaz je analogický důkazu věty 2.2. v práci Zvára (2008), uvedeme si však několik úprav, které je v důkazu nutné provést.

Označíme-li si matice $\mathbb{H} = \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top$ a $\mathbb{M} = \mathbb{I}_n - \mathbb{H}$, pak tyto matice jsou symetrické a idempotentní se stopami po řadě r a $(n-r)$ a platí

$$\hat{\mathbf{Y}} = \mathbb{H}\mathbf{Y}, \quad \mathbf{u} = \mathbb{M}\mathbf{Y} = \mathbb{M}\boldsymbol{\varepsilon}.$$

Mimochodem se potom dá z těchto vyjádření získat

$$\mathbb{E} [\mathbf{u} | \mathbb{X}] = \mathbf{0} \quad \text{s.j.}, \quad \text{var} [\mathbf{u} | \mathbb{X}] = \sigma^2 \mathbb{M} \quad \text{s.j.}$$

Chceme získat podmíněnou střední hodnotu RSS při \mathbb{X} , k tomu využijeme ekvivalentních zápisů

$$RSS = \|\mathbf{u}\|^2 = \mathbf{u}^\top \mathbf{u} = \boldsymbol{\varepsilon}^\top \mathbb{M}^\top \mathbb{M} \boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}^\top \mathbb{M} \boldsymbol{\varepsilon}.$$

Nyní už se jen využije vlastností funkcionálu stopy matice Tr , abychom mohli provést následující úpravy

$$\begin{aligned} \mathbb{E}[RSS|\mathbb{X}] &= \mathbb{E}[\boldsymbol{\varepsilon}^\top \mathbb{M} \boldsymbol{\varepsilon} | \mathbb{X}] = \mathbb{E}[\text{Tr}(\boldsymbol{\varepsilon}^\top \mathbb{M} \boldsymbol{\varepsilon}) | \mathbb{X}] = \text{Tr}(\mathbb{M} \mathbb{E}[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top | \mathbb{X}]) \\ &= \text{Tr}(\mathbb{M} \text{var}[\boldsymbol{\varepsilon} | \mathbb{X}]) = \text{Tr}(\mathbb{M} \sigma^2 \mathbb{I}_n) = \sigma^2 \text{Tr}(\mathbb{M}) = (n-r)\sigma^2 \quad \text{s.j.} \end{aligned}$$

Odtud již snadno plyne i druhé tvrzení. \square

Při práci s lineárními modely se často předpokládá normalita podmíněného rozdělení $Y|\mathbf{X}$.

Definice 5. Řekneme, že (Y, \mathbf{X}) se řídí normálním lineárním modelem, jestliže $Y|\mathbf{X}$ má normální rozdělení s parametrem střední hodnoty $\mathbf{X}^\top \boldsymbol{\beta}$ a rozptylem σ^2 . Píšeme $(Y, \mathbf{X}) \sim \text{NLM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2)$.

Máme-li náhodný výběr $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$ z takového rozdělení (Y, \mathbf{X}) , pak říkáme, že se (\mathbf{Y}, \mathbb{X}) řídí normálním lineárním modelem a píšeme $(\mathbf{Y}, \mathbb{X}) \sim \text{NLM}(\mathbb{X}^\top \boldsymbol{\beta}, \sigma^2 \mathbb{I}_n)$.

Podmíněná hustota $f_{Y|\mathbf{X}}$ vzhledem k Lebesgueově míře λ je tedy tvaru

$$f_{Y|\mathbf{X}}(y|\mathbf{x}; \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y - \mathbf{x}^\top \boldsymbol{\beta})^2}{2\sigma^2}\right\}.$$

Podmíněná hustota $f_{\mathbf{Y}|\mathbb{X}}$ rozdělení $(\mathbf{Y}|\mathbb{X})$ vzhledem k Lebesgueově n -rozměrné míře λ^n je potom rovna

$$f_{\mathbf{Y}|\mathbb{X}}(\mathbf{y}|\mathbb{X}; \boldsymbol{\beta}, \sigma^2) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2\right\}.$$

Věta 3. V modelu $(\mathbf{Y}, \mathbb{X}) \sim \text{NLM}(\mathbb{X}^\top \boldsymbol{\beta}, \sigma^2 \mathbb{I}_n)$ platí

$$\left(\frac{1}{\sigma^2} RSS \middle| \mathbb{X}\right) = \left(\frac{1}{\sigma^2} \|\mathbf{u}\|^2 \middle| \mathbb{X}\right) \sim \chi_{n-r}^2,$$

kde r je hodnota matice \mathbb{X} .

Důkaz. Lze provést analogicky jako v práci Zvára (2008, Věta 2.6.). \square

1.3 Heteroskedastické modely

Doposud jsme pracovali s klasickým tzv. *homoskedastickým* modelem, kde $\text{var}[Y|\mathbf{X}]$ je rovno skoro jistě konstantě. Tento předpoklad nám značně usnadňuje práci. Obecně je však nutno uvážit situaci, kdy toto není nutně pravda a rozptyl závisí na regresorech. Naším cílem je tedy nalézt nějaké testy, které by dokázaly určit, zda naše data jsou v souladu s předpokladem *homoskedasticity* nebo naopak tento předpoklad silně vyvracejí.

Existuje řada obecných způsobů, jak v lineárním modelu zavést heteroskedasticitu. V této práci budeme říkat, že náhodný vektor (Y, \mathbf{X}) se řídí heteroskedastickým lineárním modelem, pokud platí

$$\begin{aligned} \mathbb{E}[Y|\mathbf{X}] &= \mathbf{X}^\top \boldsymbol{\beta} & \text{s.j.}, \\ \text{var}[Y|\mathbf{X}] &= \sigma^2 h(\mathbf{X}; \boldsymbol{\tau}) & \text{s.j.}, \end{aligned} \tag{1.5}$$

kde $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$, $\sigma^2 > 0$ a $\boldsymbol{\tau} \in T \subset \mathbb{R}^q$, $q \in \mathbb{N}$ jsou neznámé parametry, h je kladná měřitelná funkce, a pokud rozdělení náhodného vektoru \mathbf{X} na parametrech $\boldsymbol{\beta}$, σ^2 a $\boldsymbol{\tau}$ nezávisí. Obecně by se dalo předpokládat, že funkce h může navíc záviset na $\boldsymbol{\beta}$, tuto situaci však zkoumat nebudeme.

Nastane-li tato situace, tak budeme analogicky jako v klasickém lineárním modelu psát

$$(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2 h(\mathbf{X}; \boldsymbol{\tau})) \quad \text{anebo} \quad (Y, \mathbf{X}) \sim \text{NLM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2 h(\mathbf{X}; \boldsymbol{\tau}))$$

v případě, že podmíněné rozdělení $(Y|\mathbf{X})$ je normální. Pro hustotu vzhledem k míře ν tohoto rozdělení platí

$$f(y, \mathbf{x}; \boldsymbol{\tau}, \boldsymbol{\beta}, \sigma^2) = f_{Y|\mathbf{X}}(y|\mathbf{x}; \boldsymbol{\tau}, \boldsymbol{\beta}, \sigma^2) \cdot f_{\mathbf{X}}(\mathbf{x}) \quad [\nu] \text{ s.v.}$$

Když budeme opět uvažovat náhodný výběr $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$ z tohoto rozdělení, tak ponecháváme značení \mathbf{Y} a \mathbb{X} a budeme říkat, že (\mathbf{Y}, \mathbb{X}) se řídí heteroskedasticitním modelem s podmíněnou rozptylovou maticí

$$\boldsymbol{\Sigma}(\mathbb{X}) := \text{var}[\mathbf{Y}|\mathbb{X}] = \sigma^2 \text{diag}(h(\mathbf{X}_1; \boldsymbol{\tau}), \dots, h(\mathbf{X}_n; \boldsymbol{\tau})) \quad \text{s.j.}$$

a psát budeme $(\mathbf{Y}, \mathbb{X}) \sim \text{LM}(\mathbb{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\mathbb{X}))$. V případě, že podmíněné rozdělení $(\mathbf{Y}|\mathbb{X})$ je normální, budeme analogicky psát $(\mathbf{Y}, \mathbb{X}) \sim \text{NLM}(\mathbb{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\mathbb{X}))$.

Vzhledem k tomu, že funkce $h(\mathbf{X}_i; \boldsymbol{\tau})$, $i \in \{1, \dots, n\}$ jsou kladné, tak existují matice

$$\begin{aligned} \boldsymbol{\Sigma}^{-1}(\mathbb{X}) &= \frac{1}{\sigma^2} \text{diag}\left(\frac{1}{h(\mathbf{X}_1; \boldsymbol{\tau})}, \dots, \frac{1}{h(\mathbf{X}_n; \boldsymbol{\tau})}\right), \\ \boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbb{X}) &= \frac{1}{\sigma} \text{diag}\left(\frac{1}{\sqrt{h(\mathbf{X}_1; \boldsymbol{\tau})}}, \dots, \frac{1}{\sqrt{h(\mathbf{X}_n; \boldsymbol{\tau})}}\right), \\ \boldsymbol{\Sigma}^{\frac{1}{2}}(\mathbb{X}) &= \sigma \text{diag}\left(\sqrt{h(\mathbf{X}_1; \boldsymbol{\tau})}, \dots, \sqrt{h(\mathbf{X}_n; \boldsymbol{\tau})}\right). \end{aligned}$$

Lemma 4. *Nechť $(\mathbf{Y}, \mathbb{X}) \sim \text{LM}(\mathbb{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\mathbb{X}))$. Potom*

$$\mathbf{Y}^* = \boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbb{X})\mathbf{Y} \sim \text{LM}\left(\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbb{X})\mathbb{X}\boldsymbol{\beta}, \mathbb{I}_n\right).$$

Důkaz. Jelikož je $\mathbf{Y}^* = \boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbb{X})\mathbf{Y} = \left(\frac{Y_1}{\sigma\sqrt{h(\mathbf{X}_1; \boldsymbol{\tau})}}, \dots, \frac{Y_n}{\sigma\sqrt{h(\mathbf{X}_n; \boldsymbol{\tau})}}\right)^\top$ a pro každé $i \in \{1, \dots, n\}$ platí

$$\begin{aligned} \mathbb{E}\left[\frac{Y_i}{\sigma\sqrt{h(\mathbf{X}_i; \boldsymbol{\tau})}} \middle| \mathbf{X}_i\right] &= \frac{\mathbb{E}[Y_i | \mathbf{X}_i]}{\sigma\sqrt{h(\mathbf{X}_i; \boldsymbol{\tau})}} = \frac{\mathbf{X}_i^\top \boldsymbol{\beta}}{\sigma\sqrt{h(\mathbf{X}_i; \boldsymbol{\tau})}} \quad \text{s.j.}, \\ \text{var}\left[\frac{Y_i}{\sigma\sqrt{h(\mathbf{X}_i; \boldsymbol{\tau})}} \middle| \mathbf{X}_i\right] &= \frac{\text{var}[Y_i | \mathbf{X}_i]}{\sigma^2 h(\mathbf{X}_i; \boldsymbol{\tau})} = \frac{\sigma^2 h(\mathbf{X}_i; \boldsymbol{\tau})}{\sigma^2 h(\mathbf{X}_i; \boldsymbol{\tau})} = 1 \quad \text{s.j.}, \end{aligned}$$

tak platí tvrzení lemmatu. □

Na vektor \mathbf{Y}^* tedy můžeme aplikovat poznatky z podkapitol 1.1 a 1.2 a pomocí tohoto lemmatu je převést do heteroskedastického modelu.

Kapitola 2

Maximální věrohodnost

Důležitou roli při odvozování celé řady testů homoskedasticity hraje teorie maximální věrohodnosti. V této kapitole si představíme základní pojmy z této teorie a následně si ukážeme, jak se aplikuje v heteroskedastickém lineárním modelu.

2.1 Metoda maximální věrohodnosti

Definice 6. Mějme $\mathbb{V} = (\mathbf{V}_1, \dots, \mathbf{V}_n)^\top$ náhodnou matici, kde $\mathbf{V}_1, \dots, \mathbf{V}_n$ jsou nezávislé stejně rozdělené k -složkové náhodné vektory s hustotou $f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})$ vůči σ -konečné míře μ a $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^m$ je parametr. Pro pevnou realizaci $\mathbb{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)^\top$ nazveme funkci $L(\boldsymbol{\theta}) := \prod_{i=1}^n f_{\mathbf{V}}(\mathbf{v}_i; \boldsymbol{\theta})$ věrohodnostní funkcí.

Definice 7. Hodnota $\hat{\boldsymbol{\theta}}_n \in \Theta$, která maximalizuje věrohodnostní funkci $L(\boldsymbol{\theta})$ pro danou realizaci náhodné matice \mathbb{V} , se nazývá maximálně věrohodný odhad.

Při hledání maximálně věrohodného odhadu se využívá toho, že libovolná ryze rostoucí transformace věrohodnostní funkce je maximalizována ve stejném bodě. Nejčastěji se využívá logaritmická transformace. Potom při dané realizaci \mathbb{V} nazveme *logaritmickou věrohodností* funkci

$$\ell(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta}).$$

Tato funkce se dá nadále chápat jako náhodná veličina, uvažujeme-li ji také jako funkci obecné náhodné matice \mathbb{V} .

Existují-li derivace věrohodnostní funkce podle všech složek parametru $\boldsymbol{\theta}$ pro každé $\boldsymbol{\theta} \in \Theta$, tak se za dalších předpokladů (které nám zajistí existenci a maximalitu řešení) dá $\hat{\boldsymbol{\theta}}_n$ nalézt jako řešení takzvaného *systému věrohodnostních rovnic*:

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_j} = 0 \quad \forall j \in \{1, \dots, m\}.$$

Definice 8. Nechť \mathbf{V} je k -složkový náhodný vektor s hustotou $f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})$ vzhledem k σ -konečné míře μ a $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^m$. A nechť dále platí podmínky regularity:

(R1) Θ je neprázdná otevřená množina.

(R2) $M = \{\mathbf{v} \in \mathbb{R}^k : f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}) > 0\}$ nezávisí na parametru $\boldsymbol{\theta}$.

(R3) Pro $[\mu]$ skoro všechna $\mathbf{v} \in M$ a pro každé $j \in \{1, \dots, m\}$ existují

$$f'_j(\mathbf{v}; \boldsymbol{\theta}) := \frac{\partial f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})}{\partial \theta_j}.$$

(R4) Pro všechna $\boldsymbol{\theta} \in \Theta$ a všechna $j \in \{1, \dots, m\}$ platí $\int_M f'_j(\mathbf{v}; \boldsymbol{\theta}) d\mu(\mathbf{v}) = 0$.

(R5) Pro všechny dvojice $i, j \in \{1, \dots, m\}$ existuje konečný integrál

$$J_{ij}(\boldsymbol{\theta}) := \int_M \frac{f'_i(\mathbf{v}; \boldsymbol{\theta})}{f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})} \frac{f'_j(\mathbf{v}; \boldsymbol{\theta})}{f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})} f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}) d\mu(\mathbf{v}).$$

(R6) Matice $\mathbb{J}(\boldsymbol{\theta}) = (J_{ij}(\boldsymbol{\theta}))_{i,j=1}^m$ je pozitivně definitní pro každé $\boldsymbol{\theta} \in \Theta$.

Potom se systém hustot $\mathcal{F}_m = \{f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ nazývá regulární.

Matice $\mathbb{J}(\boldsymbol{\theta})$ se nazývá Fisherova informační matice pro náhodný vektor \mathbf{V} .

Věta 5. Nechť systém hustot $\mathcal{F}_m = \{f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ je regulární. Předpokládejme, že pro $[\mu]$ skoro všechna $\mathbf{v} \in M$ existují derivace

$$f''_{ij}(\mathbf{v}; \boldsymbol{\theta}) = \frac{\partial^2 f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}, \quad i, j \in \{1, \dots, m\},$$

a že pro všechna $\boldsymbol{\theta} \in \Theta$ platí

$$\int_M f''_{ij}(\mathbf{v}; \boldsymbol{\theta}) d\mu(\mathbf{v}) = 0, \quad i, j \in \{1, \dots, m\}.$$

Pak jednotlivé prvky Fisherovy informační matice lze počítat jako

$$J_{ij}(\boldsymbol{\theta}) = - \int_M \frac{\partial^2 \log f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}) d\mu(\mathbf{v}),$$

kde $i, j \in \{1, \dots, m\}$.

