



**MATEMATICKO-FYZIKÁLNÍ
FAKULTA**
Univerzita Karlova

BAKALÁŘSKÁ PRÁCE

Julie Váňová

**Testy dobré shody s Poissonovým
rozdělením založené na nulovém indexu**

Katedra pravděpodobnosti a matematické statistiky

Vedoucí bakalářské práce: RNDr. Šárka Hudecová, Ph.D.

Studijní program: Obecná matematika

Praha 2023

Prohlašuji, že jsem tuto bakalářskou práci vypracoval(a) samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů. Tato práce nebyla využita k získání jiného nebo stejného titulu.

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

Na tomto místě bych chtěla poděkovat vedoucí své práce RNDr. Šárce Hudecové, Ph.D. za velkou ochotu, vstřícnost a také čas, který mi věnovala, ať už během konzultací, nebo při pečlivém pročitání jednotlivých verzí této práce.

Název práce: Testy dobré shody s Poissonovým rozdělením založené na nulovém indexu

Autor: Julie Váňová

Katedra: Katedra pravděpodobnosti a matematické statistiky

Vedoucí bakalářské práce: RNDr. Šárka Hudecová, Ph.D., Katedra pravděpodobnosti a matematické statistiky

Abstrakt: Tato bakalářská práce se zabývá metodou testování dobré shody s Poissonovým rozdělením, která je založená na takzvaném nulovém indexu. Nejprve je zaveden pojem Poissonova nulového indexu a jsou diskutovány jeho základní vlastnosti. Následně se práce soustředí na odvození asymptotického rozdělení nulových indexů a jeho využití ke konstrukci testů dobré shody, přičemž jsou uvedeny i konkrétní příklady nulových indexů a následných testů. V další části jsou potom stručně popsány některé další způsoby, kterými lze testovat dobrou shodu s Poissonovým rozdělením, konkrétně jde o χ^2 -testy dobré shody a testy založené na indexu disperze. Zmíněné metody testování jsou poté porovnány v simulační studii.

Klíčová slova: Poissonovo rozdělení, testy dobré shody, asymptotické testy, nulový index

Title: Goodness-of-fit tests for Poisson distribution based on zero index

Author: Julie Váňová

Department: Department of Probability and Mathematical Statistics

Supervisor: RNDr. Šárka Hudecová, Ph.D., Department of Probability and Mathematical Statistics

Abstract: This bachelor thesis deals with goodness-of-fit tests for Poisson distribution that are based on so-called zero index. In the first part, Poisson zero index is defined and some of its basic properties are discussed. Further, asymptotic distribution of zero indexes is derived and it is used to construct asymptotic goodness-of-fit tests. Particular examples of zero indexes and related tests are included. In the following part, other types of goodness-of-fit tests for Poisson distribution are briefly described, in particular χ^2 -tests and tests based on index of dispersion. All mentioned methods are then compared in a simulation study.

Keywords: Poisson distribution, goodness-of-fit tests, asymptotic tests, zero index

Obsah

Úvod	2
1 Uvedení do problému	3
1.1 Testy dobré shody	3
1.2 Poissonovo rozdělení	3
1.3 Některé pomocné věty a tvrzení	5
2 Testy založené na nulovém indexu	7
2.1 Nulový index – definice a vlastnosti	7
2.2 Příklady Poissonových nulových indexů	11
2.3 Testování pomocí nulových indexů	16
2.4 Úskalí testování založeného na nulových indexech	20
3 Další varianty testů dobré shody	23
3.1 χ^2 -testy dobré shody	23
3.2 Testy založené na indexu disperze	25
4 Simulační studie	27
4.1 Hladina testů	27
4.2 Síla testů	30
Závěr	34
Seznam použité literatury	35

Úvod

Testy dobré shody s nějakým předem určeným rozdělením patří mezi nejznámější statistické testy. V této práci se budeme zabývat testováním dobré shody s Poissonovým rozdělením, tedy rozdělením, které nachází mnoho využití i v praxi (např. v modelování počtu pojistných událostí). Zaměříme se při tom především na metodu, která je založena na tzv. nulovém indexu. Její hlavní myšlenka spočívá v ověřování, zda platí určitý vztah mezi relativním počtem nul v realizacích nějakého náhodného výběru a jeho výběrovým průměrem.

Idea tohoto způsobu testování je popsána v článku Weiß, Homburg a Puig (2019), který je hlavním výchozím bodem této práce. V souladu s tímto článkem také definujeme již zmíněný nulový index. Výše uvedený článek ovšem zavádí tento index v trochu odlišném kontextu (INAR(1) modely časových řad). My se zaměříme na situaci, kdy máme k dispozici náhodný výběr z nějakého diskrétního rozdělení, a zajímá nás, zda se jedná o rozdělení Poissonovo.

V první kapitole formalizujeme zkoumaný problém jako testování určité nulové hypotézy a uvedeme základní vlastnosti Poissonova rozdělení. Ve druhé se potom budeme podrobně zabývat testy dobré shody založenými na nulových indexech. Odvodíme si asymptotické rozdělení těchto indexů (za platnosti hypotézy) a ukážeme, jak lze pomocí nich testovat dobrou shodu s Poissonovým rozdělením. To vše provedeme nejdříve obecně, a následně i pro konkrétní případy. Ve třetí kapitole stručněji popíšeme i další typy testů dobré shody. Jednotlivé způsoby testování poté porovnáme v simulační studii – podíváme se, jak dobře dané testy dodržují stanovenou hladinu významnosti a „spočteme“ jejich sílu proti konkrétní alternativě (negativně binomickému rozdělení).

Vlastní přínos této práce spočívá především v uvádění vysvětlujících komentářů, formulování vlastních vět a tvrzení v kapitole týkající se nulových indexů (kapitola 2) a podrobném sepsání jejich důkazů. Dále potom v odvození asymptotického rozdělení konkrétních nulových indexů, z nichž dva ($\hat{I}_{3,n}$ a $\hat{I}_{4,n}$) jsme si sami zkonstruovali. Částečně také v zamyšlení se nad různými úskalími testování založeného na nulových indexech (sekce 2.4). Za vlastní příspěvek lze rovněž považovat simulační studii ve čtvrté kapitole.

1. Uvedení do problému

Počet dopravních nehod v Praze za jeden týden, množství gólů během jednoho fotbalového zápasu, počet nových případů nákazy danou nemocí za den – to vše jsou příklady dat, na něž můžeme nahlížet jako na realizace určitých náhodných veličin s diskrétním rozdělením. Přestože jejich konkrétní rozdělení neznáme, máme k dispozici metody, které nám umožní udělat si o něm lepší představu.

1.1 Testy dobré shody

Jednou z možností je testovat, do jaké míry naše data odpovídají hypotetickým realizacím nějakého předem stanoveného rozdělení. K tomu slouží tzv. testy dobré shody.

Uvedený problém nyní trochu zúžíme. Budeme předpokládat, že data, která máme k dispozici, představují náhodný výběr z nějakého diskrétního rozdělení. To znamená, že události, které zkoumáme, nastávají nezávisle na sobě a pravděpodobnost výskytu je pro všechny z nich stejná. Zde se navíc nebudeme zabývat obecnými testy dobré shody, ale pouze jejich verzí pro Poissonovo rozdělení.

Matematicky formulováno, budeme uvažovat náhodný výběr X_1, X_2, \dots, X_n pocházející z určitého diskrétního rozdělení. V naší situaci dává navíc smysl omezit se pouze na rozdělení, jejichž nosič je podmnožina nezáporných celých čísel (jak později uvidíme, takový nosič má i Poissonovo rozdělení). Budeme tedy pracovat s neparametrickým modelem

$$\mathcal{F} = \{\text{všechna diskrétní rozdělení, jejichž nosič je podmnožina } \mathbb{N}_0\}.$$

Předmětem testování bude celé pravděpodobnostní rozdělení našeho náhodného výběru. Zajímá nás, jak moc se toto rozdělení shoduje s Poissonovým rozdělením s parametrem λ , pro nějaké blíže nespecifikované (pevné) $\lambda > 0$. Budeme tedy testovat hypotézu H_0 proti alternativě H_1 , kde

$$H_0 : \exists \lambda \in (0, \infty) \text{ takové, že } X_1 \sim \text{Po}(\lambda), \quad (1.1)$$

$$H_1 : \nexists \lambda \in (0, \infty) \text{ takové, že } X_1 \sim \text{Po}(\lambda).$$

Místo $X_1 \sim \text{Po}(\lambda)$ lze ekvivalentně psát $X_i \sim \text{Po}(\lambda)$, $i = 1, \dots, n$, neboť uvažované náhodné veličiny jsou dle předpokladu stejně rozdělené.

Nulová hypotéza říká, že pro nějaké λ má náhodný výběr X_1, \dots, X_n požadované rozdělení. Konkrétní hodnota parametru λ je ovšem neznámá a její určení ani není předmětem tohoto testu. Naopak podle alternativy žádné λ , pro které by platilo $X_1 \sim \text{Po}(\lambda)$, neexistuje. Zamítnutím nulové hypotézy ve prospěch alternativy potom pouze tvrdíme, že daný náhodný výběr z Poissonova rozdělení nepochází, o jeho skutečném rozdělení toho však stále mnoho nevíme.

1.2 Poissonovo rozdělení

Jak již bylo zmíněno, budeme se zabývat především testováním dobré shody s Poissonovým rozdělením. V této sekci si připomeneme definici tohoto pravděpodobnostního rozdělení a podíváme se na některé jeho vlastnosti.

