



**MATEMATICKO-FYZIKÁLNÍ
FAKULTA**
Univerzita Karlova

BAKALÁŘSKÁ PRÁCE

Iva Hamerníková

Analýza rozptylu s náhodnými efekty

Katedra pravděpodobnosti a matematické statistiky

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

Studijní program: Matematika

Studijní obor: Obecná matematika

Praha 2018

Prohlašuji, že jsem tuto bakalářskou práci vypracoval(a) 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 má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle §60 odst. 1 autorského zákona.

V dne

Podpis autora

Chtěla bych poděkovat vedoucímu práce doc. RNDr. Arnoštu Komárkovi, Ph.D. za veškeré rady, konzultace a trpělivost při vedení práce.

Název práce: Analýza rozptylu s náhodnými efekty

Autor: Iva Hamerníková

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 popisem a odvozením metody analýzy rozptylu s náhodnými efekty. Nejprve uvedeme souhrn poznatků z teorie pravděpodobnosti, které budou důležité v dalším odvozování. Poté zavedeme model jednoduchého třídění s pevnými efekty a navrhneme testovou statistiku pro test shody středních hodnot skupin. V další části zavedeme model jednoduchého třídění s náhodnými efekty a odvodíme vlastnosti pozorování v tomto modelu. Za předpokladu vyváženého třídění definujeme součty čtverců a odvodíme jejich vlastnosti, díky kterým je pak můžeme použít k sestavení testové statistiky pro testování shody podmíněných středních hodnot skupin. Na závěr práce budeme pomocí simulací v programu R ověřovat, jak test analýzy rozptylu s náhodnými efekty dodržuje hladinu při porušení předpokladu normality.

Klíčová slova: analýza rozptylu, ANOVA, náhodné efekty

Title: Analysis of Variance with Random Effects

Author: Iva Hamerníková

Department: Department of Probability and Mathematical Statistics

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

Abstract: The aim of this thesis is to describe and derive the test of analysis of variance with random effects. At first we introduce a summary of results from the theory of probability which will be important in future derivations. Then we define the one-way classification model with fixed effects and propose the test statistics to test the equality of group means. In the following part we define the one-way classification model with random effects and derive properties of observations in this model. Under the assumption of balanced data we define sums of squares and derive their properties, which allow us to use them to create the test statistic. Finally we will use simulations in R to verify whether the ANOVA test with random effects observes the significance level when normality assumptions are violated.

Keywords: analysis of variance, ANOVA, random effects

Obsah

Seznam použitých zkratk	2
Úvod	3
1 Souhrn z teorie pravděpodobnosti	4
1.1 Podmíněná střední hodnota a rozptyl	4
1.2 χ^2 rozdělení a kvadratické formy	5
2 Analýza rozptylu s pevnými efekty	7
2.1 Jednoduché třídění	7
3 Analýza rozptylu s náhodnými efekty	11
3.1 Jednoduché třídění - zavedení modelu	11
3.2 Odvození vlastností	12
3.3 Součty čtverců	14
3.4 Testování pro analýzu rozptylu s náhodnými efekty	16
4 Simulace	17
Závěr	21
Seznam použité literatury	22

Seznam použitých zkratek

\mathbb{I}_k	jednotková matice o rozměrech $k \times k$
$\mathbf{1}_k$	sloupcový vektor tvořený k prvky 1
$\text{diag}(a_1, \dots, a_n)$	diagonální matice s prvky a_1, \dots, a_n na diagonále
$\text{rank}(\mathbb{A})$	hodnost matice \mathbb{A}
$\text{tr}(\mathbb{A})$	stopa matice \mathbb{A}
$\text{Var}(\mathbf{X})$	varianční matice pro náhodný vektor \mathbf{X}

Úvod

Analýza rozptylu, z anglického analysis of variance často zkracovaná jako ANOVA, je hojně používaná metoda pro porovnávání skupinových výběrů, respektive testování nulové hypotézy o rovnosti středních hodnot těchto skupin. Princip této metody spočívá v myšlence, že rozdělíme celkovou variabilitu pozorovaných dat na variabilitu mezi jednotlivými skupinami a na variabilitu uvnitř skupin. Tyto dvě veličiny pak využíváme k sestrojení testové statistiky, na základě které rozhodujeme o platnosti nulové hypotézy. Cílem této práce je podrobné odvození speciálního typu, analýzy rozptylu s náhodnými efekty.

V první kapitole shrneme užitečná tvrzení a vztahy z teorie pravděpodobnosti, jejichž výsledky pak budeme dále v práci používat při výpočtech a odvozeních. Nejprve zavedeme podmíněnou střední hodnotu a rozptyl, uvedeme vlastnosti střední hodnoty a větu o rozkladu nepodmíněného rozptylu. V další části pak připomeneme definici χ^2 rozdělení a zaměříme se na tvrzení o rozdělení a nezávislosti kvadratických forem.

V druhé kapitole zavedeme a okomentujeme standardní model jednoduchého třídění s pevnými efekty. Definujeme základní značení a odvodíme vlastnosti součtů čtverců, na základě kterých navrhne testovou statistiku a kritický obor pro test shody středních hodnot skupin.

V další části následuje zavedení modelu s náhodnými efekty na příkladu s populací myší. Podrobně odvodíme vztahy pro podmíněnou střední hodnotu a rozptyl v tomto rozdělení, porovnáme je s předchozím případem. Za předpokladu vyváženého třídění definujeme součty čtverců a odvodíme jejich vlastnosti, konkrétně střední hodnotu a za dodatečných předpokladů i jejich rozdělení. V závěru kapitoly využijeme předchozích výsledků k navržení testové statistiky pro test shody podmíněných středních hodnot skupin.

Na závěr práce budeme pomocí simulací v programu R zkoumat co se stane při porušení základního předpokladu normality pozorování. Provedeme simulace s různým počtem skupin a jiným než normálním rozdělením a budeme sledovat, zda test dodržuje hladinu.

1. Souhrn z teorie pravděpodobnosti

Při odvozování vlastností pozorování a testových statistik v dalších kapitolách budeme často používat různé závěry z teorie pravděpodobnosti, které pro přehlednost uvedeme v této kapitole.

1.1 Podmíněná střední hodnota a rozptyl

Nejprve shrneme základní poznatky o podmíněné střední hodnotě z teorie pravděpodobnosti. Věty a tvrzení v této části jsou převzaty z knih Lachout (2004) a Anděl (2007a).

K zavedení podmíněné střední hodnoty budeme potřebovat nejprve definovat podmíněnou hustotu.

Definice 1 (Podmíněná hustota). *Nechť X, Y jsou reálné náhodné veličiny, dále necht existuje spojitě rozdělení (X, Y) s hustotou $f(x, y)$, absolutně spojitou vůči Lebesgueově míře, necht existuje $f(y)$, marginální hustota Y . Jako podmíněnou hustotu x při y označíme takovou funkci $f(x|y)$, která pro všechna $A, B \in \mathcal{B}$, kde \mathcal{B} je borelovská σ -algebra, splňuje:*

$$P(X \in A, Y \in B) = \int_A \int_B f(x, y) dy dx = \int_A \int_B f(x|y) f(y) dy dx.$$

Podmíněnou hustotu také můžeme pro s.v. x a y taková, že $f(y) \neq 0$, zapisovat ve tvaru $f(x|y) = \frac{f(x, y)}{f(y)}$.

Nyní můžeme zavést podmíněnou střední hodnotu.

Definice 2 (Podmíněná střední hodnota). *Podmíněnou střední hodnotu veličiny X při daném Y definujeme předpisem*

$$E[X|Y] = \int_{\mathbb{R}} x f(x|y) dx$$

V následující větě shrneme některé vlastnosti podmíněné střední hodnoty, které budeme později využívat při výpočtech.

Věta 1 (Základní vlastnosti podmíněné střední hodnoty). *Budte náhodné veličiny $X, Y, Z \in L_1(\Omega, \mathcal{A}, P)$. Pak:*

i) $\forall a, b, c \in \mathbb{R} : E[aX + bY + c|Z] = a E[X|Z] + b E[Y|Z] + c$ s.j.

ii) $X \leq Y$ s.j. $\implies E[X|Z] \leq E[Y|Z]$ s.j.

iii) $E[E[X|Y]] = E[X]$

iv) Je-li $\sigma(X) \subset \sigma(Z)$ pak $E[X|Z] = X$ s.j.

v) Bud $Y \in L(\Omega, \mathcal{F})$ a $XY \in L_1(\Omega, \mathcal{A}, P)$, pak $E[XY|\mathcal{F}] = Y E[X|\mathcal{F}]$

vi) Jsou-li X a Y nezávislé veličiny, pak $E[X|Y] = E[X]$ s.j.

Důkaz lze najít např. v (Lachout, 2004).

Dále nás bude zajímat podmíněný rozptyl. Podobně jako je nepodmíněný rozptyl definován pomocí střední hodnoty, využijeme k definici podmíněného rozptylu podmíněnou střední hodnotu.

Definice 3 (Podmíněný rozptyl). *Nechť $S(Y, Z)$ je náhodná veličina a $E[S^2] < \infty$. Pak podmíněný rozptyl definujeme předpisem $\text{var}(S|Y) = E[(S - E(S|Y))^2|Y]$.*

Nyní uvedeme důležitou větu o rozkladu nepodmíněného rozptylu.

Věta 2. *Pokud $E[S^2] < \infty$, pak $\text{var}(S) = \text{var}[E[S|Y]] + E[\text{var}(S|Y)]$.*

Důkaz.