Důkaz. Důkaz je analogický důkazu verze věty s jednorozměrným parametrem (viz Anděl, 2007, Věta 7.19). \square

Máme-li náhodný výběr $\mathbb{V} = (\mathbf{V}_1, \dots, \mathbf{V}_n)^\top$, pak Fisherova informační matice náhodné matice \mathbb{V} , kterou označíme $\mathbb{J}^{(n)}(\boldsymbol{\theta})$, se dá zřejmě počítat jako

$$\mathbb{J}^{(n)}(\boldsymbol{\theta}) = n\mathbb{J}(\boldsymbol{\theta}).$$

Podle věty 5 lze potom její jednotlivé prvky počítat jako

$$J_{ij}^{(n)}(\boldsymbol{\theta}) = \mathbb{E} \left[- \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right], \quad i, j \in \{1, \dots, m\}.$$

Příklad (Maximální věrohodnost v lineárním modelu). Nyní si ukážeme, jak aplikovat maximální věrohodnost v situaci heteroskedastického lineárního modelu. Nechť tedy $(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2 h(\mathbf{X}; \boldsymbol{\tau}))$, viz (1.5).

Označme tedy $\boldsymbol{\theta} = (\tau_1, \dots, \tau_q, \beta_0, \dots, \beta_k, \sigma^2)^\top$ a $\Theta = T \times \mathbb{R}^{k+1} \times (0, \infty)$. Hustota náhodného vektoru vzhledem k σ -konečné míře \mathbf{v} má tvar

$$f(y, \mathbf{x}; \boldsymbol{\theta}) = f_{Y|\mathbf{X}}(y|\mathbf{x}; \boldsymbol{\theta}) \cdot f_{\mathbf{X}}(\mathbf{x}) \quad [\mathbf{v}] \text{ s.v.}$$

Pak pokud systém hustot

$$\mathcal{F}_{q+k+2} = \{f(y, \mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$$

je regulární, tak funkce

$$\log f(y, \mathbf{x}; \boldsymbol{\theta}) = \log f_{Y|\mathbf{X}}(y|\mathbf{x}; \boldsymbol{\theta}) + \log f_{\mathbf{X}}(\mathbf{x}) \quad [\mathbf{v}] \text{ s.v.},$$

a její derivace podle libovolného $\theta_j, j \in \{1, \dots, q+k+2\}$ je rovna

$$\frac{\partial \log f(y, \mathbf{x}, \boldsymbol{\theta})}{\partial \theta_j} = \frac{\partial \log f_{Y|\mathbf{X}}(y|\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_j} \quad [\mathbf{v}] \text{ s.v.},$$

neboť funkce $\log f_{\mathbf{X}}(\mathbf{x})$ na parametru $\boldsymbol{\theta}$ vůbec nezávisí.

Mějme náhodný výběr $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$ z rozdělení

$$(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2 h(\mathbf{X}; \boldsymbol{\tau})).$$

Potom logaritmickou věrohodností je dle (1.4) funkce

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= \log f_n(\mathbf{Y}, \mathbb{X}; \boldsymbol{\theta}) = \log f_{\mathbf{Y}|\mathbb{X}}(\mathbf{Y}|\mathbb{X}; \boldsymbol{\theta}) + \log f_{\mathbb{X}}(\mathbb{X}) \\ &= \sum_{i=1}^n \log f_{Y|\mathbf{X}}(Y_i|\mathbf{X}_i; \boldsymbol{\theta}) + \sum_{i=1}^n \log f_{\mathbf{X}}(\mathbf{X}_i) \quad [\mathbf{v}^n] \text{ s.v.} \end{aligned}$$

Pro její derivaci dle proměnné $\theta_j, j \in \{1, \dots, q+k+2\}$ platí

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_j} = \sum_{i=1}^n \frac{\partial \log f_{Y|\mathbf{X}}(Y_i|\mathbf{X}_i; \boldsymbol{\theta})}{\partial \theta_j} \quad [\mathbf{v}^n] \text{ s.v.}$$

Tedy maximálně věrohodný odhad $\hat{\boldsymbol{\theta}}_n$ parametru $\boldsymbol{\theta}$ lze najít jen ze znalosti podmíněného rozdělení $Y|\mathbf{X}$ jako řešení systému věrohodnostních rovnic

$$\sum_{i=1}^n \frac{\partial \log f_{Y|\mathbf{X}}(Y_i|\mathbf{X}_i; \boldsymbol{\theta})}{\partial \theta_j} = 0, \quad j \in \{1, \dots, q+k+2\}.$$

Dále pokud jsou navíc splněny předpoklady věty 5, tak k výpočtu Fisherovy matice stačí pracovat jen s hustotou podmíněného rozdělení. Potom Fisherova informační matice $\mathbb{J}^{(n)}(\boldsymbol{\theta})$ pro náhodnou matici (\mathbf{Y}, \mathbb{X}) se skládá z prvků

$$J_{ij}^{(n)}(\boldsymbol{\theta}) = - \sum_{l=1}^n \mathbb{E} \left(\frac{\partial^2 \log f_{Y|\mathbf{X}}(Y_l|\mathbf{X}_l; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right), \quad i, j \in \{1, \dots, q+k+2\}.$$

V této práci budeme dále předpokládat, že ono podmíněné rozdělení je normální, které zřejmě předpoklady regularity (R1)-(R6) a věty 5 splňuje. Jediné, co budeme při naší parametrizaci muset dodržet navíc, jsou předpoklady o parametrickém prostoru Θ .

2.2 Testy s rušivými parametry

V modelu $(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2 h(\mathbf{X}; \boldsymbol{\tau}))$ budeme vzhledem k naší parametrizaci (1.5) chtít testovat, zda je funkce h identicky rovna jedné, tj. $h \equiv 1$, neboť potom bude platit homoskedastický model $(Y, \mathbf{X}) \sim \text{LM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2)$. Hypotéza tvaru $h \equiv 1$ potom bude odpovídat specifické hodnotě parametru $\boldsymbol{\tau}$, popř. $\boldsymbol{\beta}$, pokud bychom uvažovali závislost na střední hodnotě.

Vyvstává nám zde ovšem problém. Chceme testovat hypotézu o parametru $\boldsymbol{\tau}$, ale v naší parametrizaci ještě vystupují další neznámé parametry $\boldsymbol{\beta}$ a σ^2 . Návod, jak testovat hypotézy jen o části neznámých parametrů, poskytuje právě teorie *testů s rušivými parametry*.

Uvažme opět obecně náhodný vektor \mathbf{V} s hustotou $f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})$ vůči σ -konečné míře $\boldsymbol{\mu}$ a $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)^\top$ je m -rozměrný parametr, kde tentokrát $m \geq 2$. Tento parametr si rozdělme na dvě části. Buď $1 \leq q < m$ a označme

$$\boldsymbol{\tau} = (\theta_1, \dots, \theta_q)^\top, \quad \boldsymbol{\psi} = (\theta_{q+1}, \dots, \theta_m)^\top.$$

Parametr $\boldsymbol{\tau}$ bude náš *cílový parametr*, o kterém budeme chtít provádět test hypotézy

$$H_0 : \boldsymbol{\tau} = \boldsymbol{\tau}_0 \quad \text{proti alternativě } H_1 : \boldsymbol{\tau} \neq \boldsymbol{\tau}_0.$$

Parametr $\boldsymbol{\psi}$ je potřebný jen k plnému popisu modelu, ale o něm žádnou hypotézu testovat nehodláme, nazveme jej tedy *rušivým parametrem*.

Naše testy budou založené na teorii maximální věrohodnosti, proto předpokládejme, že systém hustot \mathcal{F}_m je regulární (viz definice 8) a jsou splněny další předpoklady věty 5. Mějme dále náhodný výběr $\mathbb{V} = (\mathbf{V}_1, \dots, \mathbf{V}_n)^\top$ z rozdělení s hustotou $f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})$. Označme

$$\mathbf{U}_1(\boldsymbol{\theta}) = \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\tau}} = \begin{pmatrix} \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_q} \end{pmatrix}, \quad \mathbf{U}_2(\boldsymbol{\theta}) = \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\psi}} = \begin{pmatrix} \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_{q+1}} \\ \vdots \\ \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_m} \end{pmatrix}$$

a dále si Fisherovu informační matici pro celý náhodný výběr rozdělme na bloky

$$\mathbb{J}^{(n)}(\boldsymbol{\theta}) = \begin{pmatrix} \mathbb{J}_{11}^{(n)}(\boldsymbol{\theta}) & \mathbb{J}_{12}^{(n)}(\boldsymbol{\theta}) \\ \mathbb{J}_{21}^{(n)}(\boldsymbol{\theta}) & \mathbb{J}_{22}^{(n)}(\boldsymbol{\theta}) \end{pmatrix},$$

kde $\mathbb{J}_{11}^{(n)}$ je matice typu $q \times q$ a ostatní matice mají takový rozměr, aby $\mathbb{J}^{(n)}$ byla matice typu $m \times m$.

Dále bude nutné zavést značení, které nám umožní zapsat inverzi blokové diagonální matice. Vše nám shrnuje následující lemma.

Lemma 6. *Nechť*

$$\mathbb{J} = \begin{pmatrix} \mathbb{J}_{11} & \mathbb{J}_{12} \\ \mathbb{J}_{21} & \mathbb{J}_{22} \end{pmatrix}$$

je regulární matice, přičemž bloky \mathbb{J}_{11} a \mathbb{J}_{22} jsou čtvercové a regulární. Položme

$$\begin{aligned} \mathbb{J}_{11.2} &= \mathbb{J}_{11} - \mathbb{J}_{12} \mathbb{J}_{22}^{-1} \mathbb{J}_{21}, & \mathbb{J}^{11} &= \mathbb{J}_{11.2}^{-1}, & \mathbb{J}^{12} &= -\mathbb{J}_{11.2}^{-1} \mathbb{J}_{12} \mathbb{J}_{22}^{-1}, \\ \mathbb{J}_{22.1} &= \mathbb{J}_{22} - \mathbb{J}_{21} \mathbb{J}_{11}^{-1} \mathbb{J}_{12}, & \mathbb{J}^{22} &= \mathbb{J}_{22.1}^{-1}, & \mathbb{J}^{21} &= -\mathbb{J}_{22.1}^{-1} \mathbb{J}_{21} \mathbb{J}_{11}^{-1}. \end{aligned}$$

Pak

$$\mathbb{J}^{-1} = \begin{pmatrix} \mathbb{J}^{11} & \mathbb{J}^{12} \\ \mathbb{J}^{21} & \mathbb{J}^{22} \end{pmatrix}.$$

Důkaz. Snadnými algebraickými úpravami se ukáže, že součin matice \mathbb{J} a \mathbb{J}^{-1} dává jednotkovou matici. \square

Pro zavedení testových statistik potřebujeme maximálně věrohodné odhady parametru $\boldsymbol{\theta}$. Musíme ovšem rozlišit dva odhady. První z nich, který budeme značit $\hat{\boldsymbol{\theta}}_n$, bude maximálně věrohodný odhad parametru $\boldsymbol{\theta}$, který není svazován žádnými dalšími podmínkami. Tento odhad rozdělíme podle parametrů $\boldsymbol{\tau}$ a $\boldsymbol{\psi}$ na

$$\hat{\boldsymbol{\theta}}_n = \begin{pmatrix} \hat{\boldsymbol{\tau}}_n \\ \hat{\boldsymbol{\psi}}_n \end{pmatrix}.$$

Druhý z nich, který budeme značit $\tilde{\boldsymbol{\theta}}_n$, bude maximálně věrohodný odhad za platnosti nulové hypotézy $\boldsymbol{\tau} = \boldsymbol{\tau}_0$. Zde se tedy maximalizace týká jen parametru $\boldsymbol{\psi}$, jehož maximálně věrohodný odhad za platnosti hypotézy budeme značit $\tilde{\boldsymbol{\psi}}_n$. Potom $\tilde{\boldsymbol{\theta}}_n$ lze zapsat jako

$$\tilde{\boldsymbol{\theta}}_n = \begin{pmatrix} \boldsymbol{\tau}_0 \\ \tilde{\boldsymbol{\psi}}_n \end{pmatrix} = \arg \max_{\boldsymbol{\theta} \in \Theta: \boldsymbol{\tau} = \boldsymbol{\tau}_0} \ell(\boldsymbol{\theta}).$$

Nyní si zavedeme následující statistiky

$$LM = [\mathbf{U}_1(\tilde{\boldsymbol{\theta}}_n)]^\top [\mathbb{J}_{11.2}^{(n)}(\tilde{\boldsymbol{\theta}}_n)]^{-1} \mathbf{U}_1(\tilde{\boldsymbol{\theta}}_n), \quad (2.1)$$

$$W = (\hat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0)^\top \mathbb{J}_{11.2}^{(n)}(\hat{\boldsymbol{\theta}}_n) (\hat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0), \quad (2.2)$$

$$LR = 2 \left[\ell(\hat{\boldsymbol{\theta}}_n) - \ell(\tilde{\boldsymbol{\theta}}_n) \right]. \quad (2.3)$$

U těchto statistik jsme schopni nalézt jejich asymptotické rozdělení za platnosti $H_0 : \boldsymbol{\tau} = \boldsymbol{\tau}_0$, ovšem jen za platnosti dalších nutných předpokladů. Potom tedy využijeme této znalosti asymptotického rozdělení, abychom sestavili test na hladině $\alpha \in (0, 1)$ o hypotéze $H_0 : \boldsymbol{\tau} = \boldsymbol{\tau}_0$ proti alternativě, že tomu tak není.

V následující větě se neobejdeme bez těchto dvou obecných předpokladů:

- (O1) $\Theta \subset \mathbb{R}^m$ je parametrický prostor, který obsahuje takové neprázdné otevřené okolí O , že skutečná hodnota parametru $\boldsymbol{\theta}$ náleží do tohoto okolí O .
- (O2) Nechť $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Theta$. Pak $f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}_1) = f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}_2)$ $[\mu]$ s.v. platí právě tehdy, je-li $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$.

Abychom mohli odvodit asymptotické rozdělení testových statistik LM , W a LR za platnosti $H_0 : \boldsymbol{\tau} = \boldsymbol{\tau}_0$, budeme potřebovat předpoklady (A1)-(A3):

- (A1) Pro $[\mu]$ skoro všechna \mathbf{v} , pro všechna $\boldsymbol{\theta} \in O$ a pro všechna $i, j, l \in \{1, \dots, m\}$ existuje derivace

$$\frac{\partial^3 f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j \partial \theta_l}.$$

(A2) Pro všechna $\boldsymbol{\theta} \in O$ a pro každé $i, j \in \{1, \dots, m\}$ platí

$$\int_M f''_{ij}(\mathbf{v}; \boldsymbol{\theta}) d\mu(\mathbf{v}) = 0.$$

(A3) Pro všechna $i, j, l \in \{1, \dots, m\}$ existují funkce $M_{ijl}(\mathbf{v}) \geq 0$ takové, že

$$\mathbf{E}_{\boldsymbol{\theta}} M_{ijl}(\mathbf{V}) < \infty \quad \text{a} \quad \left| \frac{\partial^3 \log f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j \partial \theta_l} \right| \leq M_{ijl}(\mathbf{v})$$

pro všechna $\boldsymbol{\theta} \in O$ a $[\mu]$ skoro všechna $\mathbf{v} \in M$.