Definice 1 (Poissonovo rozdělení). Řekneme, že náhodná veličina X má **Poissonovo rozdělení** s parametrem $\lambda > 0$, jestliže platí

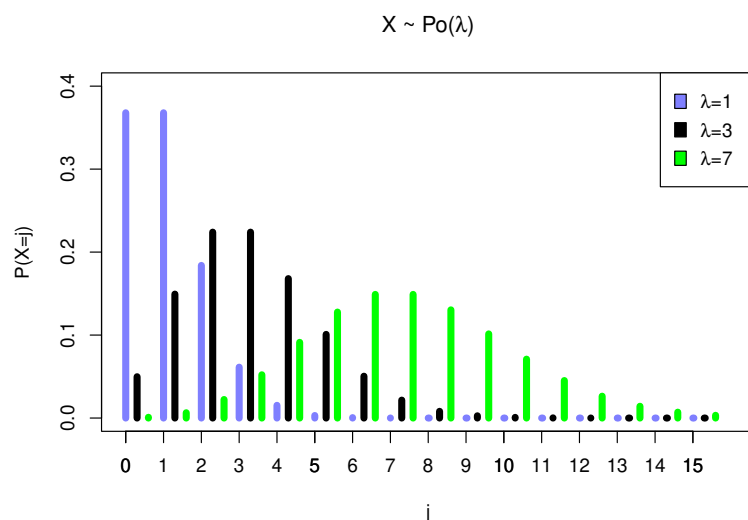
$$P(X = j) = e^{-\lambda} \cdot \frac{\lambda^j}{j!}, \quad j \in \mathbb{N}_0.$$

Značíme $X \sim \text{Po}(\lambda)$.

Můžeme si všimnout, že $P(X = j) > 0$ pro všechna $j \in \mathbb{N}_0$ a

$$\sum_{j=0}^{\infty} e^{-\lambda} \cdot \frac{\lambda^j}{j!} = e^{-\lambda} \sum_{j=0}^{\infty} \frac{\lambda^j}{j!} = e^{-\lambda} \cdot e^{\lambda} = 1,$$

tedy se skutečně jedná o pravděpodobnostní rozdělení. Z definice je dále vidět, že náhodná veličina se zmíněným rozdělením má kladnou pravděpodobnost výskytu u nekonečně mnoha různých hodnot. Poissonovo rozdělení je diskrétní rozdělení, neboť jeho nosič \mathbb{N}_0 je spočetná množina. Obrázek 1.1 zachycuje toto rozdělení pro vybrané hodnoty parametru λ .



Obrázek 1.1: Poissonovo rozdělení pro vybrané hodnoty parametru λ

Při práci s daným rozdělením se rovněž vyplatí znát některé jeho momenty, především střední hodnotu a rozptyl. Pro náhodnou veličinu $X \sim \text{Po}(\lambda)$, $\lambda > 0$, spočteme její střední hodnotu snadno:

$$EX = \sum_{k=0}^{\infty} k \cdot P(X = k) = \sum_{k=0}^{\infty} k \cdot e^{-\lambda} \cdot \frac{\lambda^k}{k!} = \sum_{k=1}^{\infty} \lambda e^{-\lambda} \cdot \frac{\lambda^{k-1}}{(k-1)!} = \lambda \sum_{j=0}^{\infty} e^{-\lambda} \cdot \frac{\lambda^j}{j!} = \lambda,$$

kde v poslední rovnosti využíváme, že Poissonovo rozdělení je pravděpodobnostní rozdělení, tedy je daná suma rovna jedné. Dále platí

$$\begin{aligned} EX^2 &= \sum_{k=0}^{\infty} k^2 \cdot P(X = k) = \sum_{k=0}^{\infty} k^2 \cdot e^{-\lambda} \cdot \frac{\lambda^k}{k!} = \sum_{k=1}^{\infty} \lambda k \cdot e^{-\lambda} \cdot \frac{\lambda^{k-1}}{(k-1)!} \\ &= \lambda \sum_{j=0}^{\infty} (j+1) \cdot e^{-\lambda} \cdot \frac{\lambda^j}{j!} = \lambda \left(\sum_{j=0}^{\infty} j \cdot e^{-\lambda} \cdot \frac{\lambda^j}{j!} + \sum_{j=0}^{\infty} e^{-\lambda} \cdot \frac{\lambda^j}{j!} \right) = \lambda(\lambda + 1). \end{aligned}$$

Potom

$$\text{var}X = \text{E}X^2 - (\text{E}X)^2 = \lambda(\lambda + 1) - \lambda^2 = \lambda.$$

Vidíme tedy, že střední hodnota a rozptyl náhodné veličiny s Poissonovým rozdělením si jsou rovny. Později uvidíme, že tuto speciální vlastnost lze použít i ke konstrukci testů dobré shody založených na tzv. indexu disperze.

Poissonovým rozdělením se řídí zejména veličiny, které udávají počty nějakých událostí (např. počty nehod, pojistných událostí, ...) za určitou časovou jednotku.

1.3 Některé pomocné věty a tvrzení

V této části uvedeme znění několika důležitých vět a tvrzení, na něž se budeme odkazovat v nadcházejících kapitolách.

Tvrzení 1. *Nechť $X_1, X_2, X_3, \dots, X_n$ je náhodný výběr z libovolného rozdělení, které má konečný první moment.*

(i) *Výběrový průměr $\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$ splňuje*

$$\text{E}(\bar{X}_n) = \text{E}X_1, \quad \bar{X}_n \xrightarrow[n \rightarrow \infty]{\text{P}} \text{E}X_1.$$

(ii) *Předpokládejme navíc, že $\text{var}X_1 \in (0, \infty)$. Výběrový rozptyl*

$$S_n^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

splňuje

$$\text{E}(S_n^2) = \text{var}X_1, \quad S_n^2 \xrightarrow[n \rightarrow \infty]{\text{P}} \text{var}X_1.$$

(iii) *Nechť $A \subset \mathbb{R}$ je borelovská množina. Označme $Y_i := \mathbf{1}_{[X_i \in A]}$, $i = 1, \dots, n$, tj.*

$$Y_i = \begin{cases} 1, & \text{pokud } X_i \in A, \\ 0, & \text{pokud } X_i \notin A. \end{cases}$$

Potom Y_i má alternativní rozdělení s pravděpodobností úspěchu $p := \text{P}(X_1 \in A)$. Dále pro relativní četnost jevu $[X_i \in A]$ definovanou jako $\bar{Y}_n = \frac{1}{n} \sum_{j=1}^n \mathbf{1}_{[X_j \in A]}$ platí

$$\text{E}(\bar{Y}_n) = p, \quad \bar{Y}_n \xrightarrow[n \rightarrow \infty]{\text{P}} p.$$

Důkaz: Části (i) a (ii) jsou dokázány v knize Dupač a Hušková (2013, Věta 5.1) s tím, že ta je oproti našemu tvrzení formulovaná pro konvergenci skoro jistě. Kýžený výsledek tak plyne například z (i) ve větě 2.7 v knize van der Vaart (1998).

(iii) je důsledkem části (i), neboť z definice střední hodnoty

$$\text{E}(Y_1) = \text{E}(\mathbf{1}_{[X_1 \in A]}) = \text{P}(X_1 \in A) = p.$$

□

Nadcházející věta je převzatá z Anděl (2007, Věta B.5). Její důkaz lze najít například v knize Anděl (1978).

Věta 2 (Mnohorozměrná Lindebergova centrální limitní věta).

Nechť $\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3, \dots$ jsou k -rozměrné nezávislé a stejně rozdělené náhodné vektory, které mají rozdělení se střední hodnotou $\boldsymbol{\mu}$ a konečnou varianční maticí $\boldsymbol{\Sigma}$. Potom pro $n \rightarrow \infty$ platí

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{X}_i - \boldsymbol{\mu}) \xrightarrow{D} \mathbf{N}_k(\mathbf{0}, \boldsymbol{\Sigma}).$$

Pro potřeby následující věty připomeneme, že pro $\mathbf{x} \in \mathbb{R}^k$ značí $\mathbb{D}g(\mathbf{x})$ Jacobiho matici, tj.

$$\mathbb{D}g(\mathbf{x}) = \begin{pmatrix} \frac{\partial g_1}{\partial x_1} & \dots & \frac{\partial g_1}{\partial x_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_p}{\partial x_1} & \dots & \frac{\partial g_p}{\partial x_k} \end{pmatrix}.$$

Věta 3 (Δ -věta).

Nechť $\{\mathbf{T}_n\}_{n=1}^{\infty}$ jsou náhodné vektory splňující $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) \xrightarrow[n \rightarrow \infty]{D} \mathbf{N}_k(\mathbf{0}, \boldsymbol{\Sigma})$ pro nějaký vektor konstant $\boldsymbol{\mu} \in \mathbb{R}^k$ a pozitivně semi-definitní matici $\boldsymbol{\Sigma} \in \mathbb{R}^{k \times k}$. Nechť dále zobrazení $g : \mathbb{R}^k \rightarrow \mathbb{R}^p$ je spojitě diferencovatelné na nějakém okolí bodu $\boldsymbol{\mu}$. Potom

$$\sqrt{n}(g(\mathbf{T}_n) - g(\boldsymbol{\mu})) \xrightarrow{D} \mathbf{N}_p(\mathbf{0}, \mathbb{D}g(\boldsymbol{\mu})\boldsymbol{\Sigma}(\mathbb{D}g(\boldsymbol{\mu}))^\top) \text{ pro } n \rightarrow \infty.$$

Důkaz: Viz van der Vaart (1998, Věta 3.1). □

2. Testy založené na nulovém indexu

V předchozí kapitole jsme představili obecný koncept testů dobré shody. Nyní se již podíváme na konkrétní případ těchto testů, který je založen na tzv. nulovém indexu. Při definování nulového indexu budeme vycházet z článku Weiß a kol. (2019), kde je ovšem tento index zaveden v o něco obecnějším kontextu. Velká část vět, tvrzení a komentářů, které budou následovat, se tak v článku nevyskytuje.