V důkazu využijeme vlastnosti podmíněné střední hodnoty z věty 1. Nejprve rozepíšeme podmíněný rozptyl do vhodnějšího tvaru.

$$\begin{aligned} \text{var}(S|Y) &= E[(S - E(S|Y))^2|Y] = E[(S - E[S] + E[S] - E(S|Y))^2|Y] = \\ &= E\{[(S - E[S])^2 + 2(S - E[S])(E[S] - E(S|Y)) + (E[S] - E(S|Y))^2|Y\} = \\ &= E[(S - E[S])^2|Y] + 2[E[S] - E(S|Y)]E[(S - E[S])|Y] + (E[S] - E(S|Y))^2 = \\ &= E[(S - E[S])^2|Y] + 2(E[S] - E(S|Y))(E(S|Y) - E[S]) + (E[S] - E(S|Y))^2 = \\ &= E[(S - E[S])^2|Y] - (E[S] - E(S|Y))^2 \end{aligned}$$

Nyní aplikujeme na výsledek operátor střední hodnoty.

$$\begin{aligned} E[\text{var}(S|Y)] &= E[E[(S - E[S])^2|Y] - (E[S] - E(S|Y))^2] = \\ &= E[E[(S - E[S])^2|Y]] - E(E[S] - E(S|Y))^2 = \\ &= E(S - E[S])^2 - E(E(S|Y) - E[S])^2 = \\ &= \text{var}(S) - E(E(S|Y) - E[E(S|Y)])^2 = \text{var}(S) - \text{var}(S|Y) \end{aligned}$$

Dohromady tedy získáváme vztah $E[\text{var}(S|Y)] = \text{var}(S) - \text{var}(S|Y)$, ze kterého plyne dokazovaný vzorec. □

1.2 χ^2 rozdělení a kvadratické formy

Nejprve připomeňme definici χ^2 rozdělení.

Definice 4 (χ^2 rozdělení). *Nechť X_1, \dots, X_n jsou nezávislé náhodné veličiny s normovaným normálním rozdělením $N(0, 1)$. Pak náhodná veličina $Y = \sum_{i=1}^n X_i^2$ má χ^2 rozdělení s n stupni volnosti, píšeme $Y \sim \chi_n^2$.*

V následujících kapitolách budeme často potřebovat zjistit rozdělení náhodných veličin ve tvaru kvadratických forem. Shrňme si proto několik užitečných vlastností, které pro ně platí, v následující větě.

Věta 3. *i) Nechť Y, Z jsou nezávislé náhodné veličiny takové, že $Y \sim \chi_{n_1}^2$, $Z \sim \chi_{n_2}^2$. Pak náhodná veličina $Y + Z \sim \chi_{n_1+n_2}^2$.
ii) Nechť $\mathbf{X} \sim N_n(\boldsymbol{\mu}, \Sigma)$, kde Σ je regulární varianční matice. Pak $\mathbf{Y} = (\mathbf{X} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{X} - \boldsymbol{\mu}) \sim \chi_n^2$.
iii) Nechť $\mathbf{X} \sim N_n(\mathbf{0}, \Sigma)$ a \mathbb{A} buď čtvercová matice taková, že $\mathbb{A}\Sigma$ je nenulová idempotentní matice. Pak $\mathbf{Y} = \mathbf{X}^T \mathbb{A} \mathbf{X} \sim \chi_{\text{tr} \mathbb{A} \Sigma}^2$.*

Důkaz.

i) Veličiny Y a Z lze přepsat jako $Y = \sum_{i=1}^{n_1} X_i^2$, $Z = \sum_{i=n_1+1}^{n_2} X_i^2$, kde X_i jsou vzájemně nezávislé veličiny s $N(0,1)$ rozdělením. Pak z definice vyplývá, že

$$Y + Z = \sum_{i=1}^{n_1+n_2} X_i^2 \sim \chi_{n_1+n_2}^2.$$

ii) Použijeme vlastnost, že \mathbf{X} má stejné rozdělení, jako $\boldsymbol{\mu} + \mathbb{B}\mathbf{X}_0$, kde $\mathbb{B}\mathbb{B}^T = \Sigma$ a $\mathbf{X}_0 \sim N_n(\mathbf{0}, \mathbb{I}_n)$. Z toho plyne, že \mathbf{Y} má stejné rozdělení jako

$$(\mathbb{B}\mathbf{X}_0)^T \Sigma^{-1} (\mathbb{B}\mathbf{X}_0) = \mathbf{X}_0^T \mathbb{B}^T \mathbb{B} \mathbb{B}^T \mathbb{B}^{-1} \mathbf{X}_0 = \mathbf{X}_0^T \mathbf{X}_0.$$

Pak $\mathbf{X}_0^T \mathbf{X}_0$ lze přepsat jako součet kvadrátů jednotlivých složek náhodného vektoru \mathbf{X}_0 , což odpovídá definici χ_n^2 rozdělení.

iii) Matice $\mathbb{A}\Sigma$ je nemulová, z toho plyne $\text{rank}(\mathbb{A}) \geq 1$, tedy existuje skeletní rozklad matice $\mathbb{A} = \mathbb{B}\mathbb{B}^T$. Pak $\mathbf{X}^T \mathbb{A} \mathbf{X} = \mathbf{X}^T \mathbb{B} \mathbb{B}^T \mathbf{X} = (\mathbb{B}^T \mathbf{X})^T (\mathbb{B}^T \mathbf{X}) = (\mathbb{B}^T \mathbf{X})^T \mathbb{I} (\mathbb{B}^T \mathbf{X})$, přičemž $\mathbb{B}^T \mathbf{X} \sim N(\mathbf{0}, \mathbb{B}^T \Sigma \mathbb{B})$. Nyní chceme použít tvrzení ii), k tomu potřebujeme ukázat, že $\mathbb{B}^T \Sigma \mathbb{B}$ je idempotentní.

Vyjděme z předpokladu, že $\mathbb{A}\Sigma$ je idempotentní a přepišme ho pomocí skeletního rozkladu \mathbb{A} , tj. $(\mathbb{B}\mathbb{B}^T \Sigma)(\mathbb{B}\mathbb{B}^T \Sigma) = (\mathbb{B}\mathbb{B}^T \Sigma)$. Nyní vynásobíme rovnici zleva maticí \mathbb{B}^{-1} , z vlastností skeletního rozkladu plyne, že $\mathbb{B}^{-1}\mathbb{B} = \mathbb{I}$. Pak $(\mathbb{B}^{-1}\mathbb{B}\mathbb{B}^T \Sigma)(\mathbb{B}\mathbb{B}^T \Sigma) = (\mathbb{B}^{-1}\mathbb{B}\mathbb{B}^T \Sigma)$, neboli $(\mathbb{B}\mathbb{B}^T \Sigma)(\mathbb{B}^T \Sigma \mathbb{B}) = (\mathbb{B}^T \Sigma \mathbb{B})$. Dokázali jsme, že je matice $\mathbb{B}\Sigma\mathbb{B}$ idempotentní, tedy její pseudoinverzí je jednotková matice \mathbb{I} . Tvrzení pak vyplývá z části ii) a faktu, že $\text{rank}(\mathbb{B}\Sigma\mathbb{B}) = \text{tr}(\mathbb{B}\Sigma\mathbb{B}) = \text{tr}(\mathbb{B}\mathbb{B}^T \Sigma) = \text{tr}(\mathbb{A}\Sigma)$. □

Na závěr uvedeme užitečné tvrzení o nezávislosti.

Věta 4. *Nechť $\mathbf{X} \sim N_n(\boldsymbol{\mu}, \Sigma)$ a \mathbb{A} je pozitivně definitní matice typu $n \times n$.*

i) Nechť \mathbb{B} je libovolná obdelníková matice splňující $\mathbb{B}\Sigma\mathbb{A} = \mathbf{0}$. Pak náhodná veličina $\mathbf{X}^T \mathbb{A} \mathbf{X}$ a náhodný vektor $\mathbb{B}\mathbf{X}$ jsou vzájemně nezávislé.

ii) Nechť \mathbb{B} je čtvercová pozitivně definitní matice splňující $\mathbb{B}\Sigma\mathbb{A} = \mathbf{0}$. Pak náhodné veličiny $\mathbf{X}^T \mathbb{A} \mathbf{X}$ a $\mathbf{X}^T \mathbb{B} \mathbf{X}$ jsou vzájemně nezávislé.

Důkaz.

i) Využijeme skeletní rozklad $\mathbb{A} = \mathbb{L}\mathbb{L}^T$. Předpokládáme $\mathbb{B}\Sigma\mathbb{L}\mathbb{L}^T = \mathbf{0} \implies \mathbb{B}\Sigma\mathbb{L}\mathbb{L}^T(\mathbb{L}^T)^{-1} = \mathbf{0}$. Z toho plyne, že $\mathbf{0} = \mathbb{B}\Sigma\mathbb{L} = \text{cov}(\mathbb{B}\mathbf{X}, \mathbb{L}^T \mathbf{X})$. Normalita vektoru \mathbf{X} implikuje, že i vektory $\mathbb{B}\mathbf{X}$ a $\mathbb{L}^T \mathbf{X}$ mají normální rozdělení a protože jsou nekorelované, jsou i nezávislé. Z toho plyne i nezávislost vektorů $\mathbb{B}\mathbf{X}$ a $(\mathbb{L}^T \mathbf{X})^T \mathbb{L}^T \mathbf{X} = \mathbf{X}^T (\mathbb{L}^T)^T \mathbb{L}^T \mathbf{X} = \mathbf{X}^T \mathbb{A} \mathbf{X}$.

ii) Opět vyjdeme ze skeletního rozkladu matic $\mathbb{A} = \mathbb{L}\mathbb{L}^T$ a $\mathbb{B} = \mathbb{P}\mathbb{P}^T$. Analogicky dosadíme do předpokladu skeletní rozkladu a získáme vztah $\mathbb{P}^T \Sigma \mathbb{L} = \mathbf{0} = \text{cov}(\mathbb{P}^T \mathbf{X}, \mathbb{L}^T \mathbf{X})$. Z normality vektorů $\mathbb{P}^T \mathbf{X}$, $\mathbb{L}^T \mathbf{X}$ a jejich nekorelovanosti plyne, že jsou nezávislé. Pak jsou nezávislé i vektory $(\mathbb{P}^T \mathbf{X})^T \mathbb{P}^T \mathbf{X} = \mathbf{X}^T \mathbb{B} \mathbf{X}$ a $(\mathbb{L}^T \mathbf{X})^T \mathbb{L}^T \mathbf{X}$. □

2. Analýza rozptylu s pevnými efekty

2.1 Jednoduché třídění

Základní situací v analýze rozptylu je takzvané jednoduché třídění. V tomto případě uvažujeme situaci, kdy máme N vzájemně nezávislých náhodných výběrů $Y_{i,1}, \dots, Y_{i,n_i}, i = 1, \dots, N$, kde $Y_{i,j}, j = 1, \dots, n_i$ pochází z rozdělení s distribuční funkcí F_i s konečnou střední hodnotou μ_i . Zajímá nás, zda je tato střední hodnota stejná pro všechny výběry. V případě zamítnutí nulové hypotézy také obvykle chceme vědět, které dvojice středních hodnot se od sebe statisticky signifikantně liší.