Věta 7. *Mějme náhodný výběr $\mathbb{V} = (\mathbf{V}_1, \dots, \mathbf{V}_n)^\top$ z rozdělení s hustotou $f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta})$ vzhledem k σ -konečné míře μ . Nechť systém $\mathcal{F}_m = \{f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ je regulární (viz definice 8).*

Nechť jsou navíc splněny předpoklady (O1), (O2), (A1), (A2), (A3) a navíc je Fisherova informační matice $\mathbb{J}^{(n)}(\boldsymbol{\theta})$ spojitá ve skutečné hodnotě parametru $\boldsymbol{\theta}$. Jestliže $n \rightarrow \infty$, pak za platnosti hypotézy $H_0 : \boldsymbol{\tau} = \boldsymbol{\tau}_0$ platí

$$LM \xrightarrow{D} \chi_q^2, \quad W \xrightarrow{D} \chi_q^2, \quad LR \xrightarrow{D} \chi_q^2.$$

Důkaz. Důkaz je proveden v Anděl (2007). □

Známe tedy asymptotické rozdělení našich statistik za platnosti nulové hypotézy $H_0 : \boldsymbol{\tau} = \boldsymbol{\tau}_0$. Můžeme tedy sestavit testy na asymptotické hladině významnosti $\alpha \in (0, 1)$. Tyto testy dosáhnou největší síly, pokud zvolíme kritický obor ve tvaru $(\chi_q^2(1 - \alpha), \infty)$, tedy když

$$H_0 : \boldsymbol{\tau} = \boldsymbol{\tau}_0 \text{ zamítáme} \iff LM, W \text{ nebo } LR \geq \chi_q^2(1 - \alpha).$$

Poznámka. Několik poznámek k těmto testům.

- a) Test založený na statistice LM (2.1) se nazývá *skórový test*, ale dříve se také používal název *test založený na Lagrangeových multiplikátorech*. Výhodou tohoto testu je, že stačí znát maximálně věrohodný odhad za platnosti nulové hypotézy. Výpočet inverzní matice nebývá náročný, neboť řád této matice odpovídá dimenzi q parametru $\boldsymbol{\tau}$ a není obvykle příliš velký.
- b) Test založený na statistice W (2.2) se nazývá *Waldův test*. Narozdíl od skórového testu nevyžaduje invertování informační matice, ale na druhou stranu vyžaduje maximálně věrohodný odhad bez předpokladu platnosti nulové hypotézy, což může být výpočetně náročné.
- c) Test založený na statistice LR (2.3) se nazývá *test založený na věrohodnostním poměru (likelihood ratio test)*. Tento test narozdíl od předchozích dvou nevyžaduje znalost Fisherovy informační matice.

Kapitola 3

Bartlettův test

V Andělově učebnici Anděl (2007, str. 210) je zaveden model analýzy rozptylu jednoduchého třídění. Základní myšlenkou je, že máme k dispozici I nezávislých náhodných výběrů z normálního rozdělení s obecně různými středními hodnotami a shodným rozptylem. Na základě analýzy rozptylu se potom odvozuje test o shodnosti středních hodnot. Ovšem při tomto testování je zásadní předpoklad shodnosti rozptylů ve všech skupinách, který ovšem nemusí být vždy splněn. Naším úkolem je tedy sestavit test, který by na hladině $\alpha \in (0, 1)$ zamítl (nebo nevyvrátil) tento předpoklad homoskedasticity.

Uvažme tedy obecně následující model:

$$Y_{ij} = \mu_i + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim N(0, \sigma_i^2), \quad i \in \{1, \dots, I\}, n_i \geq 2, j \in \{1, \dots, n_i\}, \quad (3.1)$$

kde $\mu_i \in \mathbb{R}$, $\sigma_i^2 > 0$, $i \in \{1, \dots, I\}$ jsou neznámé parametry, I nám označuje počet skupin a $n := n_1 + \dots + n_I$. Hypotéza homoskedasticity tedy odpovídá tomu, že $\sigma_1^2 = \dots = \sigma_I^2 (= \sigma^2)$. Později si ukážeme, že se tento model dá chápat jako heteroskedastický lineární model tak, jak jsme ho zavedli v podkapitole 1.3.

Nejprve se seznámíme s postupem, který navrhl Bartlett ve svém článku Bartlett (1937). Označme výběrové rozptyly v jednotlivých skupinách symboly

$$S_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2, \quad \text{kde } \bar{Y}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}.$$

Za odhad společné hodnoty rozptylů σ^2 se dá považovat statistika

$$S^2 = \frac{1}{n - I} \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2 = \sum_{i=1}^I \frac{n_i - 1}{n - I} S_i^2.$$

Jedná se o vážený průměr jednotlivých odhadů s vahami $(n_i - 1)/(n - I)$. Bartlettova testová statistika má potom tvar

$$\begin{aligned} B &= \frac{1}{c} \left((n - I) \log S^2 - \sum_{i=1}^I (n_i - 1) \log S_i^2 \right) \\ &= \frac{n - I}{c} \left(\log S^2 - \sum_{i=1}^I \frac{n_i - 1}{n - I} \log S_i^2 \right), \end{aligned} \quad (3.2)$$

kde konstanta c je dána vztahem

$$c = 1 + \frac{1}{3(I-1)} \left(\sum_{i=1}^I \frac{1}{n_i - 1} - \frac{1}{n - I} \right).$$

Bartlett (1937) odvodil, že statistika B má za platnosti nulové hypotézy asymptoticky χ_{I-1}^2 rozdělení.

Vidíme, že statistika B je založena na porovnávání logaritmu váženého průměru odhadů rozptylu pro jednotlivá i s váženým průměrem logaritmů těchto odhadů. Tato statistika je tedy založena na principech věrohodnostního poměru.

V následující části odvodíme na základě podkapitoly 2.2 test založený na testové statistice LR (viz (2.3)). Nicméně se nebude jednat přímo o Bartlettovu statistiku (3.2). Jejich podobnost a asymptotické vlastnosti budeme diskutovat na závěr této kapitoly (podkapitola 3.2).

3.1 Test založený na věrohodnostním poměru

Náš model pojmemme obecněji, než bylo uvedeno na začátku této kapitoly. Nebudeme totiž předem předpokládat, do které skupiny je jaké pozorování zařazeno. Tento proces bude řízen náhodným vektorem \mathbf{X} .

Nechť $(Y, \mathbf{X}) \sim \text{NLM}(\mathbf{X}^\top \boldsymbol{\mu}, \sigma^2 \mathbf{X}^\top \mathbf{w})$. Buď $I \geq 2$. I -složkový náhodný vektor

$$\mathbf{X} \sim \text{Mult}_I(1, \mathbf{p}), \text{ kde } \mathbf{p} = (p_1, \dots, p_I)^\top \in (0, 1)^I,$$

je neznámý parametr, který splňuje $\sum_{i=1}^I p_i = 1$. Dále $\boldsymbol{\mu} = (\mu_1, \dots, \mu_I)^\top \in \mathbb{R}^I$, $\sigma^2 > 0$ a $\mathbf{w} = (w_1, \dots, w_I)^\top \in (0, \infty)^I$ jsou neznámé parametry až na w_1 , o kterém předpokládejme, že platí $w_1 = 1$ za účelem identifikace, viz předpoklad (O2). Funkcí h ze zavedení heteroskedastického modelu je zjevně $h(\mathbf{X}; \mathbf{w}) = \mathbf{X}^\top \mathbf{w}$.

Uvědomme si, že náhodný vektor \mathbf{X} může nabývat jen kanonických vektorů

$$\mathbf{e}_i = (\delta_{ij})_{j=1}^I, \quad i \in \{1, \dots, I\},$$

a to s pravděpodobností

$$p(\mathbf{e}_i; \mathbf{p}) = \text{P}(\mathbf{X} = \mathbf{e}_i) = p_i.$$

Potom $\mathbf{X}^\top \boldsymbol{\mu} = \mu_i$ a $\mathbf{X}^\top \mathbf{w} = w_i$ pro nějaké $i \in \{1, \dots, I\}$. Tedy

$$(Y | \mathbf{X} = \mathbf{e}_i) \sim \text{N}(\mu_i, \sigma^2 w_i).$$

Tedy náhodný vektor \mathbf{X} nám určuje, do které z I skupin zařadit pozorování Y . Dále předpokládáme, že náhodná veličina Y se v i -té skupině řídí normálním rozdělením se střední hodnotou μ_i a rozptylem $\sigma^2 w_i$, což mimochodem znamená, že v první skupině je rozptyl právě σ^2 . Tedy připouštíme, že v každé skupině může být díky váhovému parametru \mathbf{w} různý rozptyl.

Rozptyl ve všech skupinách bude stejný, pokud $w_1 = 1 = w_2 = \dots = w_I$. Chceme tedy odvodit test založený na věrohodnostním poměru o hypotéze

$$H_0 : \mathbf{w} = \mathbf{1}_I \quad \text{proti alternativě } H_1 : \mathbf{w} \neq \mathbf{1}_I. \quad (3.3)$$

Postupujeme tedy dle podkapitoly 2.2, nejprve se ovšem zaměříme na parametrický prostor tohoto modelu.

Parametr \mathbf{w} bude cílovým parametrem (až na $w_1 = 1$), držíme-li se tedy značení zavedeného v sekci 2.2, tak $\boldsymbol{\tau} = (w_2, \dots, w_I)$ a $\boldsymbol{\tau}_0 = \mathbf{1}_{I-1}$. Rušivými parametry budou $\boldsymbol{\mu}, \sigma^2, \mathbf{p}$. Parametr \mathbf{p} ovšem nepochází přímo z otevřené množiny, neboť je svazován podmínkou $p_1 + \dots + p_I = 1$. Definujme si tedy

$$P = \{(p_1, \dots, p_{I-1})^\top \in (0, 1)^{I-1} : p_1 + \dots + p_{I-1} < 1\},$$

tato množina již otevřená je. Parametr p_I potom lze na základě p_1, \dots, p_{I-1} určit jako $p_I = 1 - (p_1 + \dots + p_{I-1})$. Označme $\boldsymbol{\psi} = (\mu_1, \dots, \mu_I, \sigma^2, p_1, \dots, p_{I-1})^\top$ a $\boldsymbol{\theta} = (\boldsymbol{\tau}^\top, \boldsymbol{\psi}^\top)^\top$.

Parametrický prostor tedy bude

$$\Theta = (0, \infty)^{I-1} \times \mathbb{R}^I \times (0, \infty) \times P,$$

splňuje tak náš požadavek na otevřenost (R1) a navíc skutečná hodnota parametru $\boldsymbol{\theta}$ musí ležet v nějakém otevřeném okolí (O1).

Dále ať $f(y, \mathbf{x}; \boldsymbol{\theta})$ značí hustotu náhodného vektoru (Y, \mathbf{X}) vůči součinné míře $\lambda \times \nu$, $f_{Y|\mathbf{X}}(y|\mathbf{x}; \mathbf{w}, \boldsymbol{\mu}, \sigma^2)$ značí hustotu podmíněného rozdělení $Y|\mathbf{X}$ vůči Lebesgueově míře λ a $p(\mathbf{x}; \mathbf{p})$ značí hustotu náhodného vektoru \mathbf{X} vzhledem ke sčítací míře ν .

Nechť $y \in \mathbb{R}$ a \mathbf{x} je nějaký kanonický vektor, potom pro každé $\boldsymbol{\theta} \in \Theta$ platí:

$$f_{Y|\mathbf{X}}(y|\mathbf{x}; \mathbf{w}, \boldsymbol{\mu}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2\mathbf{x}^\top\mathbf{w}}} \exp\left\{-\frac{(y - \mathbf{x}^\top\boldsymbol{\mu})^2}{2\sigma^2\mathbf{x}^\top\mathbf{w}}\right\}, \quad (3.4)$$

$$p(\mathbf{x}; \mathbf{p}) = \prod_{i=1}^I p_i^{x_i} = p_1^{x_1} \cdot \dots \cdot p_I^{x_I}, \quad (3.5)$$

$$f(y, \mathbf{x}; \boldsymbol{\theta}) = f_{Y|\mathbf{X}}(y|\mathbf{x}; \mathbf{w}, \boldsymbol{\mu}, \sigma^2) \cdot p(\mathbf{x}; \mathbf{p}) \quad [\lambda \times \nu] \text{ s.v.} \quad (3.6)$$

Dále systém hustot

$$\mathcal{F}_{I-1+I+1+I-1} = \mathcal{F}_{3I-1} = \{f(y, \mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$$

je regulární (definice 8), což plyne z vlastností normálního a multinomického rozdělení.

Uvažme náhodný výběr $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$ z rozdělení (Y, \mathbf{X}) . Pak

$$\mathbf{N} = (N_1, \dots, N_I)^\top = \sum_{l=1}^n \mathbf{X}_l \sim \text{Mult}_I(n, \mathbf{p})$$

a sdružená hustota $p_{\mathbb{X}}(\mathbb{X}; \mathbf{p})$ vzhledem k součinné míře $\nu^n = \nu \times \dots \times \nu$ náhodné matice \mathbb{X} , která představuje náhodný výběr $\mathbf{X}_1, \dots, \mathbf{X}_n$, je potom součinem hustot z rovnice (3.5). Platí tedy

$$p_{\mathbb{X}}(\mathbb{X}; \mathbf{p}) = \prod_{l=1}^n p_1^{x_{1l}} \cdot \dots \cdot p_I^{x_{Il}} = \prod_{i=1}^I p_i^{n_i}, \quad (3.7)$$

kde $(n_1, \dots, n_I)^\top$ jsou realizace náhodného vektoru \mathbf{N} , pro které $n_1 + \dots + n_I = n$. Předpokládejme, že $n_i \geq 1, i \in \{1, \dots, I\}$. Tedy n_i udává kolik pozorování \mathbf{X}_l

z náhodného výběru nabylo hodnoty e_i , neboli kolika veličinám Y_l určilo \mathbf{X}_l jejich příslušnost do skupiny $i \in \{1, \dots, I\}$.

Pro každé $l \in \{1, \dots, n\}$ určitě Y_l spadá do některé skupiny $i \in \{1, \dots, I\}$ reprezentované hodnotou náhodného vektoru \mathbf{X}_l . Přeznačme tedy veličiny Y_l na veličiny Y_{ij} , kde $i \in \{1, \dots, I\}$ označuje příslušnou skupinu a $j \in \{1, \dots, n_i\}$ bude index, pomocí kterého budu procházet veličiny ve skupinách. Používáme tedy značení shodné s tím, které jsme uvedli na začátku této kapitoly, viz (3.1).

Zavedme si pro $i \in \{1, \dots, I\}$

a) průměr hodnot Y_{ij} v i -té skupině $\bar{Y}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$ a

b) reziduální součet čtverců v i -té skupině $RSS_i = \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2$.

Tvrzení 8. V modelu uvedeném výše za platnosti $H_0 : \mathbf{w} = \mathbf{1}_I$ statistika

$$LR = n \log \left(\frac{1}{n} \sum_{i=1}^I RSS_i \right) - \sum_{i=1}^I n_i \log \left(\frac{RSS_i}{n_i} \right) \xrightarrow{D} \chi_{I-1}^2, \quad n \rightarrow \infty.$$

Důkaz. Ukážeme si, že statistika LR je statistikou pro test poměrem věrohodnosti, viz (2.3). K tomu potřebujeme nalézt maximálně věrohodné odhady parametrů $\boldsymbol{\mu}$, σ^2 , \mathbf{w} a \mathbf{p} , a to jak za platnosti nulové hypotézy H_0 , tak bez platnosti tohoto předpokladu. Nakonec tyto odhady dosadíme do logaritmické věrohodnosti ℓ a odečteme vzniklé výrazy. Začneme tedy s určením funkce ℓ .