2.1 Nulový index – definice a vlastnosti

V sekci 1.2 jsme odvodili, že pokud $X_1 \sim \text{Po}(\lambda)$, pro nějaké $\lambda > 0$, splňuje střední hodnota rovnost $\mathbf{E}X_1 = \lambda$. Z definice Poissonova rozdělení potom vyplývá, že pro náhodnou veličinu X s tímto rozdělením platí

$$\mathbf{P}(X_1 = 0) = e^{-\lambda} = e^{-\mathbf{E}X_1}.$$

Jedna z možností, jak testovat, zda rozdělení náhodného výběru odpovídá Poissonovu rozdělení, spočívá právě v ověřování platnosti tohoto vztahu.

Skutečná hodnota parametru λ i pravděpodobnosti na levé straně rovnosti je pro nás neznámá, budeme tedy potřebovat výrazy na obou stranách nějak odhadnout.

Označme $p_0 := \mathbf{P}(X_1 = 0)$. Jako odhad pravděpodobnosti se nabízí použít relativní četnost

$$\hat{p}_0 := \frac{1}{n} \sum_{j=1}^n \mathbf{1}_{[X_j=0]}$$

jevu $[X_1 = 0]$, což je podle (ii) v tvrzení 1 nestranný a konzistentní odhad p_0 .

Poznámka. Konzistencí odhadu \hat{p}_0 rozumíme platnost vztahu $\hat{p}_0 \xrightarrow[n \rightarrow \infty]{\mathbf{P}} p_0$. Někdy se v definici konzistence místo konvergence v pravděpodobnosti uvažuje konvergence skoro jistě, v této práci však budeme vždy chápat konzistenci odhadu v souvislosti s prvním zmíněným typem konvergence.

Vhodným odhadem střední hodnoty je výběrový průměr

$$\bar{X}_n := \frac{1}{n} \sum_{j=1}^n X_j,$$

neboť na základě tvrzení 1 opět disponuje žádanými vlastnostmi nestrannosti a konzistence. Zároveň se jedná i o nejlepší nestranný odhad parametru λ (viz Anděl, 2007, Příklad 7.62), tedy odhad s nejmenším rozptylem mezi všemi nestrannými odhady λ .

Budeme-li opět uvažovat vztah $p_0 = e^{-\mathbf{E}X_1}$, můžeme levou stranu rovnosti odhadnout pomocí \hat{p}_0 a pravou stranu jako $e^{-\bar{X}_n}$, což je podle věty o spojitě transformaci (viz van der Vaart (1998), Věta 2.3) konzistentní odhad $e^{-\mathbf{E}X_1}$. Obecně bude naše testová statistika nějaká funkce $f(\bar{X}_n, \hat{p}_0)$, u níž chceme, aby nějakým způsobem reflektovala, jak moc se hodnoty p_0 a $e^{-\mathbf{E}X_1}$ liší. To vede k zavedení následujícího pojmu, jenž bude jeden z klíčových pojmů této práce.

Definice 2 (Poissonův nulový index). *Nechť $(0, \infty) \times (0, 1) \subseteq G$ a $f : G \rightarrow \mathbb{R}$ je funkce, jejíž předpis nezávisí na neznámých parametrech, splňující*

- (i) $f(x, e^{-x}) = 0$ pro $\forall x > 0$;
- (ii) $f(x, y) \neq 0$, pokud $y \neq e^{-x}$.

Nechť \bar{X}_n a \hat{p}_0 jsou definované jako výše. Potom statistika

$$f(\bar{X}_n, \hat{p}_0),$$

*je-li definovaná, se nazývá **Poissonův nulový index**.*

Věta „je-li definovaná“ se v předchozí definici vyskytuje kvůli tomu, že bod $(\bar{X}_n, \hat{p}_0)^\top$ může v některých případech ležet mimo definiční obor funkce f . Jak ale uvidíme v sekci 2.2, pokud je nosičem našeho náhodného výběru podmnožina \mathbb{N}_0 a platí $\mathbf{P}(X_1 = 0) \in (0, 1)$, pravděpodobnost, že je $f(\bar{X}_n, \hat{p}_0)$ dobře definovaná, se pro $n \rightarrow \infty$ blíží k 1.

Pro funkci f splňující (i) a (ii) v definici 2 dostaneme, že pro libovolné $\lambda > 0$:

$$f(\lambda, p_0) = 0 \quad \Leftrightarrow \quad p_0 = e^{-\lambda}.$$

Dále je dobré si uvědomit, že Poissonův nulový index je skutečně statistika, neboť \bar{X}_n ani \hat{p}_0 nezávisí na neznámých parametrech. Některé jeho vlastnosti nyní zformulujeme v následujícím tvrzení.

Tvrzení 4. *Mějme náhodný výběr X_1, X_2, \dots, X_n z libovolného diskrétního rozdělení, jehož nosič je podmnožina \mathbb{N}_0 , s konečnou střední hodnotou. Označme $p_0 := \mathbf{P}(X_1 = 0)$. Nechť f je funkce splňující (i) a (ii) v definici 2, která je spojitá v bodě $(\mathbf{E}X_1, p_0)^\top \in \mathbb{R}^2$.*

1. *Platí*

$$f(\bar{X}_n, \hat{p}_0) \xrightarrow[n \rightarrow \infty]{\mathbf{P}} f(\mathbf{E}X_1, p_0),$$

tj. $f(\bar{X}_n, \hat{p}_0)$ je konzistentní odhad $f(\mathbf{E}X_1, p_0)$.

2. *Nechť navíc $\mathbf{E}X_1 > 0$ a $p_0 = e^{-\mathbf{E}X_1}$. Potom*

$$f(\bar{X}_n, \hat{p}_0) \xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0.$$

Speciálně, výše uvedený vztah platí pro $X_1 \sim \text{Po}(\lambda)$, $\lambda > 0$.

3. *Nechť $p_0 \neq e^{-\mathbf{E}X_1}$. Potom pro $n \rightarrow \infty$*

$$f(\bar{X}_n, \hat{p}_0) \not\xrightarrow{\mathbf{P}} 0.$$

Důkaz:

1. \bar{X}_n a \hat{p}_0 jsou konzistentní odhady $\mathbf{E}X_1$ a p_0 . První tvrzení tak plyne z věty o spojitě transformaci, kterou lze použít díky předpokladu spojitosti f v bodě $(\mathbf{E}X_1, p_0)^\top$.

2. Jedná se o důsledek předchozího bodu. Označíme-li $\lambda := \mathbf{E}X_1$, dostaneme $p_0 = e^{-\lambda}$. Tudíž podle (i) v definici nulového indexu

$$f(\mathbf{E}X_1, p_0) = f(\lambda, e^{-\lambda}) \stackrel{(i)}{=} 0.$$

Dále pro Poissonovo rozdělení s parametrem λ platí $p_0 = e^{-\lambda}$ a $\mathbf{E}X_1 = \lambda$, a tedy v tomto případě je rovnost $p_0 = e^{-\mathbf{E}X_1}$ zřejmě splněna.

3. Opět plyne z 1. Z části (ii) definice 2 totiž za této situace

$$f(\mathbf{E}X_1, p_0) \neq 0.$$

Z jednoznačnosti limity tak dostaneme požadovaný výsledek. □

Za platnosti (1.1) a pro dostatečně velký rozsah náhodného výběru bude tedy podle 2. části předchozího tvrzení hodnota statistiky $f(\bar{X}_n, \hat{p}_0)$ blízká nule. Z posledního bodu naopak vidíme, že pro náhodné výběry, které splňují

$$P(X_1 = 0) \neq e^{-\mathbf{E}X_1},$$

a tedy nepocházejí z Poissonova rozdělení, $f(\bar{X}_n, \hat{p}_0)$ k nule v pravděpodobnosti nekonverguje. V případě, že hodnota tohoto Poissonova nulového indexu bude výrazně vyšší (nebo nižší) než nula, budeme tedy nejspíš schopni zamítnout nulovou hypotézu. Zde je ovšem důležité podotknout, že platnost $f(\mathbf{E}X_1, p_0) = 0$ nutně neimplikuje, že X_1 má Poissonovo rozdělení (více v sekci 2.4).

Abychom mohli provést test nulové hypotézy, potřebujeme znát rozdělení naší testové statistiky. Klíčovým nástrojem k určování asymptotického rozdělení nějaké transformace náhodného vektoru je Δ -věta (věta 3 v předchozí kapitole). Při pohledu na předpoklady Δ -věty vidíme, že bychom ji mohli použít k získání asymptotického rozdělení $f(\bar{X}_n, \hat{p}_0)$, potřebujeme znát asymptotické rozdělení náhodného vektoru $(\bar{X}_n, \hat{p}_0)^\top$. Tento výsledek nám dává následující tvrzení.