Příkladem může být situace, kdy máme N odrůd ovocných stromů (např. jabloní), u kterých zjišťujeme celkovou hmotnost úrody. Chceme vědět, jestli všechny odrůdy mají shodnou střední hodnotu hmotnosti úrody (neboli zda celková hmotnost plodů nezávisí na skupině) a pokud ne, u kterých skupin nastávají významné odlišnosti.

Model 1. Pro zbytek kapitoly předpokládejme, že jednotlivé náhodné výběry jsou vzájemně nezávislé a pochází z normálních rozdělení s neznámými středními hodnotami $\mu_i, i = 1, \dots, N$ a shodným rozptylem σ^2 .

Tedy: $Y_{i,j} \stackrel{i.i.d.}{\sim} N(\mu_i, \sigma^2), j = 1, \dots, n_i$, přičemž $\mathbf{Y}_i := (Y_{i,1}, \dots, Y_{i,n_i})^T$ jsou vzájemně nezávislé pro $i = 1, \dots, N$.

Budeme testovat hypotézu $H_0: \mu_1 = \mu_2 = \dots = \mu_N$ proti alternativní hypotéze H_1 : Existuje dvojice středních hodnot, které se nerovnají.

Výše zavedený model 1 je ekvivalentní následujícímu modelu:

$$Y_{i,j} = \mu + \beta_i + \varepsilon_{i,j}, \text{ kde } \varepsilon_{i,j} \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

Toto vyjádření získáme, položíme-li $\mu_i = \mu + \beta_i$, přičemž β_i zde reprezentuje odchylku mezi μ a střední hodnotou μ_i v (i,j) - tém výběru. Veličina $\varepsilon_{i,j}$ pak představuje variabilitu dat danou přírodními podmínkami. Tento model se označuje jako přeparametrizovaný, protože máme zavedeno o parametr více než je zapotřebí.

Nulovou hypotézu pak zapíšeme v ekvivalentním tvaru:

$$H_0 : \beta_1 = \dots = \beta_N$$

. Za platnosti nulové hypotézy bude náš model tvaru: $Y_{ij} = \mu + \beta + \varepsilon_{ij}$, kde $\beta = \beta_1 = \dots = \beta_N$. Místo porovnávání střední hodnoty μ_i se tedy zaměříme na porovnávání parametrů β_i .

Abychom sestavili testovou statistiku, zavedeme napřed několik pojmů:

Definice 5 (Součty čtverců). Označme $M = \sum_{i=1}^N n_i$.

Zavedeme značení pro součet a výběrový průměr pozorování ve skupině:

$$Y_{i,+} = \sum_{j=1}^{n_i} Y_{i,j} \text{ a } \bar{Y}_{i,+} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{i,j},$$

a dále značení pro celkový součet a celkový průměr pozorování:

$$Y_{++} = \sum_{i=1}^N \sum_{j=1}^{n_i} Y_{i,j} \quad a \quad \bar{Y}_{++} = \frac{1}{M} \sum_{i=1}^N \sum_{j=1}^{n_i} Y_{i,j}.$$

Položme: $\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,n_i})^T$ vektor pozorování ve skupinách,

$\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_N^T)^T$ vektor všech pozorování a

$\bar{\mathbf{Y}} = (\bar{Y}_{1,+}, \dots, \bar{Y}_{N,+})^T$ vektor skupinových průměrů.

Nyní položme:

$SS_c = \sum_{i=1}^N \sum_{j=1}^{n_i} (Y_{i,j} - \bar{Y}_{++})^2$ celkový součet čtverců,

$SS_a = \sum_{i=1}^N n_i (\bar{Y}_{i,+} - \bar{Y}_{++})^2$ součet čtverců skupin,

$SS_e = \sum_{i=1}^N \sum_{j=1}^{n_i} (Y_{i,j} - \bar{Y}_{i,+})^2$ residuální součet čtverců.

Takto zavedený součet čtverců skupin reprezentuje rozptyl mezi jednotlivými skupinami, naopak residuální součet představuje rozptyl uvnitř skupin. Součty čtverců lze rozepsat jako kvadratické formy, které mají za platnosti uvažované nulové hypotézy a výše zavedeného modelu rozdělení chí-kvadrát a jsou nestrannými odhady rozptylu σ^2 . Podrobněji jsou tyto vztahy popsány v následující větě.

Věta 5. i) Za platnosti modelu 1 máme: $\frac{SS_e}{\sigma^2} \sim \chi_{M-N}^2$ a $E \frac{SS_e}{M-N} = \sigma^2$.

ii) Za platnosti modelu 1 a nulové hypotézy splňuje celkový součet čtverců $\frac{SS_c}{\sigma^2} \sim \chi_{M-1}^2$ a $E \frac{SS_c}{M-1} = \sigma^2$, pro součet čtverců skupin máme: $\frac{SS_a}{\sigma^2} \sim \chi_{N-1}^2$ a $E \frac{SS_a}{N-1} = \sigma^2$.

iii) SS_a a SS_e jsou nezávislé.

Důkaz.

V důkazu budeme využívat tvrzení 3 o rozdělení kvadratických forem.

i) Platí $\frac{SS_e}{\sigma^2} = \sum_{i=1}^N \frac{1}{\sigma^2} \mathbf{Y}_i^T \mathbb{A}_i \mathbf{Y}_i$, kde \mathbb{A}_i je čtvercová matice, která má na diagonále prvky $1 - \frac{1}{n_i}$ a prvky $-\frac{1}{n_i}$ mimo diagonálu. Víme, že $\frac{1}{\sigma^2} \mathbf{Y}_i^T \mathbb{A}_i \mathbf{Y}_i \sim \chi_{n_i-1}^2$, s použitím nezávislosti \mathbf{Y}_i dostáváme $\frac{SS_e}{\sigma^2} \sim \chi_{\sum_{i=1}^N (n_i-1)}^2 = \chi_{M-N}^2$.

Střední hodnota veličiny s χ^2 rozdělení je rovna počtu stupňů volnosti, tzn. $E \frac{SS_e}{\sigma^2} = M - N$, z čehož plyne, že po vydělení $M - N$ je residuální součet čtverců nestranným odhadem rozptylu σ^2 .

ii) Za platnosti nulové hypotézy lze všechna pozorování $Y_{i,j}$ uspořádat do vektoru $\mathbf{Y} \sim N_M(\boldsymbol{\mu} \mathbf{1}_M, \sigma^2 \mathbb{I}_M)$, pak jde o výběr o rozsahu M . $\frac{SS_c}{M-1}$ je jeho výběrový rozptyl, z vlastností výběrového rozptylu pak plyne tvrzení pro SS_c .

Dále definujeme vektor $\bar{\mathbf{Y}} = \text{diag}^{-1}(n_1, \dots, n_N)(Y_{1,+}, \dots, Y_{N,+})^T$. Pak

$$\frac{SS_a}{\sigma^2} = \frac{1}{\sigma^2} (\bar{\mathbf{Y}} - \boldsymbol{\mu} \mathbf{1}_N)^T \mathbb{C} (\bar{\mathbf{Y}} - \boldsymbol{\mu} \mathbf{1}_N),$$

přičemž $\mathbb{C} = \text{diag}(\mathbf{n}) - \frac{1}{M} \mathbf{n} \mathbf{n}^T$. Vektor $(\bar{\mathbf{Y}} - \boldsymbol{\mu} \mathbf{1}_N) \sim N_N(\mathbf{0}, \sigma^2 \text{diag}^{-1}(\mathbf{n}))$, kde $\mathbf{n} = (n_1, \dots, n_N)$, takže pokud dokážeme, že je matice $\mathbb{C} \text{diag}^{-1}(\mathbf{n})$ idempotentní, bude platit $\frac{SS_a}{\sigma^2} \sim \chi_{\text{tr} \mathbb{C} \text{diag}^{-1}(\mathbf{n})}^2$. Idempotenci snadno ověříme s použitím vztahu

$$\mathbb{C} \text{diag}^{-1}(\mathbf{n}) = \mathbb{I}_N - \frac{1}{M} \mathbf{n} \mathbf{n}^T \text{diag}^{-1}(\mathbf{n}) = \mathbb{I}_N - \frac{1}{M} \mathbf{n} \mathbf{1}_N.$$

Tento přepis je také vhodný pro výpočet stopy matice,

$$\text{tr}(\mathbb{I}_N - \frac{1}{M} \mathbf{n} \mathbf{n}^T \text{diag}^{-1}(\mathbf{n})) = \text{tr} \mathbb{I}_N - \frac{1}{M} \text{tr}(\mathbf{n} \mathbf{1}_N) = N - 1.$$