Sdruženou hustotu podmíněného rozdělení $f_{\mathbf{Y}|\mathbb{X}}(\mathbf{y}|\mathbb{X}; \mathbf{w}, \boldsymbol{\mu}, \sigma^2)$ lze odvodit pomocí vzorce (3.4) a dostaneme

$$f_{\mathbf{Y}|\mathbb{X}}(\mathbf{y}|\mathbb{X}; \mathbf{w}, \boldsymbol{\mu}, \sigma^2) = (2\pi\sigma^2)^{-\frac{n}{2}} \cdot \prod_{i=1}^I w_i^{\frac{n_i}{2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^I \frac{1}{w_i} \sum_{j=1}^{n_i} (y_{ij} - \mu_i)^2 \right\}. \quad (3.8)$$

Označme nyní

$$\ell_1(\mathbf{w}, \boldsymbol{\mu}, \sigma^2) = \log f_{\mathbf{Y}|\mathbb{X}}(\mathbf{Y}|\mathbb{X}; \mathbf{w}, \boldsymbol{\mu}, \sigma^2) \text{ a } \ell_2(\mathbf{p}) = \log p_{\mathbb{X}}(\mathbb{X}; \mathbf{p}).$$

Pak logaritmováním výrazů (3.8) a (3.7) dostaneme

$$\begin{aligned} \ell_1(\mathbf{w}, \boldsymbol{\mu}, \sigma^2) &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \\ &\quad - \sum_{i=1}^I \frac{n_i}{2} \log w_i - \frac{1}{2\sigma^2} \sum_{i=1}^I \frac{1}{w_i} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2, \end{aligned} \quad (3.9)$$

$$\ell_2(\mathbf{p}) = \sum_{i=1}^I n_i \log p_i. \quad (3.10)$$

Směřujeme nyní k logaritmické věrohodnostní funkci $\ell(\boldsymbol{\theta}) = \log f(\mathbf{Y}, \mathbb{X}; \boldsymbol{\theta})$. Rozšířením rovnosti (3.6) na sdružené hustoty dostáváme spolu s rovnostmi (3.9)

a (3.10) následující

$$\begin{aligned}\ell(\boldsymbol{\theta}) &= \log f(\mathbf{Y}, \mathbb{X}; \boldsymbol{\theta}) = \ell_1(\mathbf{w}, \boldsymbol{\mu}, \sigma^2) + \ell_2(\mathbf{p}), \quad [\lambda^n \times \mathbf{v}^n] \text{ s.v.}, \\ \ell(\boldsymbol{\theta}) &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \sum_{i=1}^I \frac{n_i}{2} \log w_i - \\ &\quad - \frac{1}{2\sigma^2} \sum_{i=1}^I \frac{1}{w_i} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2 + \sum_{i=1}^I n_i \log p_i.\end{aligned}\tag{3.11}$$

Druhá rovnost také platí jen $[\lambda^n \times \mathbf{v}^n]$ skoro všude, ovšem od této chvíle si dovolíme tento dovětek nadále nepsat.

Nyní můžeme přistoupit k hledání maximálně věrohodných odhadů. Začneme s hledáním odhadů $\hat{\mathbf{p}}_n$ a $\tilde{\mathbf{p}}_n$. Všimneme si, že parametr \mathbf{p} není nijak svázán s ostatními parametry, tedy platnost nulové hypotézy se na našich odhadech nijak neprojeví. Ukáží konkrétně na výpočtech.

Jedná se nám v podstatě o maximalizaci funkce ℓ_2 za dodatečné podmínky $p_1 + \dots + p_I = 1$, což vyřešíme metodou Lagrangeových multiplikátorů. Definujme si pomocnou funkci $g(\mathbf{p}, \lambda)$ jako

$$g(\mathbf{p}, \lambda) = \sum_{i=1}^I n_i \log p_i - \lambda \left(\sum_{i=1}^I p_i - 1 \right).$$

Potom spočtíme derivace podle proměnných λ a $p_i, i \in \{1, \dots, I\}$.

$$\begin{aligned}\frac{\partial g(\mathbf{p}, \lambda)}{\partial \lambda} &= 1 - \sum_{i=1}^I p_i, \\ \frac{\partial g(\mathbf{p}, \lambda)}{\partial p_i} &= \frac{n_i}{p_i} - \lambda, \quad i \in \{1, \dots, I\}.\end{aligned}$$

Položíme-li tyto derivace rovny nule, tak pro každé $i \in \{1, \dots, I\}$ dostáváme, že musí platit $\lambda = \frac{n_i}{p_i}$, tedy také $\hat{p}_i = \frac{n_i}{\lambda}$. Potom musí také dle první rovnosti platit

$$1 = \sum_{i=1}^I \frac{n_i}{\lambda} = \frac{n_1 + \dots + n_I}{\lambda} = \frac{n}{\lambda}.$$

Tím tedy dostáváme $\lambda = n$ a také odhady

$$\hat{\mathbf{p}}_n = \tilde{\mathbf{p}}_n = \left(\frac{n_1}{n}, \dots, \frac{n_I}{n} \right)^\top.\tag{3.12}$$

U ostatních parametrů budeme opravdu muset rozlišit případy, kdy nulová hypotéza (3.3) platí a kdy ne. Začneme s jednodušším případem, tedy s tím, když tato hypotéza platí.

Pak tedy víme, že $w_1 = 1 = w_2 = \dots = w_I$ a proto má funkce ℓ z (3.11) tvar

$$\ell(\boldsymbol{\tau}_0, \boldsymbol{\psi}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2 + \sum_{i=1}^I n_i \log p_i.$$

Spočtěme tedy derivace podle proměnných $\boldsymbol{\mu}$ a σ^2 (podle \boldsymbol{p} už nemusíme):

$$\frac{\partial \ell(\boldsymbol{\tau}_0, \boldsymbol{\psi})}{\partial \mu_i} = \frac{1}{\sigma^2} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i), \quad i \in \{1, \dots, I\}, \quad (3.13)$$

$$\frac{\partial \ell(\boldsymbol{\tau}_0, \boldsymbol{\psi})}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2. \quad (3.14)$$

Pokud položíme výraz (3.13) roven nule, tak pro každé $i \in \{1, \dots, I\}$ dostaneme odhady

$$\tilde{\boldsymbol{\mu}}_n = \left(\frac{1}{n_1} \sum_{j=1}^{n_1} Y_{1j}, \dots, \frac{1}{n_I} \sum_{j=1}^{n_I} Y_{Ij} \right) = (\bar{Y}_{1\bullet}, \dots, \bar{Y}_{I\bullet}). \quad (3.15)$$

Dále pokud položíme výraz (3.14) roven nule a dosadíme odhad $\tilde{\boldsymbol{\mu}}_n$, obdržíme

$$\tilde{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2 = \frac{1}{n} \sum_{i=1}^I RSS_i. \quad (3.16)$$

Získali jsme tedy maximálně věrohodné odhady za platnosti nulové hypotézy. Dále pokračujeme nalezením takových odhadů bez tohoto předpokladu. Postupujeme zcela obdobně. Funkce ℓ má tedy tvar

$$\begin{aligned} \ell(\boldsymbol{\theta}) = & -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \sum_{i=1}^I \frac{n_i}{2} \log w_i - \\ & - \frac{1}{2\sigma^2} \sum_{i=1}^I \frac{1}{w_i} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2 + \sum_{i=1}^I n_i \log p_i. \end{aligned}$$

Derivováním této funkce dostáváme

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \mu_i} = \frac{1}{\sigma^2 w_i} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i), \quad i \in \{1, \dots, I\}, \quad (3.17)$$

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^I \frac{1}{w_i} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2, \quad (3.18)$$

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial w_i} = -\frac{n_i}{2w_i} + \frac{1}{2\sigma^2 w_i^2} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2, \quad i \in \{1, \dots, I\}. \quad (3.19)$$

Položíme-li výraz (3.17) roven nule, tak dostaneme ten samý odhad pro $\boldsymbol{\mu}$ jako v předešlém případě, tedy odhad $\hat{\boldsymbol{\mu}}_n = \tilde{\boldsymbol{\mu}}_n$, viz (3.15). Jelikož předpokládáme znalost $w_1 = 1$, tak můžeme položit $\hat{w}_1 = 1$. Rovnice (3.19) pro $i = 1$ je tedy nadbytečná. Ovšem vyjdeme-li z ní, tak dostaneme

$$\hat{\sigma}_n^2 = \frac{1}{n_1} \sum_{j=1}^{n_1} (Y_{1j} - \bar{Y}_{1\bullet})^2 = \frac{RSS_1}{n_1}.$$

Potom ze vztahu (3.19) pro každé $i \in \{2, \dots, I\}$ dostaneme, že musí platit

$$n_i \hat{w}_i = \frac{1}{\hat{\sigma}_n^2} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2,$$

tedy jednoduchou úpravou dostáváme, že platí

$$\widehat{\sigma}_n^2 \widehat{w}_i = \frac{RSS_i}{n_i}, \quad \forall i \in \{1, \dots, I\}. \quad (3.20)$$

Odtud bychom dostali, že

$$\widehat{w}_i = \frac{RSS_i n_1}{RSS_1 n_i}, \quad \forall i \in \{1, \dots, I\}.$$

Zbývá ověřit, že takové odhady také nulují výraz (3.18), aby byla splněna celá soustava věrohodnostních rovnic. Jelikož

$$\begin{aligned} \sum_{i=1}^I \frac{1}{\widehat{w}_i} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2 &= \sum_{i=1}^I \frac{RSS_i}{\widehat{w}_i} = \sum_{i=1}^I \frac{n_i}{n_1} RSS_1 \\ &= RSS_1 \frac{n_1 + \dots + n_I}{n_1} = \frac{n RSS_1}{n_1} = n \widehat{\sigma}_n^2, \end{aligned}$$

tak dostáváme, že opravdu

$$-\frac{n}{2\widehat{\sigma}_n^2} + \frac{1}{2\widehat{\sigma}_n^4} \sum_{i=1}^I \frac{1}{\widehat{w}_i} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2 = -\frac{n}{2\widehat{\sigma}_n^2} + \frac{1}{2\widehat{\sigma}_n^4} n \widehat{\sigma}_n^2 = 0.$$

Nyní již jsme schopni určit testovou statistiku

$$LR = 2 \left(\ell(\widehat{\boldsymbol{\theta}}_n) - \ell(\widetilde{\boldsymbol{\theta}}_n) \right).$$

Díky výše uvedeným vzorcům (3.11), (3.15), (3.16) a (3.20) platí:

$$\begin{aligned} \ell(\widehat{\boldsymbol{\theta}}_n) &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \widehat{\sigma}_n^2 - \sum_{i=1}^I \frac{n_i}{2} \log \widehat{w}_i - \frac{1}{2\widehat{\sigma}_n^2} \sum_{i=1}^I \frac{RSS_i}{\widehat{w}_i} + \ell_2(\widehat{\boldsymbol{p}}_n) \\ &= -\frac{n}{2} \log(2\pi) - \sum_{i=1}^I \frac{n_i}{2} \log(\widehat{\sigma}_n^2 \widehat{w}_i) - \frac{1}{2} \sum_{i=1}^I \frac{RSS_i}{\widehat{\sigma}_n^2 \widehat{w}_i} + \ell_2(\widehat{\boldsymbol{p}}_n) \\ &= -\frac{n}{2} \log(2\pi) - \sum_{i=1}^I \frac{n_i}{2} \log\left(\frac{RSS_i}{n_i}\right) - \frac{1}{2} \sum_{i=1}^I \frac{n_i RSS_i}{RSS_i} + \ell_2(\widehat{\boldsymbol{p}}_n) \\ &= -\frac{n}{2} \log(2\pi) - \sum_{i=1}^I \frac{n_i}{2} \log\left(\frac{RSS_i}{n_i}\right) - \frac{n}{2} + \ell_2(\widehat{\boldsymbol{p}}_n), \\ \ell(\widetilde{\boldsymbol{\theta}}_n) &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \widetilde{\sigma}_n^2 - \frac{1}{2\widetilde{\sigma}_n^2} \sum_{i=1}^I RSS_i + \ell_2(\widetilde{\boldsymbol{p}}_n) \\ &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log\left(\frac{1}{n} \sum_{i=1}^I RSS_i\right) - \frac{n}{2} + \ell_2(\widetilde{\boldsymbol{p}}_n). \end{aligned}$$

Nyní si už jen uvědomíme, že jsme odvodili, že $\widetilde{\boldsymbol{p}}_n = \widehat{\boldsymbol{p}}_n$, viz (3.12), a odečteme dvojnásobek výrazů výše, čímž dostaneme statistiku

$$\begin{aligned} LR &= n \log\left(\frac{1}{n} \sum_{i=1}^I RSS_i\right) - \sum_{i=1}^I n_i \log\left(\frac{RSS_i}{n_i}\right) \\ &= n \left[\log\left(\frac{1}{n} \sum_{i=1}^I RSS_i\right) - \sum_{i=1}^I \frac{n_i}{n} \log\left(\frac{RSS_i}{n_i}\right) \right]. \end{aligned} \quad (3.21)$$

Vidíme tedy, že naše testová statistika je založená na rozdílu logaritmu celkového reziduálního součtu čtverců děleného n a součtu vážených logaritmů reziduálních součtů čtverců dělených n_i v rámci i -té skupiny.

Zbývá ověřit předpoklady věty 7, které, jak víme, platí pro normální i multinomické rozdělení. Proto již není těžké odvodit, že jsou splněny i v tomto modelu.

Podle věty 7 má tedy naše statistika asymptotické rozdělení chí kvadrát o takovém počtu stupňů volnosti, kolik bylo cílových parametrů. Jelikož jsme předpokládali znalost $w_1 = 1$, tak mluvíme jen o zbylých w_2, \dots, w_I . Tedy počet stupňů volnosti je $I - 1$. Celkově tedy

$$LR \xrightarrow{D} \chi_{I-1}^2, \quad n \rightarrow \infty.$$

□

Vidíme, že výsledná statistika LR (viz (3.21)) je velmi podobná statistice B (viz (3.2)), dokonce mají i stejné asymptotické rozdělení za platnosti nulové hypotézy. Náš test tedy pro obě statistiky zamítne na asymptotické hladině významnosti $\alpha \in (0, 1)$ hypotézu $H_0 : w_1 = w_2 = \dots = w_I = 1$ právě tehdy, když realizovaná hodnota statistiky B nebo LR převyší $\chi_{I-1}^2(1 - \alpha)$. Jak moc se ale použití statistik B a LR liší, shrneme v následující podkapitole.

3.2 Porovnání testových statistik

Nejprve se podívejme, v čem si jsou statistiky (3.2) a (3.21) podobné. Uvědomme si, že statistika B se také sestává z jednotlivých $RSS_i, i \in \{1, \dots, I\}$. Jednotlivé odhady rozptylu v i -té skupině, kde $i \in \{1, \dots, I\}$, vypadají takto:

$$S_i^2 = \frac{RSS_i}{n_i - 1}, \quad \hat{\sigma}_n^2 \hat{w}_i = \frac{RSS_i}{n_i}.$$

Liší se tedy jen ve jmenovateli, kde v prvním případě se RSS_i dělí počtem stupňů volnosti $n_i - 1$, kdežto v té druhé jen počtem pozorování v i -té skupině.

Celkové odhady jsou potom také podobného tvaru

$$S^2 = \frac{1}{n - I} \sum_{i=1}^I RSS_i, \quad \tilde{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^I RSS_i.$$

Ještě jednou uvádím statistiky, aby byl vidět jejich rozdíl:

$$\begin{aligned} LR &= n \left[\log \left(\frac{1}{n} \sum_{i=1}^I RSS_i \right) - \sum_{i=1}^I \frac{n_i}{n} \log \left(\frac{RSS_i}{n_i} \right) \right], \\ B &= \frac{n - I}{c} \left[\log \left(\frac{1}{n - I} \sum_{i=1}^I RSS_i \right) - \sum_{i=1}^I \frac{n_i - 1}{n - I} \log \left(\frac{RSS_i}{n_i - 1} \right) \right], \\ \text{kde } c &= 1 + \frac{1}{3(I - 1)} \left(\sum_{i=1}^I \frac{1}{n_i - 1} - \frac{1}{n - I} \right). \end{aligned} \quad (3.22)$$

Bartlettova statistika tedy používá počty stupňů volnosti namísto počtu pozorování a navíc je podělena konstantou c , která zpravidla bývá jen nepatrně větší než 1.

Všechny tyto úpravy způsobují, že statistika B má lepší asymptotické vlastnosti než LR , proto se také na rozdíl od LR v praxi používá k testování homoskedasticity. Toto tvrzení prozkoumáme dále v numerických simulacích v kapitole 5.