Tvrzení 5. *Nechť X_1, X_2, \dots, X_n je náhodný výběr z rozdělení $\text{Po}(\lambda)$, $\lambda > 0$. Potom platí*

$$\sqrt{n} \left(\begin{pmatrix} \bar{X}_n \\ \hat{p}_0 \end{pmatrix} - \begin{pmatrix} \lambda \\ e^{-\lambda} \end{pmatrix} \right) \xrightarrow[n \rightarrow \infty]{D} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix} \right).$$

Důkaz: Předpokládáme, že $X_i \sim \text{Po}(\lambda)$, $i = 1, \dots, n$. Použijeme mnohorozměrnou verzi centrální limitní věty (jejíž znění jsme si uvedli ve větě 2) na náhodné vektory $(X_i, \mathbf{1}_{[X_i=0]})^\top$, $i = 1, \dots, n$.

Z předpokladu nezávislosti a stejného rozdělení X_i , $i = 1, \dots, n$, dostáváme, že i náhodné vektory $(X_i, \mathbf{1}_{[X_i=0]})^\top$, $i = 1, \dots, n$, jsou nezávislé a mají stejné rozdělení, neboť obě jejich složky jsou měřitelné funkce X_i . Zbývá tedy spočítat střední hodnotu těchto vektorů a určit jejich varianční matici.

Pro libovolné $i \in \{1, \dots, n\}$ platí

$$\mathbf{E}X_i = \lambda, \quad \text{var}X_i = \lambda.$$

Podle tvrzení 1 má $\mathbf{1}_{[X_i=0]}$ alternativní rozdělení s pravděpodobností úspěchu $P(X_1 = 0) = e^{-\lambda}$, a tedy

$$\mathbf{E}(\mathbf{1}_{[X_i=0]}) = e^{-\lambda}, \quad \text{var}(\mathbf{1}_{[X_i=0]}) = e^{-\lambda}(1 - e^{-\lambda}).$$

Nyní spočteme kovarianci X_i a $\mathbf{1}_{[X_i=0]}$ použitím známého vzorce $\text{cov}(X, Y) = \mathbf{E}(XY) - \mathbf{E}X \cdot \mathbf{E}Y$.

$$\text{cov}(X_i, \mathbf{1}_{[X_i=0]}) = \mathbf{E}(X_i \cdot \mathbf{1}_{[X_i=0]}) - \mathbf{E}X_i \cdot \mathbf{E}(\mathbf{1}_{[X_i=0]}) = \mathbf{E}(0) - \lambda \cdot e^{-\lambda} = -\lambda e^{-\lambda}.$$

Zde využíváme, že pro libovolné $\omega \in \Omega$ buď $X_i(\omega) = 0$, nebo $X_i(\omega) \neq 0$. Druhý případ implikuje $\mathbf{1}_{[X_i=0]}(\omega) = 0$, tedy jeden člen v součinu $X_i \cdot \mathbf{1}_{[X_i=0]}$ je vždy nulový.

Z předchozích výpočtů plyne, že náhodný vektor $\begin{pmatrix} X_i \\ \mathbf{1}_{[X_i=0]} \end{pmatrix}$, $i = 1, \dots, n$, má vektor středních hodnot $\begin{pmatrix} \lambda \\ e^{-\lambda} \end{pmatrix}$ a varianční matici $\begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix}$.
Dále víme, že

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \lambda) = \sqrt{n}(\bar{X}_n - \lambda) \quad \text{a} \quad \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{1}_{[X_i=0]} - e^{-\lambda}) = \sqrt{n}(\hat{p}_0 - e^{-\lambda}).$$

Potom z centrální limitní věty dostáváme dokazovaný vztah. □

Nyní již můžeme formulovat větu, která říká, jak vypadá za (1.1) asymptotické rozdělení obecného Poissonova nulového indexu.

Věta 6. *Nechť X_1, X_2, \dots, X_n je náhodný výběr z rozdělení $\text{Po}(\lambda)$ pro nějaké $\lambda > 0$. Nechť dále $\hat{I}_n := f(\bar{X}_n, \hat{p}_0)$ je Poissonův nulový index a f má spojitě parciální derivace na okolí bodu $(\lambda, e^{-\lambda})^\top$. Potom platí*

$$\sqrt{n} \cdot \hat{I}_n \xrightarrow{D} \mathbf{N}\left(0, \mathbb{D}f(\lambda, e^{-\lambda}) \Sigma (\mathbb{D}f(\lambda, e^{-\lambda}))^\top\right) \text{ pro } n \rightarrow \infty,$$

kde

$$\Sigma = \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix} \quad \text{a} \quad \mathbb{D}f(x, y) = \left(\frac{\partial f(x, y)}{\partial x}, \frac{\partial f(x, y)}{\partial y} \right).$$

Důkaz: Tvrzení věty je důsledkem věty 3 (Δ -věta), předchozího tvrzení a toho, že za daných předpokladů

$$f(\lambda, e^{-\lambda}) = 0$$

díky (i) v definici Poissonova nulového indexu. □

Poznámka. Rozptýl normálního rozdělení v předchozí větě je zapsán ve tvaru lineární formy. Pomocí úpravy

$$\begin{pmatrix} A & B \end{pmatrix} \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix} \begin{pmatrix} A \\ B \end{pmatrix} = \begin{pmatrix} A & B \end{pmatrix} \begin{pmatrix} \lambda A - \lambda e^{-\lambda} B \\ -\lambda e^{-\lambda} A + e^{-\lambda}(1 - e^{-\lambda}) B \end{pmatrix} = \lambda A^2 - 2\lambda e^{-\lambda} AB + e^{-\lambda}(1 - e^{-\lambda}) B^2$$

tak získáme i jeho explicitní vyjádření. Dostaneme

$$\mathbb{D}f(\lambda, e^{-\lambda}) \Sigma (\mathbb{D}f(\lambda, e^{-\lambda}))^\top = \lambda A^2 - 2\lambda e^{-\lambda} AB + e^{-\lambda}(1 - e^{-\lambda}) B^2,$$

kde

$$A = \left. \frac{\partial f(x, y)}{\partial x} \right|_{\substack{x=\lambda \\ y=e^{-\lambda}}} \quad \text{a} \quad B = \left. \frac{\partial f(x, y)}{\partial y} \right|_{\substack{x=\lambda \\ y=e^{-\lambda}}}.$$

2.2 Příklady Poissonových nulových indexů

Uvedeme nyní několik způsobů, které vedou ke konstrukci konkrétních Poissonových nulových indexů. Nejprve se zaměříme na dva nulové indexy, které byly zavedeny ve výchozím článku Weiß a kol. (2019). Pro srovnání poté zmíníme i některé další varianty.

Vyjdeme ze vztahu

$$p_0 = e^{-\lambda}, \quad (2.1)$$

který platí pro náhodnou veličinu X s rozdělením $\text{Po}(\lambda)$. Podle definice Poissonova nulového indexu potřebujeme najít funkci f , jejíž definiční obor obsahuje množinu $(0, \infty) \times (0, 1)$ a která pro libovolné $\lambda > 0$ splňuje

$$f(\lambda, p_0) = 0 \quad \Leftrightarrow \quad p_0 = e^{-\lambda}.$$

Chceme tedy provést vhodné ekvivalentní úpravy (2.1), tak aby na pravé straně zůstala nula.

Pokud jsou obě strany (2.1) kladné, můžeme je zlogaritmovat. Potom dostaneme

$$\ln(p_0) = -\lambda.$$

Nyní již stačí vydělit výraz na levé i na pravé straně rovnosti konstantou $\lambda > 0$, a poté k oběma částem přičíst číslo jedna.

$$\begin{aligned} \frac{\ln(p_0)}{\lambda} &= -1 && / + 1 \\ \frac{\ln(p_0)}{\lambda} + 1 &= 0 \end{aligned}$$

Označme

$$I_1 := f_1(\lambda, p_0),$$

kde f_1 je zobrazení definované jako

$$\begin{aligned} f_1 : (0, \infty) \times (0, \infty) &\rightarrow \mathbb{R}, \\ (x, y) &\mapsto 1 + \frac{\ln(y)}{x}. \end{aligned}$$

Potom pro libovolné $\lambda > 0$ a $p_0 > 0$ platí

$$f_1(\lambda, p_0) = 0 \quad \Leftrightarrow \quad \frac{\ln(p_0)}{\lambda} = -1 \quad \Leftrightarrow \quad \ln(p_0) = -\lambda \quad \Leftrightarrow \quad p_0 = e^{-\lambda},$$

tedy f_1 splňuje (i) a (ii) v definici 2. Z toho plyne, že statistika

$$\hat{I}_{1,n} := f_1(\bar{X}_n, \hat{p}_0)$$

je Poissonův nulový index.

V praxi se může stát, že pro náš náhodný výběr získáme odhady $\bar{X}_n = 0$, nebo $\hat{p}_0 = 0$. První možnost nastane právě tehdy, když $X_1 = X_2 = \dots = X_n = 0$, a tedy bude zároveň platit $\hat{p}_0 = 1$. Druhou rovnost naopak dostaneme v situaci, kdy ani

jedno pozorování není nulové. Pro tyto případy je vhodné „spojitě“ dodefinovat $\hat{I}_{1,n}$, například jako

$$\hat{I}_{1,n} = \begin{cases} 1 + \frac{\ln(\hat{p}_0)}{\bar{X}_n} & \text{pro } \bar{X}_n > 0, \hat{p}_0 > 0, \\ 0 & \text{pro } \bar{X}_n = 0, \\ -\infty & \text{pro } \hat{p}_0 = 0, \end{cases}$$

kde v části pro $\bar{X}_n = 0$ využíváme, že v tomto případě $\hat{p}_0 = 1 = e^0 = e^{-\bar{X}_n}$. Zároveň platí $\lim_{x \rightarrow 0} (1 + \frac{\ln(e^{-x})}{x}) = \lim_{x \rightarrow 0} (1 + \frac{-x}{x}) = 0$.