Dohromady dostáváme $\frac{SS_a}{\sigma^2} \sim \chi_{N-1}^2$. Pro střední hodnotu opět platí, že $E \frac{SS_a}{\sigma^2} = N - 1$, z čehož plyne $E \frac{SS_a}{N-1} = \sigma^2$.

iii) V předchozích částech důkazu jsme ukázali, že SS_a i SS_e lze přepsat jako kvadratické formy, konkrétně

$$SS_e = \mathbf{Y}^T \mathbb{A} \mathbf{Y},$$

$$SS_a = \mathbf{Y}^T \mathbb{B} \mathbf{Y},$$

kde $B = \text{diag}(\frac{1}{n_1} \mathbf{1}_{n_1}, \dots, \frac{1}{n_N} \mathbf{1}_{n_N}) \mathbb{C} \text{diag}(\frac{1}{n_1} \mathbf{1}_{n_1}^T, \dots, \frac{1}{n_N} \mathbf{1}_{n_N}^T)$. K ověření nezávislosti nám tedy podle věty 4 stačí ukázat, že $\mathbb{B} \text{Var} \mathbf{Y} \mathbb{A} = \sigma^2 \mathbb{B} \mathbb{A} = \mathbf{0}$. Nemusíme provádět celý výpočet, stačí si všimnout, že $\text{diag}(\frac{1}{n_1} \mathbf{1}_{n_1}^T, \dots, \frac{1}{n_N} \mathbf{1}_{n_N}^T) \mathbb{A} = \mathbf{0}$, pak je i celý součin nulový a nezávislost veličin SS_a a SS_e je dokázána. \square

Na základě vztahů z předchozí věty navrhneme testovou statistiku jako podíl nestranných odhadů rozptylu:

$$F_A = \frac{\frac{SS_a}{N-1}}{\frac{SS_e}{M-N}}$$

Z věty 5 dále plyne, že za platnosti nulové hypotézy má F_A Fisherovo rozdělení s $N - 1$, $M - N$ stupni volnosti.

Navrhli jsme testovou statistiku a nyní nás zajímá, jak volit kritický obor. Pokud platí nulová hypotéza o rovnosti středních hodnot skupin, pak skupinový průměr $\bar{Y}_{i,+}$ představuje odhad středních hodnot skupin a celkový průměr \bar{Y}_{++} představuje odhad celkové střední hodnoty, které by měly být stejné. Očekáváme tedy, že rozdíl $(\bar{Y}_{i,+} - \bar{Y}_{++})$ bude nabývat velmi malých hodnot, a proto i skupinový součet čtverců SS_a bude nabývat menší hodnotu, než residuální součet čtverců SS_e . Platnost nulové hypotézy by se měla projevit v malých hodnotách testové statistiky F_A , velké hodnoty naznačují platnost alternativní hypotézy. To nás vede k myšlence zvolit kritický obor tak, aby test zamítal nulovou hypotézu pro vysoké hodnoty F_A . Za účelem získání testu s hladinou α pak H_0 zamítáme $\iff F_A \geq F_{N-1, M-N}(1 - \alpha)$, kde $F_{N-1, M-N}(1 - \alpha)$ je $(1 - \alpha)$ kvantil Fisherova rozdělení s $N - 1$, $M - N$ stupni volnosti.

Pro zápis výsledků jednoduchého třídění se často používá následující tabulka:

Zdroj měnlivosti	Součet čtverců	Počet stupňů volnosti	Podíl	Testová statistika
Skupiny	SS_a	$N - 1$	$\frac{SS_a}{N-1}$	$\frac{SS_a}{N-1} / \frac{SS_e}{M-N}$
Residuální	SS_e	$M - N$	$\frac{SS_e}{M-N}$	
Celkový	SS_c	$M - 1$		

Tabulka 2.1: Analýza rozptylu jednoduchého třídění

Alternativně se model jednoduchého třídění zavádí s využitím teorie lineárních modelů. Tento přístup je popsán např. v knize (Anděl, 2007b, kapitola 9.2).

3. Analýza rozptylu s náhodnými efekty

3.1 Jednoduché třídění - zavedení modelu

V minulé kapitole jsme zkoumali případ, kdy máme N vzájemně nezávislých náhodných výběrů, o kterých za platnosti nulové hypotézy předpokládáme, že pochází z rozdělení se stejnou distribuční funkcí. Zavedli jsme model, ve kterém jsme tyto výběry popsali pomocí normálního rozdělení s neznámou střední hodnotou a pro všechny výběry shodným neznámým rozptylem σ^2 . Jednotlivá pozorování jsme popsali pomocí působení několika veličin, které byly všechny až na člen $\varepsilon_{i,j}$ pevnými neznámými konstanty. Tento případ se nazýval model s pevnými efekty.

Nyní zavedeme mírně odlišnou situaci. Uvažujme, že skupin, ze kterých bereme naše pozorování (různé odrůdy jabloní v příkladu z předchozí kapitoly), není jen pevně dané N , ale může jich být až spočetně mnoho. Nepůjde o pevně daný počet skupin, ale o náhodný výběr z větší populace. Tato úvaha nás vede k zavedení jiného typu modelu (bude se jednat o speciální případ lineárního smíšeného modelu). Nový model vysvětlíme na příkladu.

Mějme k laboratorních myší, každá z nich má $n_i, i = 1, \dots, k$ mláďat. Budeme měřit rozdíly ve hmotnosti mláďat po dvou týdnech věku a zajímá nás, jestli mláďata mají stejnou střední hodnotu hmotnosti nezávisle na matce.

V tomto případě předpokládáme, že našich k myší bylo náhodně vybráno z populace myší. Hmotnost mláďate pro danou i -tou myš má střední hodnotu $\mu + B_i$. Předpokládáme, že, náhodný výběr myši se projeví jako jistá variabilita B_i . Dále předpokládejme, že $E B_i = 0, \text{var } B_i = d^2$.

Zároveň hmotnosti všech mláďat dané matky mají také určitou variabilitu. Stejně jako v předchozím modelu jako původce této variability označíme veličinu $\varepsilon_{i,j}$, která reprezentuje odchylku hmotnosti konkrétního mláďate od střední hodnoty hmotnosti pro mláďata dané myši a budeme předpokládat, že $E \varepsilon_{i,j} = 0, \text{var } \varepsilon_{i,j} = \sigma^2$. Hmotnost j -tého mláďate i -té myši pak popíšeme vztahem $Y_{i,j} = \mu + B_i + \varepsilon_{i,j}$.

Narozdíl od modelu s pevnými efekty, kdy B_i , resp. β_i , byla neznámou konstantou, jde nyní o náhodnou proměnnou. Tato situace se proto označuje jako model s pevnými efekty. Obecně se podle knihy (Khuri, 2010, kapitola 1.5) jako model s náhodnými efekty označuje model, kde jsou všechny členy kromě prvního náhodné.

Výše zmíněné požadavky společně s předpoklady o rozdělení zformulujeme jako následující model.

Model 2. Mějme náhodný vektor $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_k^T)^T$, kde pro všechna $i = 1, \dots, k$ jsou $\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,n_i})^T$ vzájemně nezávislé skupiny pozorování. Řekneme, že \mathbf{Y} se řídí modelem s náhodnými efekty, pokud lze psát $Y_{i,j} = \mu + B_i + \varepsilon_{i,j}$, kde $\varepsilon_{i,j} \stackrel{i.i.d.}{\sim} N(0, \sigma^2), 0 < \sigma^2 < \infty$ a $B_i \stackrel{i.i.d.}{\sim} N(0, d^2), 0 < d^2 < \infty$, přičemž μ, σ^2 a d^2 jsou neznámé parametry.

Pro zbytek textu předpokládejme, že B_i a $\boldsymbol{\varepsilon}_i = (\varepsilon_{i,1}, \dots, \varepsilon_{i,n_i})^T$ jsou nezávislé.

3.2 Odvození vlastností

Nejprve odvodíme vlastnosti pozorování při použití našeho modelu, budou nás zajímat hlavně střední hodnota a rozptyl (i,j) - tého pozorování. Kromě toho potřebujeme znát tyto hodnoty i pro (i,j) - té pozorování podmíněné veličinou B_i .

V průběhu odvozování vycházíme z modelu 2, budeme využívat předpokládané vlastnosti složek $Y_{i,j}$ (známé rozdělení a nezávislost B_i a $\varepsilon_{i,j}$). Dále budeme používat vlastnosti podmíněné střední hodnoty a podmíněného rozptylu, které byly uvedeny v kapitole 1.