Jak již bylo řečeno, statistika B má za platnosti nulové hypotézy asymptoticky rozdělení χ_{I-1}^2 . Udává se, že tato vlastnost lze použít k testování, pokud platí $n_i \geq 7, i \in \{1, \dots, I\}$.

Tento Bartlettův test se zdá být při splnění předpokladu normality tím nejsilnějším z dostupných testů. Ovšem je velmi citlivý na porušení předpokladu o normálním rozdělení. Existují proto i další možné modifikace uvedené například v Zvára (2008, str. 120).

Kapitola 4

Breusch-Paganův test

Uvažme nyní následující normální heteroskedastický lineární model

$$(Y, \mathbf{X}) \sim \text{NLM}(\mathbf{X}^\top \boldsymbol{\beta}, \sigma^2 \exp(\mathbf{Z}^\top \boldsymbol{\tau})),$$

kde $\boldsymbol{\beta} = (\beta_0, \dots, \beta_k)^\top \in \mathbb{R}^{k+1}$, $\sigma^2 > 0$ a $\boldsymbol{\tau} = (\tau_1, \dots, \tau_q)^\top \in \mathbb{R}^q$ jsou neznámé parametry a $k \in \mathbb{N}_0, q \in \mathbb{N}$. Dále vektor \mathbf{Z} je nějakou známou transformací vektoru \mathbf{X} , tedy $\mathbf{Z} = g(\mathbf{X})$ pro známou funkci $g: \mathbb{R}^{k+1} \rightarrow \mathbb{R}^q$. Lze také místo $\exp(\mathbf{Z}^\top \boldsymbol{\tau})$ uvažovat obecně nějakou diferencovatelnou kladnou reálnou funkci h od skalárního součinu $\mathbf{Z}^\top \boldsymbol{\tau}$ s vlastností $h(0) = 1$, zde volíme $h = \exp$ pro jednoduchost.

Mějme k dispozici náhodný výběr $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$ z rozdělení (Y, \mathbf{X}) . Na základě tohoto náhodného výběru máme sestavit test o homoskedasticitě. Jak je vidět z naší parametrizace, tak toto nastává, pokud $\boldsymbol{\tau} = \mathbf{0}^q$.

Jak už jsme ospravedlnili v sekci 2.1, lze pracovat jen s podmíněným rozdělením $Y|\mathbf{X}$, které je dle předpokladu normální. Předpokládejme tedy, že

$$\mathbb{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} = (\mathbf{x}_{\bullet 0}, \dots, \mathbf{x}_{\bullet k}), \quad \mathbb{Z} = \begin{pmatrix} \mathbf{z}_1^\top \\ \vdots \\ \mathbf{z}_n^\top \end{pmatrix} = \begin{pmatrix} g(\mathbf{x}_1)^\top \\ \vdots \\ g(\mathbf{x}_n)^\top \end{pmatrix} = (\mathbf{z}_{\bullet 1}, \dots, \mathbf{z}_{\bullet q})$$

jsou matice reálných konstant s plnou sloupcovou hodnotí. Dále označme

$$\mathbb{W}(\boldsymbol{\tau}) = \text{diag}(w_1(\boldsymbol{\tau}), \dots, w_n(\boldsymbol{\tau})) = \text{diag}(\exp(\mathbf{z}_1^\top \boldsymbol{\tau}), \dots, \exp(\mathbf{z}_n^\top \boldsymbol{\tau})).$$

Vzhledem k našemu předpokladu má náhodný vektor $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ při pevných \mathbb{X} a \mathbb{Z} rozdělení

$$\mathbf{Y} \sim \text{N}_n(\mathbb{X}\boldsymbol{\beta}, \sigma^2 \mathbb{W}(\boldsymbol{\tau})). \quad (4.1)$$

Chceme tedy odvodit test pro následující hypotézu

$$H_0 : \boldsymbol{\tau} = \mathbf{0}^q \quad \text{proti alternativě} \quad H_1 : \boldsymbol{\tau} \neq \mathbf{0}^q. \quad (4.2)$$

Jelikož máme neznámé parametry $\boldsymbol{\tau}$, σ^2 a $\boldsymbol{\beta}$ a chceme se věnovat testování $\boldsymbol{\tau}$, tak použijeme teorie testů s rušivými parametry uvedené v podkapitole 2.2. Označme $\boldsymbol{\theta} = (\boldsymbol{\tau}, \sigma^2, \boldsymbol{\beta})$, příslušný parametrický prostor je tedy

$$\Theta = \mathbb{R}^q \times (0, \infty) \times \mathbb{R}^{k+1}.$$

Rušivé parametry tedy budou $\boldsymbol{\psi} = (\sigma^2, \beta_0, \dots, \beta_k)^\top$ a cílovým parametrem bude $\boldsymbol{\tau}$. Naším cílem bude odvodit skórový test založený na testové statistice LM (2.1).

Buďte

$$\mathbb{A} = \mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top, \quad \tilde{\sigma}_n^2 = \frac{RSS}{n}, \quad \mathbf{v} = \begin{pmatrix} u_1^2 - \tilde{\sigma}_n^2 \\ \vdots \\ u_n^2 - \tilde{\sigma}_n^2 \end{pmatrix},$$

kde \mathbf{u} je vektor reziduí, RSS je reziduální součet čtverců a $\mathbf{1}_n$ značí n -složkový sloupcový vektor $\mathbf{1}_n = (1, 1, \dots, 1)^\top$.

Tvrzení 9. *V modelu popsaném výše za platnosti $H_0 : \boldsymbol{\tau} = \mathbf{0}^q$ statistika*

$$LM = \frac{1}{2(\tilde{\sigma}_n^2)^2} \mathbf{v}^\top \mathbb{Z} (\mathbb{Z}^\top \mathbb{A} \mathbb{Z})^{-1} \mathbb{Z}^\top \mathbf{v} \xrightarrow{D} \chi_q^2, \quad n \rightarrow \infty.$$

Důkaz. Potřebujeme tedy nalézt Fisherovu informační matici a její inverzi, dále maximálně věrohodné odhady $\tilde{\boldsymbol{\beta}}_n$ a $\tilde{\sigma}_n^2$ za platnosti nulové hypotézy (4.2) a ověřit předpoklady věty 7.

Začneme se sdruženou hustotou $f_n(\mathbf{y}|\mathbb{X}; \boldsymbol{\theta})$ vektoru \mathbf{Y} při dané realizaci \mathbb{X} vzhledem k Lebesgueově míře λ^n . Platí

$$\begin{aligned} f_n(\mathbf{y}|\mathbb{X}; \boldsymbol{\theta}) &= \frac{(2\pi)^{-\frac{n}{2}}}{\sqrt{\det(\sigma^2 \mathbb{W}(\boldsymbol{\tau}))}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbb{X}\boldsymbol{\beta})^\top \frac{\mathbb{W}^{-1}(\boldsymbol{\tau})}{\sigma^2} (\mathbf{y} - \mathbb{X}\boldsymbol{\beta}) \right\} \\ &= \frac{(2\pi)^{-\frac{n}{2}}}{(\sigma^2)^{\frac{n}{2}} \left(\prod_{i=1}^n w_i(\boldsymbol{\tau}) \right)^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\sigma^2 w_i(\boldsymbol{\tau})} \right\} \\ &= \frac{(2\pi)^{-\frac{n}{2}}}{(\sigma^2)^{\frac{n}{2}} \exp \left\{ \frac{1}{2} \sum_{i=1}^n \mathbf{z}_i^\top \boldsymbol{\tau} \right\}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\sigma^2 \exp(\mathbf{z}_i^\top \boldsymbol{\tau})} \right\}. \end{aligned}$$

První rovnost je zřejmá z vlastnosti našeho modelu (4.1) a definice mnohorozměrného normálního rozdělení. Druhá rovnost plyne z diagonality matice $\sigma^2 \mathbb{W}(\boldsymbol{\tau})$, neboť pak její determinant je jen součin prvků na diagonále a inverzní matice má na diagonále jen převrácené prvky. Třetí rovnost nám plyne z vlastnosti prvků $w_i(\boldsymbol{\tau})$ a exponenciály.

Uvědomíme si, že systém hustot $\mathcal{F}_{q+k+2} = \{f(y|\mathbf{X}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ je regulární, viz definice 8. Všechny parametry pocházejí z nějakého otevřeného intervalu v parametrickém prostoru Θ . Navíc $M = \{y \in \mathbb{R} : f(y|\mathbf{X}; \boldsymbol{\theta}) > 0\}$ nezávisí na parametru $\boldsymbol{\theta}$. Díky vlastnostem normálního rozdělení jsou také splněny všechny ostatní předpoklady a také předpoklad z věty 5, tedy lze Fisherovu informační matici počítat pomocí střední hodnoty z druhých parciálních derivací logaritmské věrohodnosti ℓ .

Dále tedy budeme pracovat s přirozeným logaritmem hustoty. Platí

$$\ell(\boldsymbol{\theta}) = \log f_n(\mathbf{Y}|\mathbb{X}; \boldsymbol{\theta}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2} \sum_{i=1}^n \mathbf{z}_i^\top \boldsymbol{\tau} - \frac{1}{2\sigma^2} \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})}.$$

Dalším krokem je spočtení všech derivací podle neznámých parametrů.

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \tau_m} = -\frac{1}{2} \sum_{i=1}^n z_{im} + \frac{1}{2\sigma^2} \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} z_{im}, \quad m \in \{1, \dots, q\}, \quad (4.3)$$

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})}, \quad (4.4)$$

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \beta_j} = \frac{1}{\sigma^2} \sum_{i=1}^n \frac{x_{ij} (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})}, \quad j \in \{0, \dots, k\}. \quad (4.5)$$

Dále si můžeme vyřešit soustavu věrohodnostních rovnic a nalézt si tak maximálně věrohodné odhady. Ovšem pro naše potřeby stačí jen nalézt odhady σ^2 a $\boldsymbol{\beta}$ za platnosti nulové hypotézy (4.2). Takže položíme výrazy z rovnic (4.4) a (4.5) rovny nule, dále položíme $\boldsymbol{\tau} = \mathbf{0}^q$ a nalezneme odhady $\tilde{\sigma}_n^2$ a $\tilde{\boldsymbol{\beta}}_n$.

V soustavě věrohodnostních rovnic tedy platí

$$\begin{aligned} \frac{1}{2(\tilde{\sigma}_n^2)^2} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \tilde{\boldsymbol{\beta}}_n)^2 &= \frac{n}{2\tilde{\sigma}_n^2} \iff \tilde{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \tilde{\boldsymbol{\beta}}_n)^2. \\ \forall j \in \{0, \dots, k\} : \sum_{i=1}^n x_{ij} (Y_i - \mathbf{x}_i^\top \tilde{\boldsymbol{\beta}}_n) &= 0 \iff \forall j : \mathbf{x}_{\cdot j}^\top (\mathbf{Y} - \mathbb{X} \tilde{\boldsymbol{\beta}}_n) = 0 \\ &\iff \mathbb{X}^\top (\mathbf{Y} - \mathbb{X} \tilde{\boldsymbol{\beta}}_n) = \mathbf{0}^{k+1} \\ &\iff \mathbb{X}^\top \mathbb{X} \tilde{\boldsymbol{\beta}}_n = \mathbb{X}^\top \mathbf{Y} \\ &\iff \tilde{\boldsymbol{\beta}}_n = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y}. \end{aligned}$$

Zavedeme-li si podobně jako v sekci 1.2 vektor $\tilde{\mathbf{Y}} = \mathbb{X} \tilde{\boldsymbol{\beta}}_n$, dostaneme odhady tvaru

$$\tilde{\boldsymbol{\beta}}_n = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y} \quad \text{a} \quad \tilde{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \tilde{Y}_i)^2 = \frac{RSS}{n} \quad (4.6)$$

jako v klasickém modelu $\mathbf{Y} \sim \mathbf{N}_n(\mathbb{X}\boldsymbol{\beta}, \sigma^2 \mathbb{I}_n)$.

Pokračujme dále ve výpočtech druhých derivací funkce ℓ . Dostáváme

$$\begin{aligned} \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \tau_{m_1} \partial \tau_{m_2}} &= -\frac{1}{2\sigma^2} \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} z_{im_1} z_{im_2}, \quad m_1, m_2 \in \{1, \dots, q\}, \\ \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial (\sigma^2)^2} &= \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})}, \\ \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \beta_{j_1} \partial \beta_{j_2}} &= -\frac{1}{\sigma^2} \sum_{i=1}^n \frac{x_{ij_1} x_{ij_2}}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})}, \quad j_1, j_2 \in \{0, \dots, k\}, \\ \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \sigma^2 \partial \beta_j} &= -\frac{1}{\sigma^4} \sum_{i=1}^n \frac{x_{ij} (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})}, \quad j \in \{0, \dots, k\}, \\ \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \sigma^2 \partial \tau_m} &= -\frac{1}{2\sigma^4} \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} z_{im}, \quad m \in \{1, \dots, q\}, \\ \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \beta_j \partial \tau_m} &= -\frac{1}{2\sigma^2} \sum_{i=1}^n \frac{x_{ij} (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} z_{im}, \quad j \in \{0, \dots, k\}, m \in \{1, \dots, q\}. \end{aligned}$$

Využijme tedy těchto druhých derivací k tomu, abychom podle věty 5 spočetli Fisherovu informační matici za platnosti naší hypotézy (4.2). Přičemž tento předpoklad je v naší parametrizaci nutný jen u jedné rovnosti, proto jsem nad onu rovnost uvedl symbol H_0 . Připomenu-li, že $\mathbb{E} Y_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ a $\text{var} Y_i = \sigma^2 \exp(\mathbf{z}_i^\top \boldsymbol{\tau})$, tak snadno dostaneme, že

$$\begin{aligned} \mathbb{E} \left[-\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \tau_{m_1} \partial \tau_{m_2}} \right] &= \frac{1}{2\sigma^2} \sum_{i=1}^n \frac{\text{var} Y_i}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} z_{im_1} z_{im_2} = \frac{1}{2} \sum_{i=1}^n z_{im_1} z_{im_2} = \frac{\mathbf{z}_{\bullet m_1}^\top \mathbf{z}_{\bullet m_2}}{2}, \\ \mathbb{E} \left[-\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial (\sigma^2)^2} \right] &= -\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} \sum_{i=1}^n \frac{\text{var} Y_i}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} = -\frac{n}{2\sigma^4} + \frac{n\sigma^2}{\sigma^6} = \frac{n}{2\sigma^4}, \\ \mathbb{E} \left[-\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \beta_{j_1} \partial \beta_{j_2}} \right] &= \frac{1}{\sigma^2} \sum_{i=1}^n \frac{x_{ij_1} x_{ij_2}}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} \stackrel{H_0}{=} \frac{\mathbf{x}_{\bullet j_1}^\top \mathbf{x}_{\bullet j_2}}{\sigma^2}, \\ \mathbb{E} \left[-\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \sigma^2 \partial \beta_j} \right] &= \frac{1}{\sigma^4} \sum_{i=1}^n \frac{x_{ij} \mathbb{E}(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} = 0, \\ \mathbb{E} \left[-\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \sigma^2 \partial \tau_m} \right] &= \frac{1}{2\sigma^4} \sum_{i=1}^n \frac{\text{var} Y_i}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} z_{im} = \frac{1}{2\sigma^2} \sum_{i=1}^n z_{im} = \frac{\mathbf{z}_{\bullet m}^\top \mathbf{1}_n}{2\sigma^2}, \\ \mathbb{E} \left[-\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \beta_j \partial \tau_m} \right] &= \frac{1}{2\sigma^2} \sum_{i=1}^n \frac{x_{ij} \mathbb{E}(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})}{\exp(\mathbf{z}_i^\top \boldsymbol{\tau})} z_{im} = 0. \end{aligned}$$