Je ovšem dobré si uvědomit, že při dostatečném rozsahu výběru obvykle k žádné z těchto „extrémních“ situací nedojde. Protože X_1, X_2, \dots, X_n jsou nezávislé a stejně rozdělené, snadno spočteme, že

$$P(\bar{X}_n = 0) = P(X_1 = 0, \dots, X_n = 0) = (P(X_1 = 0))^n,$$

$$P(\hat{p}_0 = 0) = P(X_1 \neq 0, \dots, X_n \neq 0) = (P(X_1 \neq 0))^n = (1 - P(X_1 = 0))^n.$$

Pokud je tedy $P(X_1 = 0) \in (0, 1)$, obě dvě pravděpodobnosti výše konvergují pro $n \rightarrow \infty$ k nule.

Další možnost je vynásobit obě strany (2.1) nenulovou hodnotou e^λ a následně odečíst číslo jedna.

$$\begin{aligned} p_0 &= e^{-\lambda} && / \cdot e^\lambda \\ p_0 e^\lambda &= 1 && / - 1 \\ p_0 e^\lambda - 1 &= 0 \end{aligned}$$

Položíme $I_2 := f_2(\lambda, p_0)$, kde

$$\begin{aligned} f_2 : \mathbb{R}^2 &\rightarrow \mathbb{R}, \\ (x, y) &\mapsto ye^x - 1. \end{aligned}$$

Dostaneme opět

$$f_2(\lambda, p_0) = 0 \Leftrightarrow p_0 e^\lambda - 1 = 0 \Leftrightarrow p_0 e^\lambda = 1 \Leftrightarrow p_0 = e^{-\lambda},$$

tedy

$$\hat{I}_{2,n} := f_2(\bar{X}_n, \hat{p}_0) = \hat{p}_0 e^{\bar{X}_n} - 1$$

je podle definice 2 rovněž Poissonův nulový index.

Můžeme si všimnout, že zde, na rozdíl od předchozího příkladu, je zobrazení f_2 dobře definované pro libovolné hodnoty \bar{X}_n a \hat{p}_0 .

V předchozím případě jsme (2.1) dělili číslem $e^{-\lambda}$. Alternativně bychom mohli tuto rovnost vydělit hodnotou p_0 za předpokladu, že je kladná.

$$\begin{aligned} p_0 &= e^{-\lambda} && / : p_0 \\ 1 &= \frac{e^{-\lambda}}{p_0} && / - 1 \\ 0 &= \frac{e^{-\lambda}}{p_0} - 1 \end{aligned}$$

Pro tyto účely definujeme funkci

$$\begin{aligned} f_3 : \mathbb{R} \times (0, \infty) &\rightarrow \mathbb{R}, \\ (x, y) &\mapsto \frac{e^{-x}}{y} - 1. \end{aligned}$$

Znovu obdržíme vztah

$$I_3 := f_3(\lambda, p_0) = 0 \Leftrightarrow 1 = \frac{e^{-\lambda}}{p_0} \Leftrightarrow p_0 = e^{-\lambda}.$$

Získáme tak třetí nulový index $\hat{I}_{3,n} := f_3(\bar{X}_n, \hat{p}_0)$, který podobně jako $\hat{I}_{1,n}$ ošetříme i pro případ $\hat{p}_0 = 0$:

$$\hat{I}_{3,n} = \begin{cases} \frac{e^{-\bar{X}_n}}{\hat{p}_0} - 1 & \text{pro } \hat{p}_0 > 0, \\ \infty & \text{pro } \hat{p}_0 = 0. \end{cases}$$

Jiný způsob konstrukce nulového indexu opět využívá rovnost, kterou jsme získali aplikací přirozeného logaritmu na obě strany (2.1). Pro $p_0 \in (0, 1)$ bude logaritmus z p_0 různý od nuly, a budeme jím tedy moci celou rovnost vydělit.

$$\begin{aligned} \ln(p_0) &= -\lambda && / : \ln(p_0) \\ 1 &= \frac{-\lambda}{\ln(p_0)} && / + \frac{\lambda}{\ln(p_0)} \\ 1 + \frac{\lambda}{\ln(p_0)} &= 0 \end{aligned}$$

V tomto případě budeme uvažovat funkci

$$\begin{aligned} f_4 : \mathbb{R} \times (0, 1) &\rightarrow \mathbb{R}, \\ (x, y) &\mapsto 1 + \frac{x}{\ln(y)}. \end{aligned}$$

Dostaneme opět

$$I_4 := f_4(\lambda, p_0) = 0 \Leftrightarrow p_0 = e^{-\lambda},$$

což nám dává další Poissonův nulový index $\hat{I}_{4,n} := f_4(\bar{X}_n, \hat{p}_0)$. Stejně jako u předchozích indexů ho dodefinujeme i pro hodnoty na hranici definičního oboru.

$$\hat{I}_{4,n} = \begin{cases} 1 + \frac{\bar{X}_n}{\ln(\hat{p}_0)} & \text{pro } \hat{p}_0 \in (0, 1), \\ 1 & \text{pro } \hat{p}_0 = 0, \\ 0 & \text{pro } \hat{p}_0 = 1. \end{cases}$$

Tvrzení 4 v sekci 2.1 nám říká, že $\hat{I}_{k,n} = f_j(\bar{X}_n, \hat{p}_0)$ jsou konzistentními odhady $f_k(\mathbf{E}X_1, p_0)$, za předpokladu, že f_k jsou spojité v bodě $(\mathbf{E}X_1, p_0)^\top$, pro $k \in \{1, 2, 3, 4\}$. Protože zmíněné funkce jsou spojité na celých svých definičních oborech, stačí požadovat, aby $(\mathbf{E}X_1, p_0)^\top$ bylo prvkem jejich definičního oboru. Podle 2. bodu tvrzení 4 navíc $\hat{I}_{k,n}$; $k = 1, 2, 3, 4$; za platnosti hypotézy konvergují v pravděpodobnosti k nule.

Tvrzení 7. *Nechť X_1, X_2, \dots, X_n je náhodný výběr z libovolného diskrétního rozdělení se střední hodnotou $\mathbf{E}X_1 \in (0, \infty)$. Pokud $p_0 = \mathbf{P}(X_1 = 0) \in (0, 1)$, potom*

$$\begin{aligned} \hat{\mathbf{I}}_{1,n} &\xrightarrow[n \rightarrow \infty]{\mathbf{P}} 1 + \frac{\ln(p_0)}{\mathbf{E}X_1}, & \hat{\mathbf{I}}_{2,n} &\xrightarrow[n \rightarrow \infty]{\mathbf{P}} p_0 e^{\mathbf{E}X_1} - 1, \\ \hat{\mathbf{I}}_{3,n} &\xrightarrow[n \rightarrow \infty]{\mathbf{P}} \frac{e^{-\mathbf{E}X_1}}{p_0} - 1, & \hat{\mathbf{I}}_{4,n} &\xrightarrow[n \rightarrow \infty]{\mathbf{P}} 1 + \frac{\mathbf{E}X_1}{\ln(p_0)}, \end{aligned}$$

kde $\hat{\mathbf{I}}_{k,n}$, $k = 1, 2, 3, 4$, jsou Poissonovy nulové indexy definované výše. Speciálně, pokud $X_1 \sim \text{Po}(\lambda)$ pro nějaké $\lambda > 0$, platí

$$\hat{\mathbf{I}}_{k,n} \xrightarrow[n \rightarrow \infty]{\mathbf{P}} 0, \quad k = 1, 2, 3, 4.$$

Pro další práci s nulovými indexy se hodí znát jejich asymptotické rozdělení za platnosti hypotézy (1.1). V následující větě se proto budeme zabývat jeho odvozením.

Věta 8. *Uvažujme náhodný výběr X_1, X_2, \dots, X_n z rozdělení $\text{Po}(\lambda)$ pro nějaké $\lambda > 0$. Potom*

1.

$$\sqrt{n} \cdot \hat{\mathbf{I}}_{1,n} \xrightarrow[n \rightarrow \infty]{\mathbf{D}} \mathbf{N}\left(0, \frac{e^\lambda - \lambda - 1}{\lambda^2}\right),$$

2.

$$\sqrt{n} \cdot \hat{\mathbf{I}}_{2,n} \xrightarrow[n \rightarrow \infty]{\mathbf{D}} \mathbf{N}(0, e^\lambda - \lambda - 1),$$

3.

$$\sqrt{n} \cdot \hat{\mathbf{I}}_{3,n} \xrightarrow[n \rightarrow \infty]{\mathbf{D}} \mathbf{N}(0, e^\lambda - \lambda - 1),$$

4.

$$\sqrt{n} \cdot \hat{\mathbf{I}}_{4,n} \xrightarrow[n \rightarrow \infty]{\mathbf{D}} \mathbf{N}\left(0, \frac{e^\lambda - \lambda - 1}{\lambda^2}\right).$$

Důkaz:

1.