Nejprve odvodíme střední hodnotu pro (i,j) - té pozorování. Využijeme linearity střední hodnoty a toho, že už známe střední hodnotu B_i a $\varepsilon_{i,j}$.

$$\mathbf{E}(Y_{i,j}) = \mathbf{E}(\mu + B_i + \varepsilon_{i,j}) = \mu + \mathbf{E}(B_i) + \mathbf{E}(\varepsilon_{i,j}) = \mu$$

Dále odvodíme střední hodnotu pro (i,j) - té pozorování podmíněné veličinou B_i . Využíváme znalostí o podmíněných veličinách, zejména vztahů o linearitě podmíněné střední hodnoty. Důležitý je také předpoklad, že $\varepsilon_{i,j}$ a B_i jsou na sobě nezávislé veličiny.

$$\begin{aligned} \mathbf{E}(Y_{i,j}|B_i) &= \mathbf{E}(\mu + B_i + \varepsilon_{i,j}|B_i) = \mathbf{E}(\mu|B_i) + \mathbf{E}(B_i|B_i) + \mathbf{E}(\varepsilon_{i,j}|B_i) = \\ &= \mu + B_i + \mathbf{E}(\varepsilon_{i,j}) = \mu + B_i \end{aligned}$$

Zajímá nás také podmíněný rozptyl (i,j) - tého pozorování. Při odvození použijeme invariantnost rozptylu vůči konstantě (při podmínění veličinou B_i vnímáme B_i jako zafixovanou na určité hodnotě) a dále předpoklad o nezávislosti $\varepsilon_{i,j}$ a B_i .

$$\text{var}(Y_{i,j}|B_i) = \text{var}(\mu + B_i + \varepsilon_{i,j}|B_i) = \text{var}(B_i + \varepsilon_{i,j}|B_i) = \text{var}(\varepsilon_{i,j}|B_i) = \text{var}(\varepsilon_{i,j}) = \sigma^2$$

Když jsme odvodili podmíněný rozptyl, můžeme použít vzorec 2 o rozkladu rozptylu a zjistit hodnotu nepodmíněného rozptylu.

$$\begin{aligned} \text{var}(Y_{i,j}) &= \mathbf{E}(\text{var}(Y_{i,j}|B_i)) + \text{var}(\mathbf{E}(Y_{i,j}|B_i)) = \mathbf{E}(\sigma^2) + \text{var}(\mu + B_i) = \\ &= \sigma^2 + \text{var}(B_i) = \sigma^2 + d^2 \end{aligned}$$

Nyní jsme odvodili vlastnosti pro jednotlivá pozorování, pro srovnání výsledků s předchozím modelem uvádíme následující tabulku:

Model 1	Model 2
$Y_{i,j} = \mu + \beta_i + \varepsilon_{i,j}, \beta_i \in \mathbb{R}$	$Y_{i,j} = \mu + B_i + \varepsilon_{i,j}, B_i \sim N(0, d^2)$
$\mathbf{E}(Y_{i,j}) = \mu + \beta_i$	$\mathbf{E}(Y_{i,j}) = \mu$
$\text{var}(Y_{i,j}) = \sigma^2$	$\text{var}(Y_{i,j}) = \sigma^2 + d^2$

Tabulka 3.1: Srovnání vlastností $Y_{i,j}$

Obdobně jako v přechozí kapitole položíme $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_k^T)^T$ vektor všech pozorování. Chtěli bychom s použitím předchozích výsledků sestavit rozptylovou matici $\text{Var}(\mathbf{Y})$.

Nejprve zjistíme rozptylovou matici pro vektor \mathbf{Y}_i , který reprezentuje i -tý náhodný výběr o rozsahu n_i . Struktura rozptylové matice je známá.

$$\text{Var}(\mathbf{Y}_i) = \begin{pmatrix} \text{var}(Y_{i,1}), & \text{cov}(Y_{i,1}, Y_{i,2}), & \dots, & \text{cov}(Y_{i,1}, Y_{i,n_i}) \\ \text{cov}(Y_{i,2}, Y_{i,1}), & \text{var}(Y_{i,2}), & \dots, & \text{cov}(Y_{i,2}, Y_{i,n_i}) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(Y_{i,n_i}, Y_{i,1}), & \text{cov}(Y_{i,n_i}, Y_{i,2}), & \dots, & \text{var}(Y_{i,n_i}) \end{pmatrix}$$

Z předchozích výpočtů už víme, že na diagonále budou hodnoty $\sigma^2 + d^2$ (tuto hodnotu jsme odvodili pro libovolné (i, j) - té pozorování). Zbývá už zjistit jen příslušné hodnoty kovariance. Pro následující výpočet nás zajímají pouze hodnoty mimo diagonálu, tedy zkoumáme jen situace kdy $j \neq l$.

$$\begin{aligned} \text{cov}(Y_{i,j}, Y_{i,l}) &= \mathbf{E}(Y_{i,j}, Y_{i,l}) - \mathbf{E}(Y_{i,j}) \mathbf{E}(Y_{i,l}) = \mathbf{E}((\mu + B_i + \varepsilon_{i,j})(\mu + B_i + \varepsilon_{i,l})) - \mu^2 = \\ &= \mathbf{E}(\mu^2) + \mathbf{E}(2\mu B_i) + \mathbf{E}(\mu \varepsilon_{i,j}) + \mathbf{E}(\mu \varepsilon_{i,l}) + \mathbf{E}(B_i \varepsilon_{i,l}) + \mathbf{E}(B_i \varepsilon_{i,j}) + \mathbf{E}(\varepsilon_{i,j} \varepsilon_{i,l}) - \mu^2 = \\ &= \mu^2 + \text{var}(B_i) - \mu^2 = d^2 \end{aligned}$$

V průběhu odvození jsme opakovaně využívali nezávislost veličin B_i a $\varepsilon_{i,j}$ a z toho plynoucího vztahu $\mathbf{E}(B_i \varepsilon_{i,j}) = 0$. Díky zavedenému modelu jsme také mohli využít znalost momentů (nulová střední hodnota a známý rozptyl) u obou veličin.

Dohromady dostáváme, že

$$\text{Var}(\mathbf{Y}_i) = \begin{pmatrix} \sigma^2 + d^2, & d^2, & \dots, & d^2 \\ d^2, & \sigma^2 + d^2, & \dots, & d^2 \\ \vdots & \vdots & \ddots & \vdots \\ d^2, & d^2, & \dots, & \sigma^2 + d^2 \end{pmatrix}.$$

Rozptylovou matici můžeme přepsat i pomocí vztahu $\text{Var}(\mathbf{Y}_i) = \sigma^2 \mathbb{I}_{n_i} + d^2 \mathbb{J}_{n_i}$, kde \mathbb{I}_{n_i} je jednotková matice o délce hrany n_i a \mathbb{J}_{n_i} je čtvercová matice o délce hrany n_i , tvořená samými jedničkami. Tento úspornější zápis využijeme pro zápis matice celkového rozptylu.

Nyní můžeme přistoupit k odvození matice rozptylu všech pozorování $\text{Var}(\mathbf{Y})$. Známe už rozptyl a kovariance pro jeden náhodný výběr, zbývá nám zjistit mezivýběrové kovariance. V modelu 2 jsme definovali, že všechny náhodné výběry jsou vzájemně nezávislé, z toho vyplývá důležitý vztah: $\text{cov}(Y_{i,j}, Y_{k,l}) = 0$, kdykoliv $i \neq k$.

Vidíme, že matice $\text{Var}(\mathbf{Y})$ bude blokově diagonální, diagonála obsahuje rozptylové matice jednotlivých výběrů $\text{Var}(\mathbf{Y}_i)$, $i = 1, 2, \dots, k$, což jsou čtvercové matice s délkou hrany n_i , $i = 1, 2, \dots, k$ a strukturou popsanou výše. Zbytek matice je tvořen nulami. Můžeme si rovněž povšimnout, že hlavní diagonála je tvořena hodnotami $\sigma^2 + d^2$.

Matice vypadá následovně (indexy n_1, \dots, n_k značí rozměry bloků na diagonále, symboly 0 reprezentují nulové bloky vhodných rozměrů):

$$\text{Var}(\mathbf{Y}) = \begin{pmatrix} \sigma^2 \mathbb{I}_{n_1} + d^2 \mathbb{J}_{n_1} & 0 & \dots & 0 \\ 0 & \sigma^2 \mathbb{I}_{n_2} + d^2 \mathbb{J}_{n_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \mathbb{I}_{n_k} + d^2 \mathbb{J}_{n_k} \end{pmatrix}.$$

3.3 Součty čtverců

Obdobně jako v kapitole 2 nejprve zavedeme součty čtverců, nyní ovšem budeme předpokládat tzv. vyvážené třídění, tedy že všechny skupiny obsahují stejný počet (n) vzorků.

Definice 6 (Součty čtverců). *Položme:*

$$SS_c = \sum_{i=1}^k \sum_{j=1}^n (Y_{i,j} - \bar{Y}_{++})^2 \text{ celkový součet čtverců,}$$

$$SS_a = \sum_{i=1}^k n(\bar{Y}_{i,+} - \bar{Y}_{++})^2 \text{ součet čtverců skupin,}$$

$$SS_e = \sum_{i=1}^k \sum_{j=1}^n (Y_{i,j} - \bar{Y}_{i,+})^2 \text{ residuální součet čtverců,}$$

kde $\bar{Y}_{i,+}$ a \bar{Y}_{++} značí skupinový a celkový průměr jako v definici 5.

Opět platí, že skupinový součet čtverců reprezentuje variabilitu mezi jednotlivými skupinami a residuální součet čtverců celkovou variabilitu uvnitř skupin. Chceme zjistit, jaká je v novém modelu střední hodnota součtů čtverců.