Nyní již máme vše potřebné pro sestavení Fisherovy informační matice. Bude to matice řádu $q + 1 + (k + 1)$, kterou zapíšu blokově tak, aby cílový parametr $\boldsymbol{\tau}$ byl v levém horním rohu. Matice má tedy tvar

$$\mathbb{J}^{(n)}(\boldsymbol{\tau}, \sigma^2, \boldsymbol{\beta}) \stackrel{H_0}{=} \begin{pmatrix} \frac{\mathbb{Z}^\top \mathbb{Z}}{2} & \frac{\mathbb{Z}^\top \mathbf{1}_n}{2\sigma^2} & \mathbb{O}_{q \times (k+1)} \\ \frac{\mathbf{1}_n^\top \mathbb{Z}}{2\sigma^2} & \frac{n}{2\sigma^4} & \mathbf{0}_{1 \times (k+1)} \\ \mathbb{O}_{(k+1) \times q} & \mathbf{0}_{(k+1) \times 1} & \frac{\mathbb{X}^\top \mathbb{X}}{\sigma^2} \end{pmatrix}. \quad (4.7)$$

Matice (4.7) si rozdělme na bloky a použijeme dále značení zavedené v lemmatu 6. V našem případě máme takovouto situaci:

$$\begin{aligned} \mathbb{J}_{11}^{(n)} &= \frac{\mathbb{Z}^\top \mathbb{Z}}{2}, & \mathbb{J}_{12}^{(n)} &= \begin{pmatrix} \frac{\mathbb{Z}^\top \mathbf{1}_n}{2\sigma^2} & \mathbb{O}_{q \times (k+1)} \end{pmatrix}, \\ \mathbb{J}_{21}^{(n)} &= \begin{pmatrix} \frac{\mathbf{1}_n^\top \mathbb{Z}}{2\sigma^2} \\ \mathbb{O}_{(k+1) \times q} \end{pmatrix}, & \mathbb{J}_{22}^{(n)} &= \begin{pmatrix} \frac{n}{2\sigma^4} & \mathbf{0}_{1 \times (k+1)} \\ \mathbf{0}_{(k+1) \times 1} & \frac{\mathbb{X}^\top \mathbb{X}}{\sigma^2} \end{pmatrix}. \end{aligned}$$

Povšimneme si, že díky našim předpokladům o plných sloupcových hodnotnostech matic \mathbb{X} a \mathbb{Z} jsou matice $\mathbb{Z}^\top \mathbb{Z}$ a $\mathbb{X}^\top \mathbb{X}$ regulární, a tedy jsou regulární i $\mathbb{J}_{11}^{(n)}$ a $\mathbb{J}_{22}^{(n)}$.

Pro odvození testové statistiky budeme potřebovat $\mathbb{J}_{11.2}^{(n)}$ a následně pracovat s její inverzní maticí. Pokračujme tedy určením této matice. Využijeme přitom blokové diagonality matice $\mathbb{J}_{22}^{(n)}$, kterou zinvertujeme „po blocích“ (na tuto matici lze aplikovat lemma 6).

Postupnými úpravami dostáváme

$$\begin{aligned}
\mathbb{J}_{11.2}^{(n)} &= \mathbb{J}_{11}^{(n)} - \mathbb{J}_{12}^{(n)} \left(\mathbb{J}_{22}^{(n)} \right)^{-1} \mathbb{J}_{21}^{(n)} \\
&= \frac{\mathbb{Z}^\top \mathbb{Z}}{2} - \left(\frac{\mathbb{Z}^\top \mathbf{1}_n}{2\sigma^2} \quad \mathbb{O}^{q \times (k+1)} \right) \begin{pmatrix} \frac{2\sigma^4}{n} & \mathbf{0}^{1 \times (k+1)} \\ \mathbf{0}^{(k+1) \times 1} & \sigma^2 (\mathbb{X}^\top \mathbb{X})^{-1} \end{pmatrix} \begin{pmatrix} \frac{\mathbf{1}_n^\top \mathbb{Z}}{2\sigma^2} \\ \mathbb{O}^{(k+1) \times q} \end{pmatrix} \\
&= \frac{\mathbb{Z}^\top \mathbb{Z}}{2} - \frac{\mathbb{Z}^\top \mathbf{1}_n}{2\sigma^2} \cdot \frac{2\sigma^4}{n} \cdot \frac{\mathbf{1}_n^\top \mathbb{Z}}{2\sigma^2} \\
&= \frac{1}{2} \mathbb{Z}^\top \left(\mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \mathbb{Z} = \frac{1}{2} \mathbb{Z}^\top \mathbb{A} \mathbb{Z}.
\end{aligned}$$

Je zřejmé, že matice \mathbb{A} , kterou jsme si zavedli již před tvrzením 9, je symetrická ($\mathbb{A} = \mathbb{A}^\top$) a idempotentní ($\mathbb{A}^2 = \mathbb{A}$).

Označme dále $\mathbf{U}_1(\boldsymbol{\tau}, \sigma^2, \boldsymbol{\beta})$ vektor parciálních derivací funkce ℓ podle cílových proměnných $\boldsymbol{\tau}$. Platí

$$\mathbf{U}_1(\boldsymbol{\tau}, \sigma^2, \boldsymbol{\beta}) = \begin{pmatrix} \frac{\partial \ell(\boldsymbol{\tau}, \sigma^2, \boldsymbol{\beta})}{\partial \tau_1} \\ \vdots \\ \frac{\partial \ell(\boldsymbol{\tau}, \sigma^2, \boldsymbol{\beta})}{\partial \tau_q} \end{pmatrix} \stackrel{(4.3)}{\underset{H_0}{=}} \begin{pmatrix} \frac{1}{2\sigma^2} \sum_{i=1}^n \left[(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 - \sigma^2 \right] z_{i1} \\ \vdots \\ \frac{1}{2\sigma^2} \sum_{i=1}^n \left[(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 - \sigma^2 \right] z_{iq} \end{pmatrix}.$$

Nás bude dle (2.1) zajímat hodnota tohoto vektoru pro $\boldsymbol{\tau} = \mathbf{0}^q$ a maximálně věrohodné odhady $\tilde{\sigma}_n^2$ a $\tilde{\boldsymbol{\beta}}_n$, což dle (4.6) dává

$$\mathbf{U}_1(\mathbf{0}^q, \tilde{\sigma}_n^2, \tilde{\boldsymbol{\beta}}_n) = \begin{pmatrix} \frac{1}{2\tilde{\sigma}_n^2} \sum_{i=1}^n (u_i^2 - \tilde{\sigma}_n^2) z_{i1} \\ \vdots \\ \frac{1}{2\tilde{\sigma}_n^2} \sum_{i=1}^n (u_i^2 - \tilde{\sigma}_n^2) z_{iq} \end{pmatrix},$$

kde $u_i = Y_i - \mathbf{x}_i^\top \tilde{\boldsymbol{\beta}}_n = Y_i - \mathbf{x}_i^\top (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y} = Y_i - \tilde{Y}_i$, $i \in \{1, \dots, n\}$ jsou rezidua tvořící dohromady vektor \mathbf{u} . Připomeňme, že druhá mocnina eukleidovské normy tohoto vektoru je právě reziduální součet čtverců, neboli $RSS = \|\mathbf{u}\|^2$.

Vektor \mathbf{U}_1 lze dále upravit do jednoduššího tvaru. Ukážeme si dva možné způsoby, jak se na toto dívat. V prvním případě si definujme vektor \mathbf{v} následovně

$$v_i := u_i^2 - \tilde{\sigma}_n^2 = u_i^2 - \frac{RSS}{n} = \frac{1}{n} \sum_{j=1}^n (u_i^2 - u_j^2), \quad i \in \{1, \dots, n\}.$$

Vidíme, že i -tá složka vektoru \mathbf{v} je průměr rozdílů čtverce rezidua u_i se všemi ostatními. Pak lze vektor \mathbf{U}_1 zapsat jako

$$\mathbf{U}_1(\mathbf{0}^q, \tilde{\sigma}_n^2, \tilde{\boldsymbol{\beta}}_n) = \begin{pmatrix} \frac{1}{2\tilde{\sigma}_n^2} \sum_{i=1}^n v_i z_{i1} \\ \vdots \\ \frac{1}{2\tilde{\sigma}_n^2} \sum_{i=1}^n v_i z_{iq} \end{pmatrix} = \begin{pmatrix} \frac{1}{2\tilde{\sigma}_n^2} \mathbf{z}_{\bullet 1}^\top \mathbf{v} \\ \vdots \\ \frac{1}{2\tilde{\sigma}_n^2} \mathbf{z}_{\bullet q}^\top \mathbf{v} \end{pmatrix} = \frac{\mathbb{Z}^\top \mathbf{v}}{2\tilde{\sigma}_n^2} = \frac{n\mathbb{Z}^\top \mathbf{v}}{2RSS}. \quad (4.8)$$

Druhá možnost je definovat si pomocné veličiny

$$\bar{z}_m := \frac{1}{n} \sum_{i=1}^n z_{im} = \frac{1}{n} \mathbf{z}_{\bullet m}^\top \mathbf{1}_n, \quad m \in \{1, \dots, q\},$$

které dohromady tvoří sloupcový vektor $\bar{\mathbf{z}} = \frac{1}{n} \mathbf{Z}^\top \mathbf{1}_n$. Potom platí

$$\begin{aligned} \sum_{i=1}^n (u_i^2 - \tilde{\sigma}_n^2) z_{im} &= \sum_{i=1}^n u_i^2 z_{im} - \frac{RSS}{n} \sum_{i=1}^n z_{im} \\ &= \sum_{i=1}^n u_i^2 z_{im} - RSS \bar{z}_m = \sum_{i=1}^n u_i^2 (z_{im} - \bar{z}_m). \end{aligned}$$

A tedy se potom dá vektor \mathbf{U}_1 zapsat jako

$$\mathbf{U}_1 \left(\mathbf{0}^q, \tilde{\sigma}_n^2, \tilde{\boldsymbol{\beta}}_n \right) = \frac{1}{2\tilde{\sigma}_n^2} \sum_{i=1}^n u_i^2 (\mathbf{z}_i - \bar{\mathbf{z}}). \quad (4.9)$$

Pomocí těchto průměrů \bar{z}_m lze také upravit naši matici $\mathbb{J}_{11.2}$, a to do tvaru

$$\mathbb{J}_{11.2}^{(n)} = \frac{1}{2} \mathbf{Z}^\top \mathbf{A} \mathbf{Z} = \frac{1}{2} (\mathbf{Z} - \mathbf{1}_n \bar{\mathbf{z}}^\top)^\top (\mathbf{Z} - \mathbf{1}_n \bar{\mathbf{z}}^\top). \quad (4.10)$$

O platnosti této rovnosti se jednoduše přesvědčíme pomocí následujících úprav:

$$\begin{aligned} (\mathbf{Z} - \mathbf{1}_n \bar{\mathbf{z}}^\top)^\top (\mathbf{Z} - \mathbf{1}_n \bar{\mathbf{z}}^\top) &= \left(\mathbf{Z} - \mathbf{1}_n \frac{1}{n} (\mathbf{Z}^\top \mathbf{1}_n)^\top \right)^\top \left(\mathbf{Z} - \mathbf{1}_n \frac{1}{n} (\mathbf{Z}^\top \mathbf{1}_n)^\top \right) \\ &= \left(\mathbf{Z} - \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} \mathbf{Z} \right)^\top \left(\mathbf{Z} - \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} \mathbf{Z} \right) \\ &= (\mathbf{A} \mathbf{Z})^\top (\mathbf{A} \mathbf{Z}) = \mathbf{Z}^\top \mathbf{A}^2 \mathbf{Z} = \mathbf{Z}^\top \mathbf{A} \mathbf{Z}. \end{aligned}$$

Nyní už si můžeme definovat testovou statistiku LM podle (2.1) následovně

$$LM := \left[\mathbf{U}_1 \left(\mathbf{0}^q, \tilde{\sigma}_n^2, \tilde{\boldsymbol{\beta}}_n \right) \right]^\top \left[\mathbb{J}_{11.2}^{(n)} \left(\mathbf{0}^q, \tilde{\sigma}_n^2, \tilde{\boldsymbol{\beta}}_n \right) \right]^{-1} \mathbf{U}_1 \left(\mathbf{0}^q, \tilde{\sigma}_n^2, \tilde{\boldsymbol{\beta}}_n \right).$$

Zde použijeme našich úvah výše, abychom viděli, jak testová statistika LM vypadá. Kombinací dvojic rovnic (4.8), (4.10) a (4.9), (4.10) dostaneme, že LM má tvar

$$\begin{aligned} LM &= \left(\frac{\mathbf{Z}^\top \mathbf{v}}{2\tilde{\sigma}_n^2} \right)^\top \left(\frac{1}{2} \mathbf{Z}^\top \mathbf{A} \mathbf{Z} \right)^{-1} \frac{\mathbf{Z}^\top \mathbf{v}}{2\tilde{\sigma}_n^2} = \frac{1}{2(\tilde{\sigma}_n^2)^2} \mathbf{v}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{A} \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{v} \\ &= \frac{1}{2(\tilde{\sigma}_n^2)^2} \left[\sum_{i=1}^n u_i^2 (\mathbf{z}_i - \bar{\mathbf{z}}) \right]^\top \left[(\mathbf{Z} - \mathbf{1}_n \bar{\mathbf{z}}^\top)^\top (\mathbf{Z} - \mathbf{1}_n \bar{\mathbf{z}}^\top) \right]^{-1} \left[\sum_{i=1}^n u_i^2 (\mathbf{z}_i - \bar{\mathbf{z}}) \right]. \end{aligned} \quad (4.11)$$

Zbývá už jen ověřit předpoklady věty 7. Už víme, že systém hustot \mathcal{F}_{q+k+2} je regulární a že skutečná hodnota parametrů leží v nějakém otevřeném intervalu obsaženém v Θ . Dále je vzhledem k naší parametrizaci splněn předpoklad (O2) z vlastnosti hustoty normálního rozdělení. Existence třetích derivací (A1) plyne

jednoduše dalším derivováním druhých derivací funkce ℓ , v čemž nám nic nebrání. Předpoklad (A2) jsme již oprávněně používali při výpočtu Fisherovy informační matice. Předpoklad (A3) plyne opět z vlastnosti normálního rozdělení. Z rovnosti (4.7) plyne, že Fisherova matice je spojitá v každém $\boldsymbol{\theta} \in \Theta$, tedy i ve skutečné hodnotě tohoto parametru. Tím jsme již ověřili předpoklady.

Věta 7 mi tedy dává, že za platnosti nulové hypotézy $H_0 : \boldsymbol{\tau} = \mathbf{0}^q$ platí

$$LM \xrightarrow{D} \chi_q^2, \quad n \rightarrow \infty, \quad (4.12)$$

kde počet stupňů volnosti q odpovídá počtu cílových parametrů. □

Náš skórový test tedy na asymptotické hladině významnosti $\alpha \in (0, 1)$ zamítne hypotézu H_0 právě tehdy, když $LM \geq \chi_q^2(1 - \alpha)$.

Tento test navrhli v roce 1979 Breusch a Pagan ve svém článku Breusch a Pagan (1979). Ukázalo se ovšem, že tento test je velmi citlivý na porušení předpokladu normálního rozdělení. Proto o dva roky později navrhl Koenker v článku Koenker (1981) úpravu, která tuto testovou statistiku *studentizovala*, takže se dá použít v případě, že si nejsme jisti s platností předpokladu normality.

Úprava spočívá v tom, že v testové statistice

$$LM = \frac{1}{2(\tilde{\sigma}_n^2)^2} \mathbf{v}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{A} \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{v}$$

nahradíme $(\tilde{\sigma}_n^2)^2$ odhadem rozptylu veličin ε_i^2 pomocí

$$\tilde{\varphi}_n = \frac{1}{n} \sum_{i=1}^n (u_i^2 - \tilde{\sigma}_n^2)^2 = \frac{1}{n} \sum_{i=1}^n v_i^2 = \frac{\|\mathbf{v}\|^2}{n}.$$

Kapitola 5

Numerické studie

V této kapitole se na základě počítačových simulací v prostředí **R** (R Core Team, 2016) podíváme, jaké vlastnosti mají testy odvozené v kapitolách 3 a 4 v závislosti na rozsahu náhodného výběru n a na splnění předpokladu normality. Veškeré testy budeme provádět na asymptotické hladině 5 %.