Víme, že $\hat{\mathbf{I}}_{1,n} = f_1(\bar{X}_n, \hat{p}_0)$ je nulový index. Funkce

$$f_1 : (x, y) \mapsto 1 + \frac{\ln(y)}{x}$$

je spojitě diferencovatelná na množině $(0, \infty) \times (0, \infty)$, a tedy i v okolí bodu $\boldsymbol{\mu} := (\lambda, e^{-\lambda})^\top$. Potom podle věty 6 platí pro $n \rightarrow \infty$

$$\sqrt{n} \cdot \hat{\mathbf{I}}_{1,n} \xrightarrow{\mathbf{D}} \mathbf{N}\left(0, \mathbb{D}f_1(\boldsymbol{\mu})\boldsymbol{\Sigma}(\mathbb{D}f_1(\boldsymbol{\mu})^\top)\right), \quad (2.2)$$

kde

$$\boldsymbol{\Sigma} = \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix} \quad \text{a} \quad \mathbb{D}f_1(x, y) = \left(-\frac{\ln(y)}{x^2}, \frac{1}{xy}\right).$$

Nyní již stačí dosadit konkrétní hodnoty.

$$\mathbb{D}f_1(\boldsymbol{\mu}) = \left(-\frac{\ln(e^{-\lambda})}{\lambda^2}, \frac{1}{\lambda e^{-\lambda}}\right) = \left(\frac{1}{\lambda}, \frac{1}{\lambda e^{-\lambda}}\right)$$

Dále platí

$$\begin{aligned} \mathbb{D}f_1(\boldsymbol{\mu})\boldsymbol{\Sigma}(\mathbb{D}f_1(\boldsymbol{\mu}))^\top &= \begin{pmatrix} \frac{1}{\lambda} & \frac{1}{\lambda e^{-\lambda}} \end{pmatrix} \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix} \begin{pmatrix} \lambda^{-1} \\ \frac{1}{\lambda e^{-\lambda}} \end{pmatrix} = \\ \begin{pmatrix} \frac{1}{\lambda} & \frac{1}{\lambda e^{-\lambda}} \end{pmatrix} \begin{pmatrix} 0 \\ -e^{-\lambda} + \frac{1 - e^{-\lambda}}{\lambda} \end{pmatrix} &= \frac{1 - e^{-\lambda}}{\lambda^2 e^{-\lambda}} - \frac{e^{-\lambda}}{\lambda e^{-\lambda}} = \frac{1 - e^{-\lambda} - \lambda e^{-\lambda}}{\lambda^2 e^{-\lambda}} = \frac{e^\lambda - 1 - \lambda}{\lambda^2}. \end{aligned}$$

Dosazením do (2.2) dostaneme $\sqrt{n} \cdot \hat{\mathbf{I}}_{1,n} \xrightarrow[n \rightarrow \infty]{\text{D}} \mathbf{N}(0, \frac{e^\lambda - \lambda - 1}{\lambda^2})$, což jsme chtěli dokázat.

2.

Postupujeme analogicky jako v předchozí části, pouze uvažujeme zobrazení f_2 . Budeme se tedy držet značení $\boldsymbol{\mu}$ a $\boldsymbol{\Sigma}$.

Funkce

$$f_2 : (x, y) \mapsto e^x y - 1$$

má spojité parciální derivace na $\mathbb{R} \times \mathbb{R}$, a tedy i v okolí bodu $\boldsymbol{\mu}$. Platí

$$\mathbb{D}f_2(x, y) = (e^x y, e^x).$$

Podle věty 6 tedy pro $n \rightarrow \infty$ máme

$$\sqrt{n} \cdot \hat{\mathbf{I}}_{2,n} \xrightarrow{\text{D}} \mathbf{N}(0, \mathbb{D}f_2(\boldsymbol{\mu})\boldsymbol{\Sigma}(\mathbb{D}f_2(\boldsymbol{\mu}))^\top). \quad (2.3)$$

Dále

$$\mathbb{D}f_2(\boldsymbol{\mu}) = (e^\lambda \cdot e^{-\lambda}, e^\lambda) = (1, e^\lambda).$$

Díky explicitnímu vyjádření $\mathbb{D}f_2(\boldsymbol{\mu})\boldsymbol{\Sigma}(\mathbb{D}f_2(\boldsymbol{\mu}))^\top$ (viz poznámka pod větou 6) tak stačí spočítat

$$\begin{aligned} \lambda \cdot 1^2 - 2\lambda e^{-\lambda} \cdot 1 \cdot e^\lambda + e^{-\lambda}(1 - e^{-\lambda})(e^\lambda)^2 &= \lambda - 2\lambda + e^{-\lambda}(1 - e^{-\lambda})e^{2\lambda} = \\ -\lambda + e^\lambda(1 - e^{-\lambda}) &= e^\lambda - \lambda - 1. \end{aligned}$$

Potom dosazením do (2.3) dostaneme $\sqrt{n}(\hat{\mathbf{I}}_{2,n} - 0) \xrightarrow[n \rightarrow \infty]{\text{D}} \mathbf{N}(0, e^\lambda - \lambda - 1)$.

3.

Zobrazení

$$f_3 : (x, y) \mapsto \frac{e^{-x}}{y} - 1$$

je spojitě diferencovatelné na intervalu $\mathbb{R} \times (0, \infty)$, a tedy opět i v okolí bodu $\boldsymbol{\mu}$.

Dále

$$\mathbb{D}f_3(x, y) = \left(-\frac{e^{-x}}{y}, -\frac{e^{-x}}{y^2} \right) \quad \text{a} \quad \mathbb{D}f_3(\boldsymbol{\mu}) = (-1, -e^\lambda).$$

Všimněme si, že

$$\mathbb{D}f_3(\boldsymbol{\mu}) = -\mathbb{D}f_2(\boldsymbol{\mu}),$$

tudíž rozptyl asymptotického rozdělení bude stejný jako v předchozím bodě.

4.

Znovu vidíme, že funkce

$$f_4 : (x, y) \mapsto 1 + \frac{x}{\ln(y)}$$

má spojité parciální derivace na $\mathbb{R} \times (0, 1)$, a tedy zřejmě i na okolí bodu $\boldsymbol{\mu}$. Dostaneme

$$\mathbb{D}f_4(x, y) = \left(\frac{1}{\ln(y)}, -\frac{x}{y \ln^2(y)} \right)$$

Potom

$$\mathbb{D}f_4(\boldsymbol{\mu}) = \left(\frac{1}{\ln(e^{-\lambda})}, -\frac{\lambda}{e^{-\lambda} \ln^2(e^{-\lambda})} \right) = \left(-\frac{1}{\lambda}, -\frac{1}{e^{-\lambda} \lambda} \right),$$

což nám dá

$$\begin{aligned} \mathbb{D}f_4(\boldsymbol{\mu}) \boldsymbol{\Sigma} (\mathbb{D}f_4(\boldsymbol{\mu}))^\top &= \lambda \left(-\frac{1}{\lambda} \right)^2 - 2\lambda e^{-\lambda} \left(-\frac{1}{\lambda} \right) \left(-\frac{1}{e^{-\lambda} \lambda} \right) + e^{-\lambda} (1 - e^{-\lambda}) \left(-\frac{1}{e^{-\lambda} \lambda} \right)^2 \\ &= \frac{\lambda}{\lambda^2} - \frac{2\lambda e^{-\lambda}}{e^{-\lambda} \lambda^2} + e^{-\lambda} (1 - e^{-\lambda}) \cdot \frac{1}{e^{-2\lambda} \lambda^2} = \frac{\lambda}{\lambda^2} - \frac{2\lambda}{\lambda^2} + e^{-\lambda} (1 - e^{-\lambda}) \cdot \frac{e^{2\lambda}}{\lambda^2} \\ &= \frac{\lambda - 2\lambda + e^\lambda (1 - e^{-\lambda})}{\lambda^2} = \frac{e^\lambda - \lambda - 1}{\lambda^2}. \end{aligned}$$

□

Věta 8 nám říká, že za nulové hypotézy (1.1) a pro velké hodnoty n mají veličiny $\sqrt{n} \cdot \hat{I}_{k,n}$ (pro $k \in \{1, 2, 3, 4\}$) rozdělení, které je blízké normálnímu rozdělení s nulovou střední hodnotou a danými rozptyly.

Můžeme si ještě všimnout, že pro libovolné $0 < \lambda < 1$ bude $\frac{1}{\lambda^2} < 1$, tudíž

$$\text{avar}(\hat{I}_{1,n}) = \text{avar}(\hat{I}_{4,n}) > \text{avar}(\hat{I}_{2,n}) = \text{avar}(\hat{I}_{3,n}),$$

kde $\text{avar}(\hat{I}_{k,n})$ značí asymptotický rozptyl nulového indexu $\hat{I}_{k,n}$ za platnosti (1.1). Naopak je-li $\lambda > 1$, platí $\frac{1}{\lambda^2} > 1$, a tedy asymptotické rozptyly nulových indexů splňují nerovnost

$$\text{avar}(\hat{I}_{1,n}) = \text{avar}(\hat{I}_{4,n}) < \text{avar}(\hat{I}_{2,n}) = \text{avar}(\hat{I}_{3,n}).$$

Speciálně, pro $\lambda = 1$ budou všechny čtyři asymptotické rozptyly stejně velké.

2.3 Testování pomocí nulových indexů

V předchozích částech této kapitoly jsme zavedli pojem Poissonova nulového indexu a odvodili jeho asymptotické rozdělení, přičemž jsme se podívali i na čtyři konkrétní případy. Nyní se zaměříme na to, jak tyto poznatky využít k našemu stanovenému cíli – k testování nulové hypotézy (1.1).