Nejprve uvedeme několik pomocných mezivýsledků. K jejich výpočtu nám pomohly předpoklady o nezávislosti veličin B a ε a z toho plynoucích vztahů pro střední hodnoty: $\mathbf{E}(B_i B_j) = 0$, $\mathbf{E}(B_i B_i) = d^2$, $\mathbf{E}(B_i \varepsilon_{k,j}) = 0$, $\mathbf{E}(\varepsilon_{i,j} \varepsilon_{i,l}) = 0$. Dále využíváme znalostí o rozptylu a kovarianci pozorování $Y_{i,j}$, odvozené v předchozích částech.

$$\mathbf{E}(Y_{i,j})^2 = \mu^2 + d^2 + \sigma^2$$

$$\mathbf{E}(Y_{i,j} Y_{i,k}) = \mu^2 + d^2$$

$$\mathbf{E}(Y_{i,j} Y_{u,v}) = \mu^2$$

$$\mathbf{E}(\bar{Y}_{i,+}^2) = \mu^2 + d^2 + \frac{1}{n}\sigma^2$$

$$\mathbf{E}(\bar{Y}_{++}^2) = \mu^2 + \frac{1}{k}d^2 + \frac{1}{kn}\sigma^2$$

Nyní přikročíme k výpočtu středních hodnot jednotlivých součtů čtverců.

$$\mathbf{E}(SS_c) = \mathbf{E}\left(\sum_{i=1}^k \sum_{j=1}^n (Y_{i,j} - \bar{Y}_{++})^2\right) = \mathbf{E}\left(\sum_{i=1}^k \sum_{j=1}^n (Y_{i,j}^2 - 2Y_{i,j}\bar{Y}_{++} + \bar{Y}_{++}^2)\right) =$$

$$= \mathbf{E}\left(\sum_{i=1}^k \sum_{j=1}^n (Y_{i,j}^2)\right) - 2\mathbf{E}\left(\sum_{i=1}^k \sum_{j=1}^n Y_{i,j}\bar{Y}_{++}\right) + \mathbf{E}\left(\sum_{i=1}^k \sum_{j=1}^n \bar{Y}_{++}^2\right) =$$

$$= kn\mathbf{E}(Y_{i,j})^2 - 2\mathbf{E}(kn\bar{Y}_{++}\bar{Y}_{++}) + kn\mathbf{E}(\bar{Y}_{++}^2) =$$

$$= kn(\mu^2 + d^2 + \sigma^2) - kn\left(\mu^2 + \frac{1}{k}d^2 + \frac{1}{kn}\sigma^2\right) = d^2(kn - n) + \sigma^2(kn - 1)$$

$$\mathbf{E}(SS_a) = \mathbf{E}\left(\sum_{i=1}^k n(\bar{Y}_{i,+} - \bar{Y}_{++})^2\right) = n\mathbf{E}\left(\sum_{i=1}^k \bar{Y}_{i,+}^2 - 2\bar{Y}_{i,+}\bar{Y}_{++} + \bar{Y}_{++}^2\right) =$$

$$= nk\mathbf{E}(\bar{Y}_{i,+}^2) - 2n\mathbf{E}(k\bar{Y}_{++}\bar{Y}_{++}) + kn\mathbf{E}(\bar{Y}_{++}^2) = nk\mathbf{E}(\bar{Y}_{i,+}^2) - kn\mathbf{E}(\bar{Y}_{++}^2) =$$

$$= kn\left(\mu^2 + d^2 + \frac{1}{n}\sigma^2\right) - kn\left(\mu^2 + \frac{1}{k}d^2 + \frac{1}{kn}\sigma^2\right) = d^2(kn - n) + \sigma^2(k - 1)$$

$$\begin{aligned}
\mathbb{E}(SS_e) &= \mathbb{E}\left(\sum_{i=1}^k \sum_{j=1}^n (Y_{i,j} - \bar{Y}_{i,+})^2\right) = \mathbb{E}\left(\sum_{i=1}^k \sum_{j=1}^n (Y_{i,j}^2 - 2Y_{i,j}\bar{Y}_{i,+} + \bar{Y}_{i,+}^2)\right) = \\
&= \mathbb{E}\left(\sum_{i=1}^k \sum_{j=1}^n (Y_{i,j}^2)\right) - 2\mathbb{E}\left(\bar{Y}_{i,+} \sum_{i=1}^k n\bar{Y}_{i,+}\right) + \mathbb{E}\left(\sum_{i=1}^k \sum_{j=1}^n \bar{Y}_{i,+}^2\right) = \\
&= kn \mathbb{E}(Y_{i,j})^2 - 2 \sum_{i=1}^k n \mathbb{E}(\bar{Y}_{i,+}^2) + kn \mathbb{E}(\bar{Y}_{i,+}^2) = kn(\mu^2 + d^2 + \sigma^2) - kn(\mu^2 + d^2 + \frac{1}{n}\sigma^2) \\
&= \sigma^2(kn - k)
\end{aligned}$$

Výsledky shrneme v tvrzení.

Tvrzení 6. i) Za platnosti zavedeného modelu je $(\frac{SS_a}{n(k-1)} - \frac{SS_e}{kn(n-1)})$ nestranným odhadem d^2 , $\frac{SS_e}{k(n-1)}$ je nestranným odhadem σ^2 a navíc $\frac{SS_e}{\sigma^2} \sim \chi_{k(n-1)}^2$.

ii) Pokud navíc platí $d^2 = 0$, pak $\mathbb{E}\frac{SS_a}{k-1} = \sigma^2$ a $\frac{SS_a}{\sigma^2} \sim \chi_{k-1}^2$.

iii) Platí-li $d^2 = 0$, jsou SS_a a SS_e nezávislé.

Důkaz.

i) Důkaz nestrannosti odhadů plyne z předcházejících výpočtů, rozdělení $\frac{SS_e}{\sigma^2}$ pak bylo odvozeno v kapitole 2.

ii) Pokud je člen d^2 nulový, pak $\mathbb{E}(SS_a) = \sigma^2(k-1)$ plyne z předchozích výpočtů.

Pokusme se najít vyjádření SS_a jako kvadratické formy a potom využít větu 3, která dává do souvislosti kvadratické formy a χ^2 rozdělení. Položme $M = kn$. Za předpokladu $d^2 = 0$ můžeme všechna pozorování $Y_{i,j}$ upořádat do vektoru \mathbf{Y} , pro který platí $\mathbf{Y} \sim N_M(\boldsymbol{\mu}\mathbf{1}_M, \sigma^2\mathbb{I}_M)$. Struktura rozptylové matice se velmi zjednodušila, neboť vymizely členy d^2 a také máme vyvážené třídění. Budeme tedy postupovat analogicky jako u pevných efektů. Nejprve zadefinujeme vektor skupinových průměrů $\bar{\mathbf{Y}} = (\bar{Y}_{1,+}, \dots, \bar{Y}_{k,+})$, pro který platí, že $\bar{\mathbf{Y}} - \boldsymbol{\mu}\mathbf{1}_k$ má rozdělení $N_k(\mathbf{0}, \sigma^2 \text{diag}(\frac{1}{n}, \dots, \frac{1}{n}))$. Pro následující výpočet označme $\mathbf{n} = n\mathbf{1}_k$.

$$\begin{aligned}
SS_a &= \sum_{i=1}^k n(\bar{Y}_{i,+} - \bar{Y}_{++})^2 = \sum_{i=1}^k n(\bar{Y}_{i,+} - (\frac{1}{M}\mathbf{n}^T\bar{\mathbf{Y}}))^2 = \\
&= (\bar{\mathbf{Y}} - \mathbf{1}_k \frac{1}{M}\mathbf{n}^T\bar{\mathbf{Y}})^T \text{diag}(\mathbf{n})(\bar{\mathbf{Y}} - \mathbf{1}_k \frac{1}{M}\mathbf{n}^T\bar{\mathbf{Y}}) = \\
&= \bar{\mathbf{Y}}^T \text{diag}(\mathbf{n})\bar{\mathbf{Y}} - \bar{\mathbf{Y}}^T \frac{1}{M}\mathbf{n}\mathbf{n}^T\bar{\mathbf{Y}} = \bar{\mathbf{Y}}^T \mathbb{C}\bar{\mathbf{Y}}, \\
&\text{kde } \mathbb{C} = \text{diag}(\mathbf{n}) - \frac{1}{M}\mathbf{n}\mathbf{n}^T.
\end{aligned}$$

Našli jsme vyjádření SS_a coby kvadratické formy, nyní stačí ukázat, že je matice $\text{diag}^{-1}(\mathbf{n})\mathbb{C}$ idempotentní. Ověříme definici idempotence.

$$\begin{aligned}
(\text{diag}^{-1}(\mathbf{n})\mathbb{C})(\text{diag}^{-1}(\mathbf{n})\mathbb{C}) &= (\mathbb{I}_k - \frac{1}{M}\mathbf{n}\mathbf{1}_k^T)(\mathbb{I}_k - \frac{1}{M}\mathbf{n}\mathbf{1}_k^T) = \\
&= \mathbb{I}_k\mathbb{I}_k - \frac{2}{M}\mathbf{n}\mathbf{1}_k^T + \frac{1}{M^2}\mathbf{n}\mathbf{1}_k^T\mathbf{n}\mathbf{1}_k^T = \mathbb{I}_k - \frac{1}{M}\mathbf{n}\mathbf{1}_k^T
\end{aligned}$$

S využitím věty o χ^2 rozdělení a stopě matice pak dostáváme následující vztah pro rozdělení SS_a .

$$\frac{SS_a}{\sigma^2} = \frac{1}{\sigma^2} \sigma^2 \sim \chi_{tr(diag^{-1}(\mathbf{n})\mathbb{C})}^2 = \chi_{k-1}^2.$$

iii) Součty čtverců SS_a i SS_e lze přepsat jako kvadratické formy, ověříme tedy předpoklad věty 4.

$$SS_e = \mathbf{Y}^T \mathbb{A} \mathbf{Y} = \mathbf{Y}^T \text{diag}(\mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) \mathbf{Y},$$

$$SS_a = \mathbf{Y}^T \mathbb{H}^T \text{diag}^{-1}(\mathbf{n}) \mathbb{C} \text{diag}^{-1}(\mathbf{n}) \mathbb{H} \mathbf{Y},$$

kde $\mathbb{H} = \text{diag}(\mathbf{1}_n^T, \dots, \mathbf{1}_n^T)$ je matice typu $k \times M$.