5.1 Experimentální porovnání statistik B a LR

V této části se vrátíme zpět k testům homoskedasticity v modelu analýzy rozptylu jednoduchého třídění s obecně různými rozptyly v jednotlivých skupinách uvedeným v kapitole 3, které byly založeny na dvojici statistik

$$LR = n \left[\log \left(\frac{1}{n} \sum_{i=1}^I RSS_i \right) - \sum_{i=1}^I \frac{n_i}{n} \log \left(\frac{RSS_i}{n_i} \right) \right],$$
$$B = \frac{n-I}{c} \left[\log \left(\frac{1}{n-I} \sum_{i=1}^I RSS_i \right) - \sum_{i=1}^I \frac{n_i-1}{n-I} \log \left(\frac{RSS_i}{n_i-1} \right) \right],$$
$$\text{kde } c = 1 + \frac{1}{3(I-1)} \left(\sum_{i=1}^I \frac{1}{n_i-1} - \frac{1}{n-I} \right).$$

Obě statistiky by se při dostatečně velkém n měly řídit rozdělením χ_{I-1}^2 , kde I je počet skupin. Na konkrétním příkladě tedy prozkoumáme, jak velké n je zapotřebí, abychom toto mohli tvrdit. Dále se podíváme, zda testy založené na těchto statistikách dodržují předepsanou hladinu. Nakonec prozkoumáme jejich schopnost rozpoznat neplatnou hypotézu, tedy experimentálně prozkoumáme jejich sílu. U těchto testů nebudeme zkoumat jejich chování při nesplnění předpokladu normality.

Naše simulace provedeme vždy pro $I = 5$ skupin. V podkapitole 3.1 jsme zavedli model tak, že přiřazování napozorovaných hodnot veličiny Y do skupin probíhá náhodně. Regresor určující příslušnost do skupiny je generován z rozdělení $\text{Mult}_5(n, \mathbf{p})$, kde \mathbf{p} si nastavíme na hodnoty

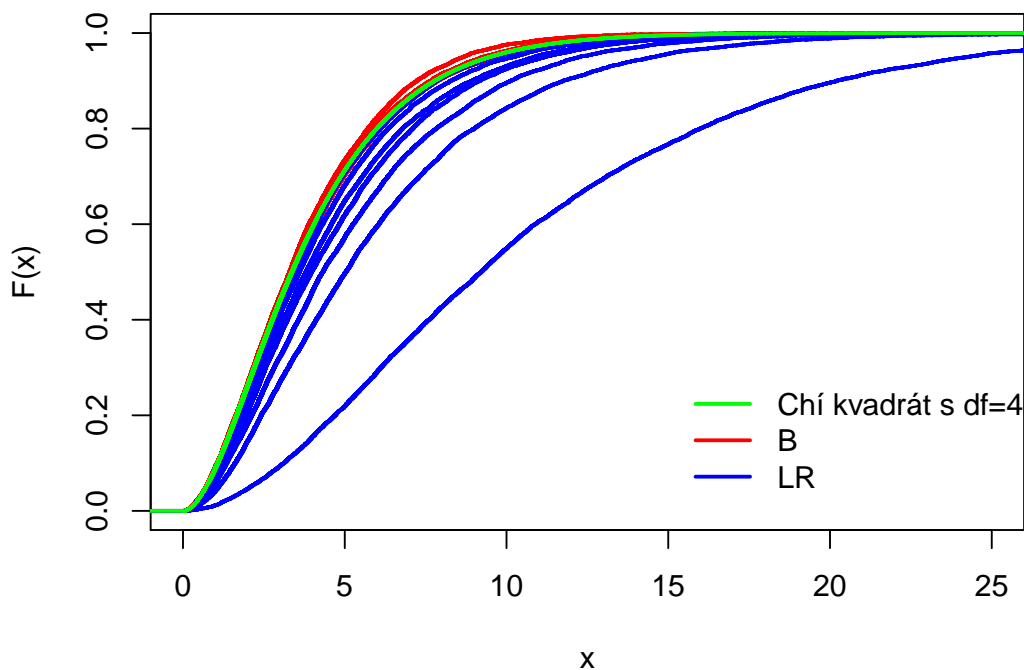
$$\mathbf{p} = \left(\frac{1}{15}, \frac{2}{15}, \frac{3}{15}, \frac{4}{15}, \frac{5}{15} \right)^\top.$$

Pro výpočet statistiky B musíme zajistit, aby v každé skupině byly alespoň dvě pozorování.

Nejprve začneme s případem, kdy je hypotéza homoskedasticity platná. Všechna pozorování y tedy budeme generovat z normálního rozdělení se stejnou hodnotou rozptylu $\sigma^2 = 1$. Střední hodnoty v jednotlivých skupinách nastavíme na různé hodnoty, v tomto případě nastavíme na hodnoty $\boldsymbol{\mu} = (-1, 0, 1, 2, 3)^\top$. Dále si určíme jednotlivé rozsahy výběrů $n \geq 5 \cdot 2$, pro které budeme chtít generovat data. Pro každé takové n potom 10 000-krát nasimulujeme náhodný výběr z tohoto rozdělení. Pro jednotlivou simulaci napočítáme hodnoty testových statistik B a LR a podíváme se, zda překročily kritickou hodnotu $\chi_4^2(0.95) \doteq 9.4877$. Statistiku LR spočteme jednoduše z našeho vzorce (3.21), na výpočet statistiky B lze použít funkce `bartlett.test` z balíčku `stats`.

Naše simulace provedeme pro $n \in \{10, 30, 50, 80, 100, 200, 400\}$. Výsledky experimentů jsou shrnuty na obrázku 5.1 a v tabulce 5.1. Obrázek 5.1 se zaměřuje na zjištění skutečného rozdělení jednotlivých statistik za nulové hypotézy. Můžeme si všimnout, že empirická distribuční funkce pro statistiku B se téměř překrývá s distribuční funkcí rozdělení χ_4^2 , a to už i pro $n = 10$. Na druhou stranu je vidět, že empirická distribuční funkce statistiky LR se pro nízké hodnoty n výrazně odlišuje od distribuční funkce rozdělení χ_4^2 , ale s rostoucím n se k této funkci blíží a při $n = 400$ už s ní také splývá.

Tabulka 5.1 uvádí, jaký byl podíl zamítnutí nulové hypotézy v 10 000 opakováních při jednotlivých rozsazích n . Vidíme, že Bartlettova statistika už při $n = 30$ hladiny 0.05 téměř dosáhla. Na druhou stranu statistika LR zamítala pro nízká n daleko častěji a hladině 0.05 se přiblížila až pro $n = 400$.



Obrázek 5.1: Grafy empirických distribučních funkcí statistik B a LR za platnosti hypotézy při rostoucím rozsahu výběru.

n	B	LR
10	0.0328	0.4815
30	0.0464	0.1783
50	0.0497	0.1228
80	0.0516	0.0899
100	0.0529	0.0822
200	0.0527	0.0632
400	0.0471	0.0528

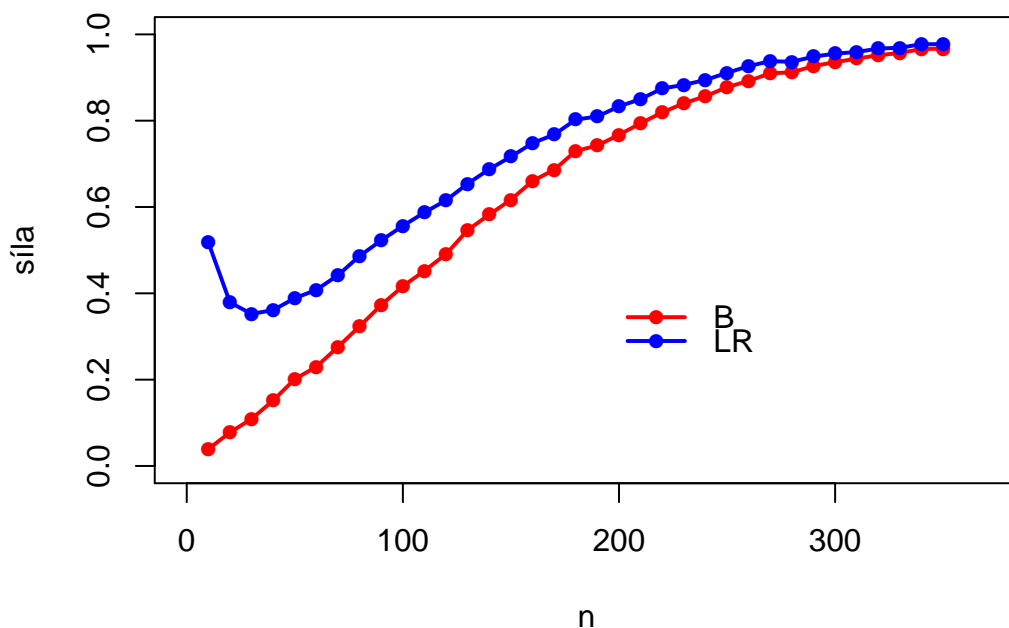
Tabulka 5.1: Naměřené hladiny statistik B a LR za platnosti nulové hypotézy.

Konvergence statistiky LR k limitnímu rozdělení χ_4^2 je viditelně pomalejší než v případě statistiky B . Toto je v souladu se zjištěními prezentovanými v literatuře, které jsme uvedli v kapitole 3.

Nyní budeme zkoumat sílu těchto statistik. Jediné, co na našem současném modelu musíme změnit, je platnost homoskedasticity. Budeme tedy generovat pozorování stejným postupem, jediné, co změníme, je nastavení rozptylu v jednotlivých skupinách. Nastavme tedy parametr směrodatné odchylky v jednotlivých skupinách na $\sigma = (1, 1.2, 1.4, 1.6, 1.8)^\top$.

Abychom prozkoumali sílu těchto statistik, provedeme podrobnější simulace. Budeme uvažovat rozsahy náhodného výběru z množiny

$$N = \{10, 20, 30, 40, 50, \dots, 340, 350\}.$$



Obrázek 5.2: Experimentální síla statistik B a LR při různém rozsahu výběru.

Pro každé $n \in N$ provedeme 10 000 simulací a zaznamenanáme, kolikrát jsme se (správně) rozhodli zamítnout nulovou hypotézu homoskedasticity. Naše výsledky jsou zaznamenány na obrázku 5.2. Můžeme vidět, že u obou statistik dochází s rostoucím n ke zvětšování podílu zamítnutých hypotéz, dokonce se pro vysoké n blíží k jedné. Statistika LR zjevně zamítá neplatnou hypotézu homoskedasticity častěji než Bartlettova statistika, ovšem zásadnější je pro nás dodržení hladiny. Tu LR v našem modelu dodržuje až při n blízkém 400, viz tabulka 5.1, jenže pro takové n už jsou síly B a LR srovnatelné.

Na tomto modelu jsme se tedy přesvědčili, že statistika B má skutečně lepší asymptotické vlastnosti než námi odvozená statistika LR .

V příloze pod položkou 1) uvádím skript, který jsem stvořil při provádění simulací. Manipulací s tímto skriptem lze nejen zreprodukovat naše výsledky, ale také zkoumat dále vlastnosti statistik B a LR za volby jiných parametrů.

5.2 Breusch-Paganův test

Nyní budeme vyšetřovat vlastnosti testové statistiky LM odvozené v kapitole 4. V tomto případě budeme zkoumat, zda naše statistika dodržuje předepsanou hladinu, už se ale nebudeme zajímat o její sílu. Navíc prověříme, jak (ne)dodržuje předepsanou hladinu, není-li splněn předpoklad o normálním rozdělení. V obou případech budeme navíc zkoumat vliv *studentizace* této statistiky, která byla navržena v práci Koenker (1981), viz poznámka v závěru kapitoly 4. Naše simulace provedeme na následujícím modelu.

Pro zvolený rozsah výběru n budeme uvažovat matici konstant \mathbb{X} o třech sloupcích o délce n , přičemž první z nich bude obsahovat pouze prvky 1, druhý bude tvořen posloupností n čísel tvořící ekvidistantní dělení na intervalu $[0, 1]$ a třetí sloupec bude střídavě obsahovat prvky 0 a 1. Neboli pro n sudé

$$\mathbb{X}^\top = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & \dots & 1 & 1 \\ 0 & \frac{1}{n-1} & \frac{2}{n-1} & \frac{3}{n-1} & \frac{4}{n-1} & \dots & \frac{n-2}{n-1} & 1 \\ 0 & 1 & 0 & 1 & 0 & \dots & 0 & 1 \end{pmatrix}.$$

Sice hodláme naše data generovat za předpokladu pevného rozptylu, ale abychom mohli testovat, tak musíme určit na jakých pomocných regresorech \mathbb{Z} by mohl obecně rozptyl záviset. My zvolíme za pomocné regresory prostřední sloupec matice \mathbb{X} , tedy vektor

$$\mathbf{z} = \left(0, \frac{1}{n-1}, \frac{2}{n-1}, \frac{3}{n-1}, \frac{4}{n-1}, \dots, \frac{n-2}{n-1}, 1 \right)^\top.$$

Tedy parametrizace rozptylu by za platnosti alternativy u veličiny Y_i vypadala jako $\sigma^2 \exp(z_i \cdot \tau)$, $i \in \{1, \dots, n\}$, kde $\tau \in \mathbb{R}$ je neznámý jednorozměrný parametr, o němž budeme testovat, zda se rovná nule. Proto asymptotické rozdělení statistiky LM by mělo být χ_1^2 . Příslušná kritická hodnota pro asymptotickou hladinu 5 % je $\chi_1^2(0.95) \doteq 3.8415$. Dále nastavíme hodnoty ostatních parametrů na

$$\boldsymbol{\beta} = (0, 1, 1)^\top, \quad \sigma = 2.$$

Nyní budeme moct přistoupit ke generování veličin $y_i, i \in \{1, \dots, n\}$, které budeme provádět následovně

$$y_i := \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, i \in \{1, \dots, n\},$$

kde \mathbf{x}_i a $\boldsymbol{\beta}$ již známe, ale veličiny ε_i si nagenerujeme. Budeme uvažovat 3 případy, podle toho, z jakého rozdělení ε_i pochází. Nejprve podle předpokladu provedeme pro $\varepsilon_i \sim \mathbf{N}(0, 2^2)$. Pro rozpor s předpokladem normality použijeme rozdělení $\text{LN}(0, 2^2)$ (logaritmicko-normální) a t_3 (studentovo rozdělení se 3 stupni volnosti). Jakmile budeme mít všechny tyto veličiny nagenerovány, tak můžeme přistoupit k samotnému testování. Využijeme k tomu balíček `lmtest` (Zeileis a Hothorn, 2002), konkrétně funkce `bptest`, která zvládá napočítat, jak obyčejnou statistiku LM , tak i její studentizovanou verzi.

Pro každé $n \in \{30, 50, 80, 100, 150, 200, 300, 400\}$ nagenerujeme 10 000 náhodných výběrů o rozsahu n pro každé z těchto tří rozdělení a pro každý náhodný výběr z daného rozdělení napočítáme, jak studentizovanou, tak obecnou verzi statistiky LM . Zaznamenáme si, kolikrát jsme celkově přesáhli kritickou hodnotu $\chi_1^2(0.95) \doteq 3.8415$ v jednotlivých případech. Naše výsledky shrnuje tabulka 5.2.