Z předešlých dvou sekcí víme, že asymptotické rozdělení nulového indexu (za platnosti hypotézy) obecně závisí na hodnotě λ , která je pro nás neznámá. Abychom mohli testovat (1.1), bude tedy nutné upravit statistiku \hat{I}_n tak, aby její rozdělení již na neznámých parametrech nezáviselo.

Připomeňme, že věta 6 nám říká, že za určitých předpokladů má nulový index za platnosti hypotézy (1.1) asymptotické rozdělení (používáme-li značení zavedené ve zmíněné větě)

$$\sqrt{n} \cdot \hat{I}_n \xrightarrow[n \rightarrow \infty]{D} \mathbf{N}\left(0, \mathbb{D}f(\boldsymbol{\mu}) \boldsymbol{\Sigma} (\mathbb{D}f(\boldsymbol{\mu}))^\top\right).$$

Jednou z vlastností normálního rozdělení je, že

$$Z \sim \mathbf{N}(\boldsymbol{\mu}, \sigma^2) \implies c \cdot Z \sim \mathbf{N}(c\boldsymbol{\mu}, c^2\sigma^2),$$

kde $c \in \mathbb{R}$ je libovolná konstanta a Z je nějaká náhodná veličina. Díky Cramér-Sluckého větě, jejíž znění lze nalézt například v knize van der Vaart (1998, Věta 2.8), tento vztah platí i v případě, kdy symbol \sim nahradíme konvergencí v distribuci. Je-li $\mathbb{D}f(\boldsymbol{\mu})\boldsymbol{\Sigma}(\mathbb{D}f(\boldsymbol{\mu}))^\top > 0$, můžeme položit

$$c := \left(\mathbb{D}f(\boldsymbol{\mu})\boldsymbol{\Sigma}(\mathbb{D}f(\boldsymbol{\mu}))^\top \right)^{-\frac{1}{2}} > 0,$$

a dostaneme tak

$$\sqrt{n} \cdot \left(\mathbb{D}f(\boldsymbol{\mu})\boldsymbol{\Sigma}(\mathbb{D}f(\boldsymbol{\mu}))^\top \right)^{-\frac{1}{2}} \hat{\mathbf{I}}_n \xrightarrow[n \rightarrow \infty]{\text{D}} \mathbf{N}(0, 1).$$

Vidíme, že za (1.1) konverguje výše uvedený výraz v distribuci k náhodné veličině s normovaným normálním rozdělením, tedy rozdělením nezávislejícím na neznámých parametrech. Stále však nulový index $\hat{\mathbf{I}}_n$ dělíme odmocninou z jeho asymptotického rozptylu, kde se objevuje neznámá konstanta λ . K získání testové statistiky budeme tedy muset tento rozptyl nějak odhadnout.

Z věty 6 víme, že za nulové hypotézy závisí rozptyl asymptotického rozdělení $\hat{\mathbf{I}}_n$ na λ a $e^{-\lambda}$, tedy celkově jde tento rozptyl vyjádřit jako $h(\lambda)$, kde

$$h : (0, \infty) \rightarrow (0, \infty),$$

a λ je parametr daného Poissonova rozdělení. Jednou z možností je tedy nahradit neznámou proměnnou λ výběrovým průměrem \bar{X}_n . Platí totiž

$$\bar{X}_n \xrightarrow[n \rightarrow \infty]{\text{P}} \mathbf{E}X_1 \stackrel{(1.1)}{=} \lambda.$$

Pokud je funkce h spojitá v bodě $\mathbf{E}X_1$, dostaneme z věty o spojitě transformaci

$$h(\bar{X}_n) \xrightarrow[n \rightarrow \infty]{\text{P}} h(\mathbf{E}X_1) \stackrel{(1.1)}{=} h(\lambda),$$

a $h(\bar{X}_n)$ je tak konzistentní odhad $h(\lambda)$.

Nyní již můžeme zformulovat klíčovou větu, která ukazuje, jakým způsobem lze pomocí nulových indexů testovat (1.1). Pro připomenutí je dobré zmínit, že pro dané $\alpha \in (0, 1)$ budeme symbolem $u_{1-\frac{\alpha}{2}}$ značit $(1 - \frac{\alpha}{2})$ -tý kvantil normovaného normálního rozdělení. Platí $u_{1-\frac{\alpha}{2}} = \Phi^{-1}(1 - \frac{\alpha}{2})$, kde Φ je distribuční funkce $\mathbf{N}(0, 1)$,

Věta 9. *Nechť X_1, X_2, \dots, X_n je náhodný výběr z nějakého diskrétního rozdělení, jehož nosič je podmnožina \mathbb{N}_0 . Nechť $\alpha \in (0, 1)$ je předem stanovené číslo. Uvažujme Poissonův nulový index*

$$\hat{\mathbf{I}}_n := f(\bar{X}_n, \hat{p}_0),$$

pro něž za nulové hypotézy (1.1) platí

$$\sqrt{n} \cdot \hat{\mathbf{I}}_n \xrightarrow[n \rightarrow \infty]{\text{D}} \mathbf{N}\left(0, h(\mathbf{E}X_1)\right),$$

kde h je nějaká funkce, $h : (0, \infty) \rightarrow (0, \infty)$. Necht \hat{h}_n je konzistentní odhad $h(\mathbf{E}X_1)$, který nezávisí na neznámých parametrech.

Potom za platnosti hypotézy (1.1) platí

$$T_n := \sqrt{n} \cdot \frac{\hat{I}_n}{\sqrt{\hat{h}_n}} \xrightarrow[n \rightarrow \infty]{D} \mathbf{N}(0, 1).$$

Položíme-li

$$C(\alpha) := \left(-\infty, -u_{1-\frac{\alpha}{2}}\right) \cup \left(u_{1-\frac{\alpha}{2}}, \infty\right),$$

potom test (1.1) řídící se rozhodovacím pravidlem:

$$\text{zamítneme } (H_0) \iff T_n \in C(\alpha)$$

má asymptotickou hladinou významnosti α .

Důkaz: T_n lze rozložit na součin dvou výrazů jako

$$T_n = \sqrt{n} \cdot \frac{\hat{I}_n}{\sqrt{\hat{h}_n}} = \sqrt{n} \frac{\hat{I}_n}{\sqrt{h(\mathbf{E}X_1)}} \cdot \frac{\sqrt{h(\mathbf{E}X_1)}}{\sqrt{\hat{h}_n}},$$

přičemž z vlastností normálního rozdělení diskutovaných před touto větou máme

$$\sqrt{n} \frac{\hat{I}_n}{\sqrt{h(\mathbf{E}X_1)}} \xrightarrow[n \rightarrow \infty]{D} \mathbf{N}(0, 1).$$

Dále

$$\frac{\sqrt{h(\mathbf{E}X_1)}}{\sqrt{\hat{h}_n}} \xrightarrow[n \rightarrow \infty]{P} 1,$$

což plyne z konzistence \hat{h}_n a věty o spojitě transformaci (druhá odmocnina je spojitá funkce). Potom podle Cramérový-Sluckého věty

$$T_n \xrightarrow[n \rightarrow \infty]{D} 1 \cdot Z = Z, \quad \text{kde } Z \sim \mathbf{N}(0, 1).$$

K důkazu zbytku věty si nejdříve uvědomme, že

$$\mathbf{P}(|T_n| < u_{1-\frac{\alpha}{2}}) = \mathbf{P}(-u_{1-\frac{\alpha}{2}} < T_n < u_{1-\frac{\alpha}{2}}) = \mathbf{P}(T_n < u_{1-\frac{\alpha}{2}}) - \mathbf{P}(T_n \leq -u_{1-\frac{\alpha}{2}}).$$

Z definice konvergence v distribuci tedy máme (za platnosti hypotézy)

$$\lim_{n \rightarrow \infty} \mathbf{P}(T_n < u_{1-\frac{\alpha}{2}}) = \Phi(u_{1-\frac{\alpha}{2}}^-) = \Phi(u_{1-\frac{\alpha}{2}}),$$

kde poslední rovnost plyne ze spojitosti funkce Φ na celé reálné ose, a podobně

$$\lim_{n \rightarrow \infty} \mathbf{P}(T_n \leq -u_{1-\frac{\alpha}{2}}) = \Phi(-u_{1-\frac{\alpha}{2}}).$$

Díky tomu, že normované normální rozdělení je symetrické kolem nuly, dostaneme

$$\Phi(-u_{1-\frac{\alpha}{2}}) = 1 - \Phi(u_{1-\frac{\alpha}{2}}).$$

Potom za platnosti (1.1)

$$\begin{aligned}\lim_{n \rightarrow \infty} \mathbb{P}\left(|T_n| < u_{1-\frac{\alpha}{2}}\right) &= \Phi(u_{1-\frac{\alpha}{2}}) - \Phi(-u_{1-\frac{\alpha}{2}}) = \Phi(u_{1-\frac{\alpha}{2}}) - (1 - \Phi(u_{1-\frac{\alpha}{2}})) \\ &= 1 - \frac{\alpha}{2} - \left(1 - \left(1 - \frac{\alpha}{2}\right)\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha.\end{aligned}$$

Za (H_0) tedy

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(T_n \in C(\alpha)\right) = \lim_{n \rightarrow \infty} \mathbb{P}\left(|T_n| \geq u_{1-\frac{\alpha}{2}}\right) = 1 - \lim_{n \rightarrow \infty} \mathbb{P}\left(|T_n| < u_{1-\frac{\alpha}{2}}\right) = \alpha,$$

což odpovídá tomu, že test $(T_n, C(\alpha))$ má hladinu α asymptoticky. \square

Zde je dobré připomenout, že Poissonův nulový index \hat{I}_n , který se vyskytuje v čitateli T_n , je za (1.1) blízky nule (tvrzení 4). Ve prospěch alternativy tak svědčí vysoké záporné a kladné hodnoty statistiky T_n .