Potom $\mathbb{H}^T \text{diag}^{-1}(\mathbf{n}) \mathbb{C} \text{diag}^{-1}(\mathbf{n}) \mathbb{H} \sigma^2 \mathbb{I}_M \mathbb{A} = \sigma^2 \mathbb{H}^T \text{diag}^{-1}(\mathbf{n}) \mathbb{C} \text{diag}^{-1}(\mathbf{n}) \mathbb{H} \mathbb{A} = 0$, neboť $\mathbf{1}_n^T (\mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) = 0$. Tím je podmínka pro nezávislost ověřena. \square

3.4 Testování pro analýzu rozptylu s náhodnými efekty

Chceme testovat, zda mezi (podmíněnými) středními hodnotami jednotlivých skupin panuje shoda, nebo se statisticky významně liší. Odvodili jsme, že pro danou i -tou skupinu je střední hodnota podmíněná veličinou B_i rovna $\mu + B_i$, kde náhodná veličina B_i nyní modeluje odlišnost mezi skupinami, navíc předpokládáme, že $E B_i = 0$ a $\text{var} B_i = d^2$. Tedy v případě, že všechny skupiny mají skoro jistě shodnou střední hodnotu, musí platit, že $B_i = 0 \forall i$, což je nyní vyjádřeno stavem, kdy $d^2 = 0$. Vyslovíme tedy nulovou hypotézu pomocí rozptylu d^2 . Budeme testovat:

$$H_0 : d^2 = 0 \text{ proti } H_1 : d^2 \neq 0.$$

Z tvrzení 6 plyne, že za platnosti H_0 mají (po vydělení σ^2) skupinový a residuální součet čtverců χ^2 rozdělení s příslušnými stupni volnosti, navíc jsou nezávislé. Můžeme je tedy využít pro konstrukci testové statistiky

$$F = \frac{\frac{SS_a}{k-1}}{\frac{SS_e}{k(n-1)}} \stackrel{H_0}{\sim} F_{k-1, k(n-1)}.$$

Stejně jako v případě s pevnými efekty budou velké hodnoty F svědčit pro neplatnost H_0 . Nulovou hypotézu zamítneme právě tehdy, když F je větší, než $(1 - \alpha)$ kvantil Fisherova rozdělení s $k - 1, k(n - 1)$ stupni volnosti.

4. Simulace

Předchozí výsledky pro model s náhodnými efekty byly odvozeny za předpokladu normality pozorování $Y_{i,j}$, kterou implikovalo normální rozdělení $\varepsilon_{i,j}$ a B_i . Díky normalitě jsme mohli odvodit χ^2 rozdělení součtů čtverců a sestavit testovou statistiku. Co když byl ale předpoklad normality chybný a $\varepsilon_{i,j}$ má jiné rozdělení? Zajímá nás, zda bude testová statistika (alespoň asymptoticky) dodržovat hladinu.

Provedeme simulace v programu R. Předpokládáme, že pro naše data platí:

$$Y_{i,j} = \mu + B_i + \varepsilon_{i,j} \stackrel{H_0}{=} \mu + \varepsilon_{i,j}, \text{ bez újmy na obecnosti položíme } \mu = 0.$$

Porovnáme případ, kdy $\varepsilon_{i,j}$ pochází z normálního rozdělení s případy, kdy pochází z rozdělení s těžším chvostem a šikmého rozdělení. Konkrétně budeme simulovat nejprve případ $\varepsilon_{i,j} \sim N(0,3)$, poté $\varepsilon_{i,j} \sim t_3$ a $\ln(\varepsilon_{i,j}) \sim N(0,3)$. Pro potřeby simulace přeškálujeme lognormální rozdělení (odečteme střední hodnotu a vydělíme směrodatnou odchylkou), aby $E \varepsilon_{i,j} = 0$ a $\text{var } \varepsilon_{i,j} = 3$.

Postupně budeme generovat data pro $k = 10, 50, 200, 500$ skupin a 10 pozorování v každé skupině, pro tato data pak spočítáme příslušnou F-statistiku. Pro každý případ vygenerujeme F-statistiku 10 000 krát, abychom měli reprezentativní počet výsledků. Dále budeme počítat odhad pravděpodobnosti, že za platnosti nulové hypotézy hodnota testové statistiky F převyší $(1 - \alpha)$ kvantil Fisherova rozdělení s danými stupni volnosti, neboli pravděpodobnost, že i za platnosti H_0 dojde k jejímu zamítnutí (chyba 1. druhu).

Použijeme hladinu $\alpha = 0.05$, výsledné empirické odhady této pravděpodobnosti jsou shrnuty v následující tabulce. Vidíme, že i pro t_3 a lognormální rozdělení jsme obdrželi výsledky blízké hodnotě 0.05, z toho usuzujeme, že test dodržuje hladinu i při námi uvažovaném porušení normality.

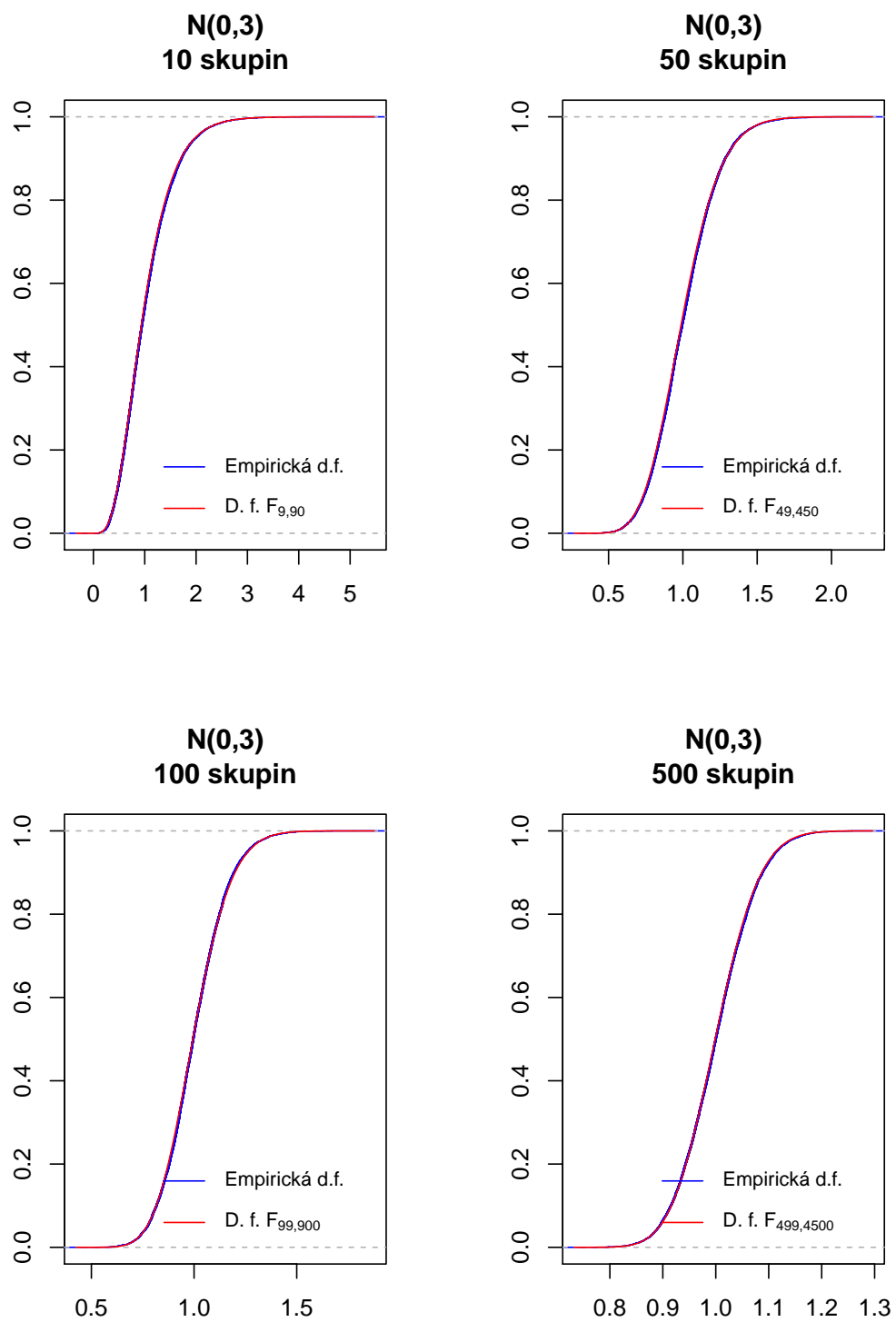
Počet skupin	Rozdělení		
	N(0,3)	t_3	LN(0,3)
10	0.0512	0.0445	0.0311
50	0.0478	0.0474	0.0485
100	0.0475	0.0441	0.0469
500	0.0512	0.0483	0.054