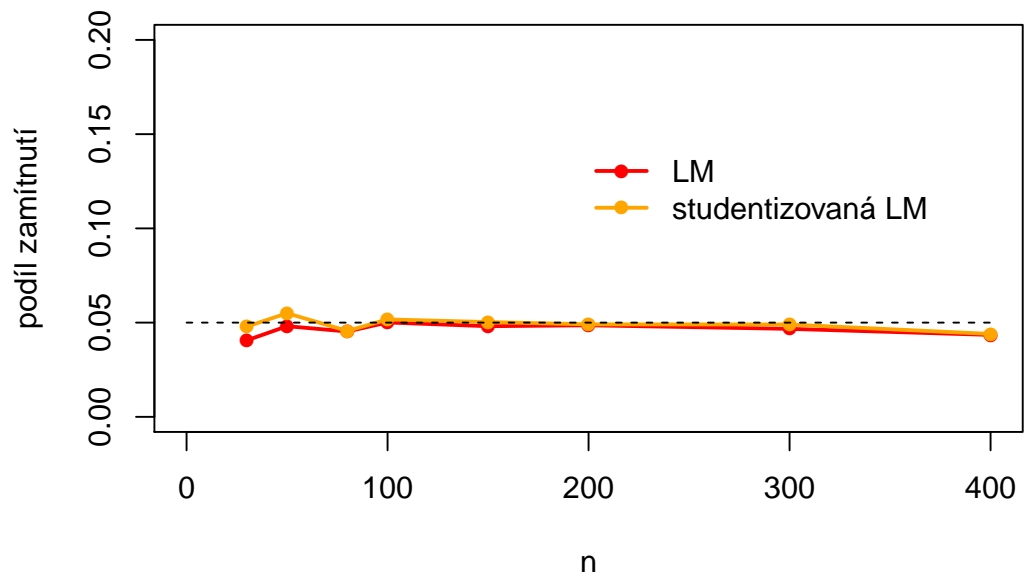
Pro jednotlivá rozdělení si tato data zakreslíme do grafu, abychom názorně viděli naše výsledky. Nejprve se podívejme na výsledky normálního rozdělení, viz obrázek 5.3. Vidíme, že v tomto případě statistika LM v obou případech vesměs dodržuje předepsanou asymptotickou hladinu 0.05.

Pro logaritmicko normální rozdělení to už neplatí. Na obrázku 5.4 můžeme vidět, že nestudentizovaná statistika překračuje kritickou hodnotu ve více jak polovině pozorování, a navíc s větším rozsahem výběru dokonce dosažená hladina roste. Na druhou stranu studentizovaná statistika už předepsanou hladinu 0.05 nepřekračuje, ale z tabulky 5.2 vidíme, že tato dosažená hladina se spíše blíží hodnotě 0.025.

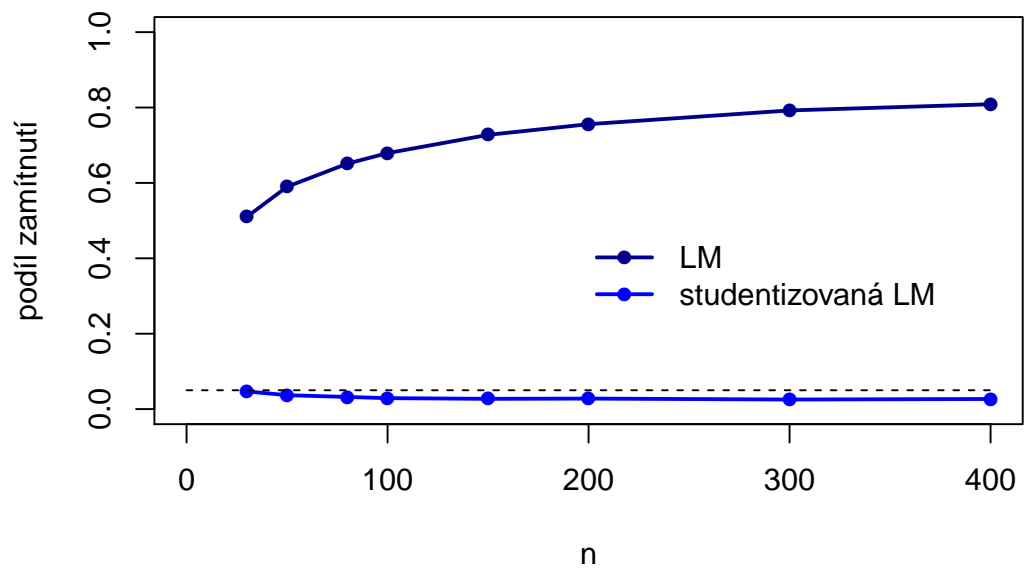
Studentovo rozdělení t_3 je na tom velmi podobně jako logaritmicko-normální rozdělení. Na obrázku 5.5 vidíme, že dosažená hladina u nestudentizované statistiky s větším n roste zhruba k hodnotě 0.45. Studentizovaná statistika už předepsanou hladinu 0.05 nepřekračuje, ale také je vidět, že se spíše blíží hodnotě 0.04 než hodnotě 0.05.

n	Nestudentizovaná LM			Studentizovaná LM		
	N	LN	t_3	N	LN	t_3
30	0.0406	0.5100	0.1844	0.0479	0.0469	0.0465
50	0.0481	0.5898	0.2440	0.0550	0.0369	0.0417
80	0.0453	0.6512	0.2947	0.0456	0.0322	0.0409
100	0.0501	0.6791	0.3145	0.0517	0.0290	0.0406
150	0.0481	0.7275	0.3505	0.0503	0.0272	0.0408
200	0.0486	0.7558	0.3866	0.0491	0.0279	0.0421
300	0.0467	0.7923	0.4158	0.0490	0.0254	0.0431
400	0.0435	0.8087	0.4396	0.0440	0.0268	0.0376

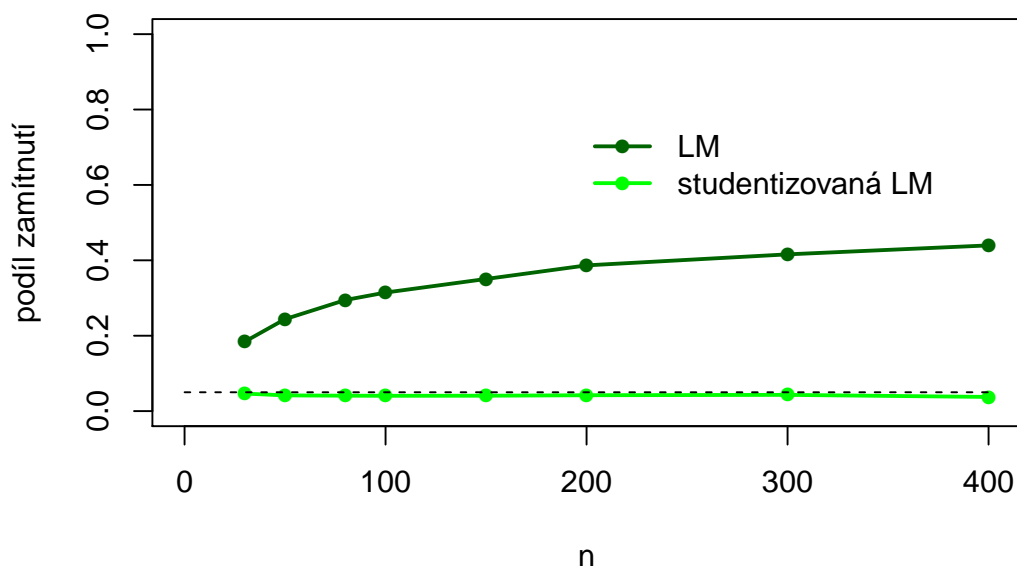
Tabulka 5.2: Dosažené hladiny testů pro jednotlivá rozdělení v závislosti na rozsahu výběru n .



Obrázek 5.3: Graf podílu zamítnutí při N rozdělení v závislosti na n .



Obrázek 5.4: Graf podílu zamítnutí při LN rozdělení v závislosti na n .

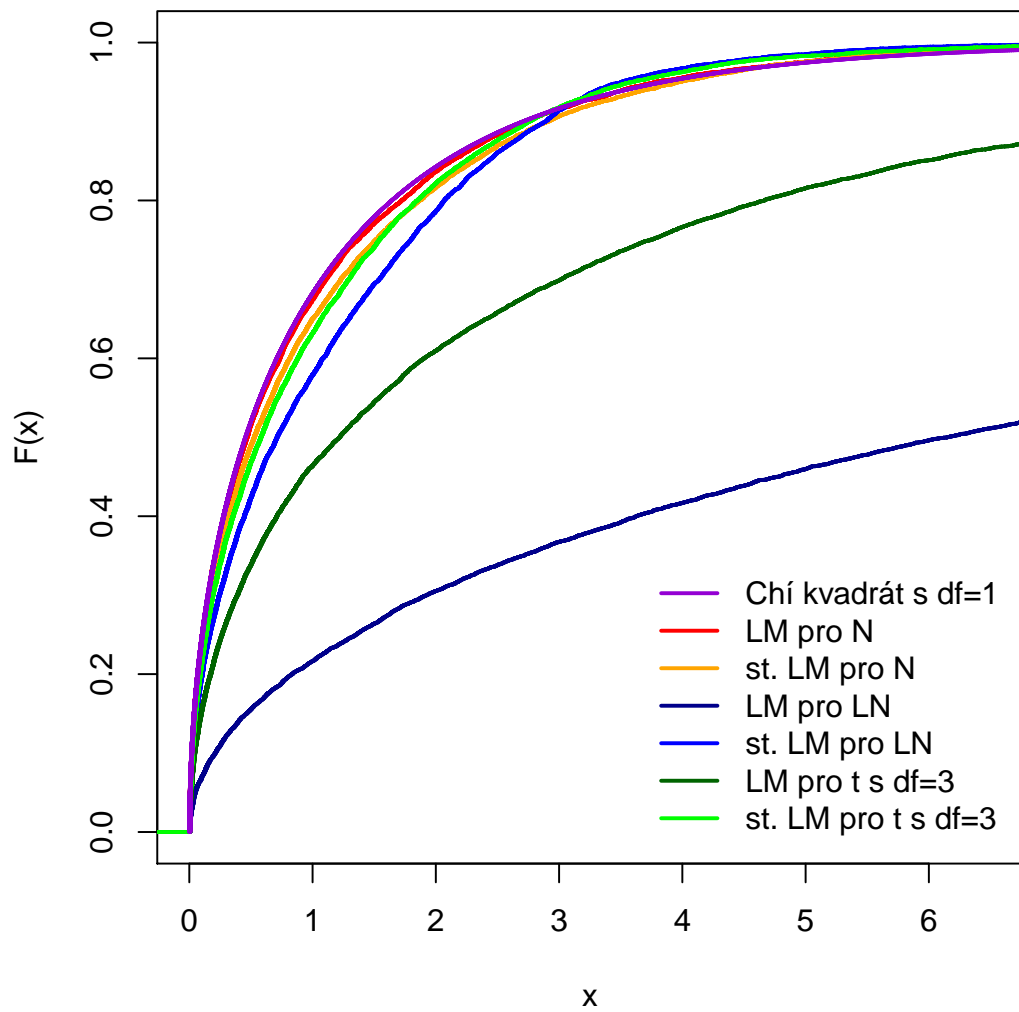


Obrázek 5.5: Graf podílu zamítnutí při t_3 rozdělení v závislosti na n .

Můžeme se také podívat, jak vypadají empirické distribuční funkce těchto statistik v našich šesti případech. Tyto funkce si ale pro přehlednost vykreslíme pro jedinou hodnotu n . Na obrázku 5.6, který je vykreslen pro $n = 50$, vidíme, že pro data vygenerovaná z normálního rozdělení se graf empirické distribuční funkce statistiky LM podobá grafu distribuční funkce rozdělení χ_1^2 . Pro nenormální data se studentizovanou statistikou už pozorujeme drobné odchylky a pro nestudentizovanou vidíme značný rozdíl.

Těmito simulacemi jsme se tedy přesvědčili, že za platného předpokladu normálního rozdělení naše testování homoskedasticity přibližně dodržuje předepsanou hladinu. Pokud tento předpoklad splněn není, tak LM sama o sobě dává nežádoucí výsledky, ale po její studentizaci jsou výsledky příznivější.

Mezi přílohami pod položkou 2) je uveden také vypracovaný skript pro prostředí [R](#). Jeho modifikací lze zkoumat vlastnosti Breusch-Paganova testu i pro jiné parametry, než s jakými jsme pracovali zde.



Obrázek 5.6: Grafy empirických distribučních funkcí statistiky LM pro různá rozdělení při rozsahu výběru $n = 50$.

Závěr

Práce představila klasický lineární model a zobecnila jej na model heteroskedastický. K nezávisle proměnným (regresorům) bylo v těchto modelech přistupováno jako k náhodným vektorům. Následně jsme si připomněli zásadní poznatky z teorie maximální věrohodnosti a ukázali si, jak je aplikovat v lineárním modelu. Zejména jsme se potom zajímali o testy s rušivými parametry, které z teorie maximální věrohodnosti vycházejí. Na základě těchto teoretických podkladů jsme si odvodili dva základní testy předpokladu homoskedasticity. Nejprve jsme si zobecnili model analýzy rozptylu jednoduchého třídění ve smyslu heteroskedastického lineárního modelu a odvodili si test založený na věrohodnostním poměru pro shodnost rozptylů. Naši výslednou statistiku jsme porovnali s lehce modifikovanou statistikou navrženou Bartletttem. Na základě provedených simulací jsme se přesvědčili, že tato modifikace skutečně disponuje lepšími asymptotickými vlastnostmi než námi odvozená statistika. Ve druhém případě jsme zkoumali obecný normální heteroskedastický model, kde rozptyl chyb tohoto modelu závisel na doprovodných veličinách. Zde jsme po dlouhé řadě výpočtů došli ke skórové testové statistice a uvedli jsme její studentizovanou modifikaci. V následných numerických studiích jsme si ověřili její asymptotické vlastnosti a zkoumali, jak se tato statistika chová, není-li splněn předpoklad normálního rozdělení. Ukázalo se, že v tomto případě má její studentizovaná verze mnohem příznivější vlastnosti.

Literatura

- ANDĚL, J. (2007). *Základy matematické statistiky*. Druhé opravené vydání. Matfyzpress, Praha. ISBN 978-80-7378-162-0.
- BARTLETT, M. S. (1937). Properties of Sufficiency and Statistical Tests. *Proceedings of the Royal Statistical Society*, **A 160**, 268–282.
- BREUSCH, T. S. a PAGAN, R. (1979). A Simple Test for Heteroscedasticity and Random Coefficient Variation. *Econometrica*, **47**(5), 1287–1294.
- KOENKER, R. (1981). A Note on Studentizing a Test for Heteroscedasticity. *Journal of Econometrics*, **17**(1), 107–112.
- R CORE TEAM (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- ZEILEIS, A. a HOTHORN, T. (2002). Diagnostic checking in regression relationships. *R News*, **2**(3), 7–10. URL <http://CRAN.R-project.org/doc/Rnews/>.
- ZVÁRA, K. (2008). *Regrese*. Druhé opravené vydání. Matfyzpress, Praha. ISBN 978-80-7378-041-8.

Seznam obrázků

5.1	Grafy empirických distribučních funkcí statistik B a LR za platnosti hypotézy při rostoucím rozsahu výběru.	34
5.2	Experimentální síla statistik B a LR při různém rozsahu výběru. .	35
5.3	Graf podílu zamítnutí při N rozdělení v závislosti na n	38
5.4	Graf podílu zamítnutí při LN rozdělení v závislosti na n	38
5.5	Graf podílu zamítnutí při t_3 rozdělení v závislosti na n	39
5.6	Grafy empirických distribučních funkcí statistiky LM pro různá rozdělení při rozsahu výběru $n = 50$	40

Seznam tabulek

5.1	Naměřené hladiny statistik B a LR za platnosti nulové hypotézy.	35
5.2	Dosažené hladiny testů pro jednotlivá rozdělení v závislosti na rozsahu výběru n .	37

Přílohy

V kapitole 5 jsme pracovali s výpočetním prostředím **R**. Příslušná implementace našich provedených simulací (zvláště pro Bartlettův test a zvláště pro Breusch-Paganův test) je uvedena v následujících dvou skriptech, které jsou k dispozici ke stažení ze Studijního informačního systému.

Názvy jednotlivých souborů:

- 1) `skript-Bartlett.R`
- 2) `skript-Breusch-Pagan.R`