Tabulka 4.1: Výsledné odhady na základě simulací

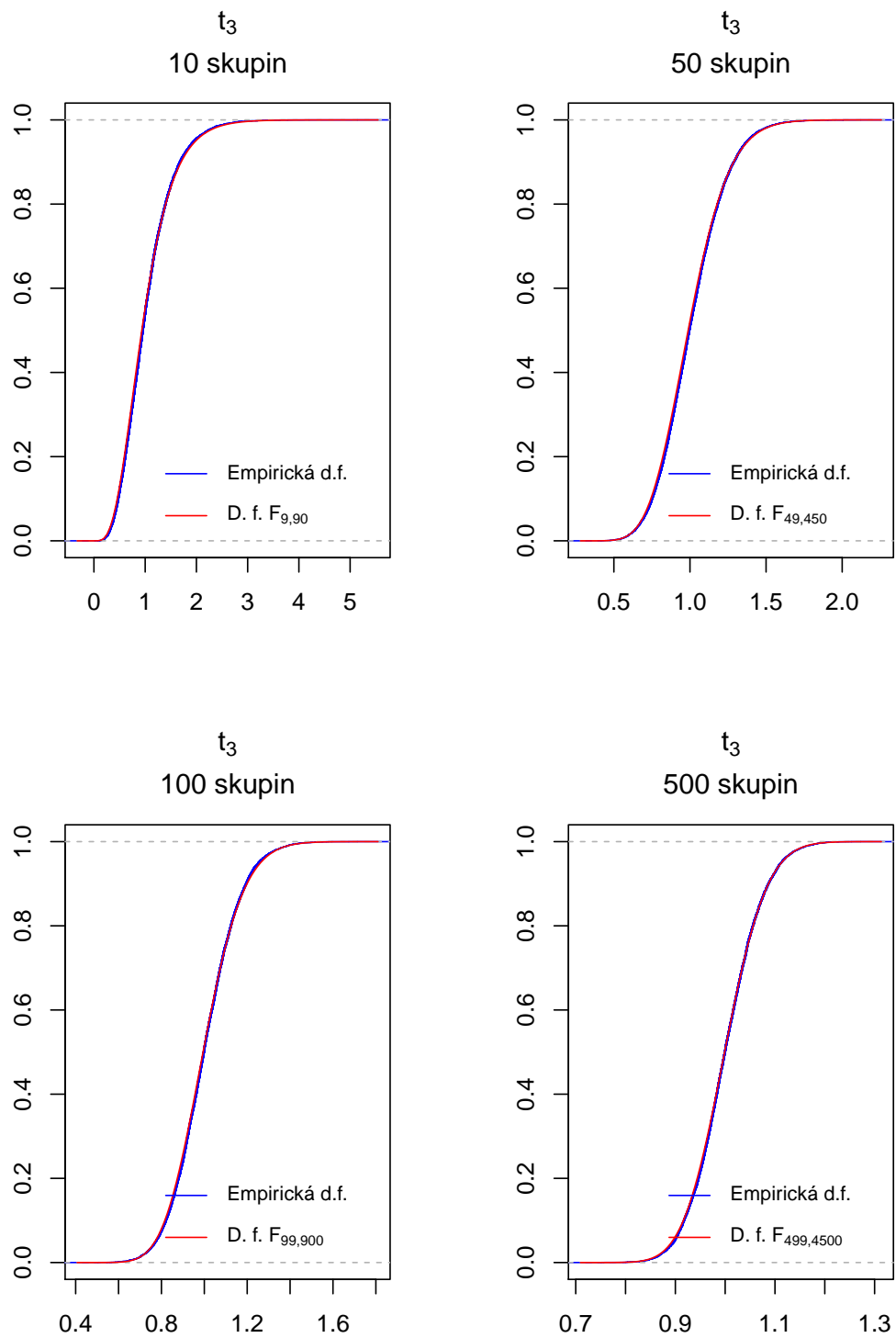
Dále pro každý případ porovnáme empirickou distribuční funkci testové statistiky F s distribuční funkcí Fisherova rozdělení s příslušnými stupni volnosti. Grafy s tímto porovnáním jsou na následujících stránkách.

Vidíme, že v případě $N(0,3)$ i t_3 rozdělení empirická distribuční funkce testové statistiky F téměř splývá se skutečnou distribuční funkcí Fisherova rozdělení, a to už pro 10 skupin. V případě lognormálního rozdělení mezi nimi pozorujeme rozdíl, ale pro rostoucí počet skupin k se empirická funkce svým tvarem více podobá skutečné distribuční funkci F rozdělení s příslušnými stupni volnosti.

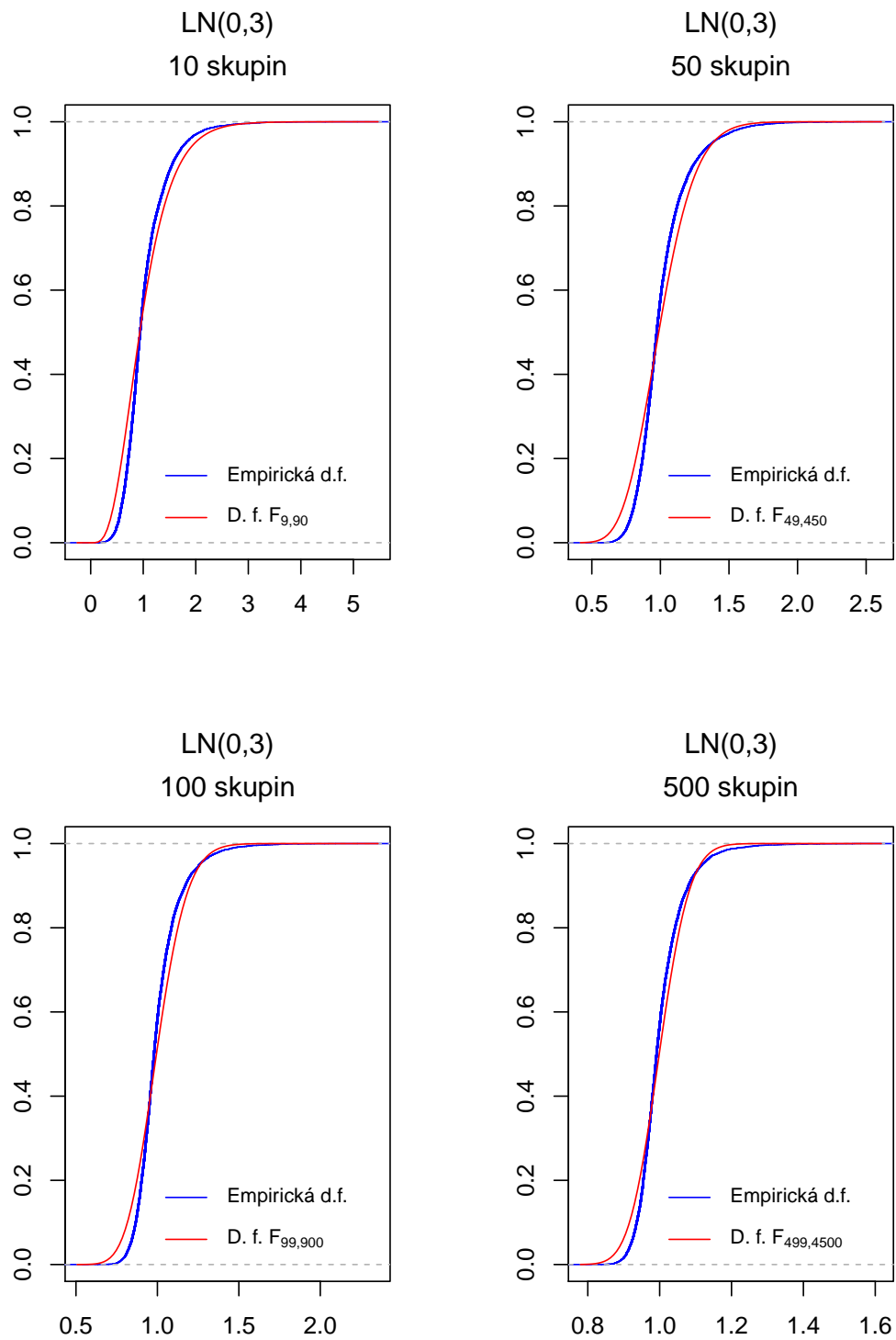
Můžeme tedy učinit závěr, že pro dostatečně velký počet pozorování dodržuje test analýzy rozptylu s náhodnými efekty hladinu i v případě, kdy pozorování nepochází z normálního rozdělení.



Obrázek 4.1: Porovnání distribučních funkcí pro normální rozdělení



Obrázek 4.2: Porovnání distribučních funkcí pro t_3 rozdělení



Obrázek 4.3: Porovnání distribučních funkcí pro lognormální rozdělení

Závěr

V této práci jsme se zabývali analýzou rozptylu s náhodnými efekty. Vyšli jsme ze základních tvrzení z teorie pravděpodobnosti, podrobněji jsem se zaměřili na vztahy platící pro rozdělení kvadratických forem.

Dále jsme zavedli případ jednoduchého třídění s pevnými efekty. Podrobně jsme definovali součty čtverců a ukázali jejich vlastnosti, díky kterým jsme je mohli použít k sestavení testové statistiky.

V další části jsme nejprve zavedli model jednoduchého třídění s náhodnými efekty a podrobně jsme odvodili vlastnosti, která plynou pro pozorování a skupiny pozorování na základě tohoto modelu. S předpokladem vyváženého třídění, tj. všechny skupiny mají stejný počet pozorování, jsme zavedli součty čtverců, pro které jsme ukázali, že (za platnosti dodatečného předpokladu na nulovost rozptylu d^2 a po vydělení příslušnými stupni volnosti) jsou skupinový a residuální součet čtverců nezávislými odhady rozptylu σ^2 . V závěru kapitoly jsme se zabývali testem na shodu podmíněných středních hodnot skupin, na základě vlastností součtů čtverců jsme setrojili testovou statistiku a kritický obor.

Na závěr jsme se zabývali otázkou, jak test dodržuje hladinu při porušení předpokladu normality. Provedli jsme simulace v programu R pro pevný počet pozorování ve skupině a různé počty skupin, data jsme nejprve generovali z normálního rozdělení a poté porovnávali jejich chování s daty generovanými z t_3 a lognormálního rozdělení, coby reprezentanty rozdělení s těžším chvostem a šikmého rozdělení. Zjistili jsme, že pro případ vyváženého třídění test dodržuje hladinu velmi dobře i v případě uvažovaného porušení normality.

Seznam použité literatury

- ANDĚL, J. (2007a). *Základy matematické statistiky*. Druhé opravené vydání. Matfyzpress, Praha. ISBN 80-7378-001-1.
- ANDĚL, J. (2007b). *Statistické metody*. Čtvrté upravené vydání. Matfyzpress, Praha. ISBN 80-7378-003-8.
- KHURI, A. I. (2010). *Linear model methodology*. Chapman and Hall/CRC. ISBN 978-1-58488-481-1.
- LACHOUT, P. (2004). *Teorie pravděpodobnosti*. Druhé vydání. Karolinum, Praha. ISBN 80-246-0872-3.