V souvislosti s testováním je rovněž dobré zmínit p-hodnotu, která je definovaná jako

$$p = \inf\{\alpha \in (0, 1) : t \in C(\alpha)\},$$

kde t značí realizovanou hodnotu testové statistiky T_n . Dosazením kritických hodnot a drobnou úpravou získáme, že p-hodnotu výše uvedeného testu jde vyjádřit ve tvaru

$$\begin{aligned}p &= \inf\{\alpha \in (0, 1) : |t| \geq u_{1-\frac{\alpha}{2}}\} = \inf\{\alpha \in (0, 1) : |t| \geq \Phi^{-1}(1 - \frac{\alpha}{2})\} = \\ &= \inf\{\alpha \in (0, 1) : \Phi(|t|) \geq 1 - \frac{\alpha}{2}\} = \inf\{\alpha \in (0, 1) : \alpha \geq 2(1 - \Phi(|t|))\} = \\ &= 2(1 - \Phi(|t|)).\end{aligned}$$

U třetí rovnosti zde využíváme, že distribuční funkce Φ je neklesající.

P-hodnota podle své definice vyjadřuje nejmenší možnou hladinu, na které bychom zamítali nulovou hypotézu. Pokud tedy pro naše data vyjde p-hodnota, která je menší než předem stanovená hladina $\alpha \in (0, 1)$, budeme moci zamítnout (1.1). Alternativně tak lze při testování hypotézy rozhodovat o jejím zamítnutí pomocí pravidla

$$\text{zamítneme } (H_0) \iff p \leq \alpha.$$

Analogii věty 9 můžeme podobně jako jiné obecné věty zformulovat ještě pro nulové indexy zavedené v sekci 2.2.

Věta 10. *Nechť X_1, X_2, \dots, X_n je náhodný výběr z nějakého diskrétního rozdělení, jehož nosič je podmnožina \mathbb{N}_0 , a $\alpha \in (0, 1)$ je předem stanovené číslo.*

Označme

$$\begin{aligned}T_{1,n} &:= \frac{\hat{I}_{1,n}}{\sqrt{\frac{e^{\bar{X}_n} - \bar{X}_n - 1}{n(\bar{X}_n)^2}}}, & T_{2,n} &:= \sqrt{n} \cdot \frac{\hat{I}_{2,n}}{\sqrt{e^{\bar{X}_n} - \bar{X}_n - 1}}, \\ T_{3,n} &:= \sqrt{n} \cdot \frac{\hat{I}_{3,n}}{\sqrt{e^{\bar{X}_n} - \bar{X}_n - 1}}, & T_{4,n} &:= \frac{\hat{I}_{4,n}}{\sqrt{\frac{e^{\bar{X}_n} - \bar{X}_n - 1}{n(\bar{X}_n)^2}}}.\end{aligned}$$

Potom za platnosti nulové hypotézy (1.1) platí pro $k \in \{1, 2, 3, 4\}$

$$T_{k,n} \xrightarrow[n \rightarrow \infty]{D} N(0, 1).$$

Dále statistický test (1.1) řídící se pravidlem:

$$\text{zamítneme } (H_0) \iff |T_{k,n}| \geq u_{1-\frac{\alpha}{2}}$$

má asymptotickou hladinu významnosti α , $k = 1, 2, 3, 4$.

2.4 Úskalí testování založeného na nulových indexech

Testy dobré shody popsané v této kapitole jsou takzvané asymptotické testy, což znamená, že dosahují požadované hladiny významnosti α pouze asymptoticky. Skutečná hladina při daném rozsahu výběru tak může být výrazně vyšší.

Z tohoto hlediska by variantou na možné vylepšení mohlo být nahrazení kvantilu normovaného normálního rozdělení kvantilem Studentova t-rozdělení s $n - 1$ stupni volnosti, kde n je rozsah výběru. Pro zamítání nulové hypotézy bychom potom používali pravidlo:

$$\text{zamítneme } (H_0) \iff |T_n| \geq t_{n-1}(1 - \frac{\alpha}{2}),$$

kde $t_{n-1}(1 - \frac{\alpha}{2})$ značí $(1 - \frac{\alpha}{2})$ -tý kvantil t-rozdělení s $n - 1$ stupni volnosti. Platí totiž

$$t_{n-1}(1 - \frac{\alpha}{2}) \searrow u_{1-\frac{\alpha}{2}},$$

pro $n \rightarrow \infty$ a libovolné $\alpha \in (0, 1)$. Protože podle předchozího je pro každé přirozené $n \geq 2$

$$t_{n-1}(1 - \frac{\alpha}{2}) \geq u_{1-\frac{\alpha}{2}},$$

zamítali bychom tak nulovou hypotézu pro menší škálu hodnot statistiky T_n .

Výhodou tohoto přístupu je, že daný test by byl konzervativnější než test využívající kvantily $N(0, 1)$, a tedy by měl, obzvláště pro malé hodnoty rozsahu n , menší skutečnou hladinu významnosti a s tím související pravděpodobnost chyby I. druhu (tj. pravděpodobnost zamítnutí platné hypotézy). Naproti tomu, test pracující s kvantily normovaného normálního rozdělení má vyšší sílu, a tedy i menší pravděpodobnost chyby II. druhu (pravděpodobnost nezamítnutí neplatné hypotézy).

Jedna z dalších nevýhod testů dobré shody založených na nulovém indexu je, že testují pouze jeden aspekt pravděpodobnostního rozdělení. Z tohoto důvodu je lze použít pouze k odhalení toho, že dané rozdělení neodpovídá Poissonovu. Naopak nezamítnutí nulové hypotézy rozhodně nestačí k ukázání, že nějaký náhodný výběr pochází z Poissonova rozdělení. Důvodem je, že existují případy (viz příklad níže), kdy $P(X = 0) = e^{-EX}$, ale náhodná veličina X nemá Poissonovo rozdělení.

Příklad. Necht rozdělení náhodné veličiny X je dané vztahy

$$P(X = 0) = \frac{1}{e}, \quad P(X = 1) = 1 - \frac{2}{e}, \quad P(X = 2) = \frac{1}{e}.$$

Potom X má diskrétní rozdělení s nosičem $\{0, 1, 2\}$ a střední hodnotou

$$EX = \sum_{k=0}^2 k \cdot P(X = k) = 1 - \frac{2}{e} + 2 \cdot \frac{1}{e} = 1.$$

Dále platí

$$P(X = 0) = e^{-1} = e^{-EX}.$$

Zároveň však v tomto případě neexistuje konstanta $\lambda > 0$, pro kterou by X mělo rozdělení $\text{Po}(\lambda)$, neboť například

$$P(X = 3) = 0,$$

ovšem pro veličinu s Poissonovým rozdělením by tato pravděpodobnost byla kladná.

Můžeme si všimnout, že v situaci, kdy by rozdělení našeho náhodného výběru X_1, X_2, \dots, X_n odpovídalo rozdělení náhodné veličiny X z předchozího příkladu, by nulová hypotéza (1.1) neplatila (a tedy by platila alternativa (H_1)). Zároveň však pro libovolný Poissonův nulový index

$$\hat{I}_n := f(\bar{X}_n, \hat{p}_0),$$

kde funkce f je spojitá v bodě $(EX_1, p_0)^T$, platí podle druhého bodu tvrzení 4

$$\hat{I}_n \xrightarrow[n \rightarrow \infty]{P} 0,$$

neboť v tomto případě $EX_1 > 0$ a

$$p_0 = e^{-EX_1}.$$

Vzhledem k tomu, že u testů dobré shody založených na nulovém indexu svědčí proti (1.1) hodnoty statistiky \hat{I}_n významně vzdálené od 0, by zmíněné testy v tomto případě nejspíš porušení nulové hypotézy neodhalily.

Problém může nastat rovněž v situaci, kdy skutečné rozdělení náhodného výběru je Poissonovo s vyšší hodnotou parametru λ . Pokud totiž není rozsah našeho výběru dostatečně velký, může se stát, že žádná z napozorovaných hodnot není nulová. Potom získáme odhad

$$\hat{p}_0 = 0.$$

U některých nulových indexů (např. $\hat{I}_{1,n}$) přitom pracujeme s logaritmem této relativní četnosti nebo jí dělíme. V takových případech by hodnota testové statistiky T_n byla $-\infty$ nebo ∞ , a tedy bychom zamítli platnou hypotézu.

V tabulce níže můžeme vidět hodnoty $e^{-\lambda}$ pro různé hodnoty λ . Ty odpovídají pravděpodobnostem $P(X_1 = 0)$, kde $X_1 \sim \text{Po}(\lambda)$. V pravém sloupci je potom pro každé λ uveden průměrný počet pozorování, na které připadne jedna nulová realizace.

λ	$e^{-\lambda}$	e^{λ}
1	0,37	2,72
2	0,14	7,39
3	$4,98 \cdot 10^{-2}$	20,09
4	$1,83 \cdot 10^{-2}$	54,60
5	$6,74 \cdot 10^{-3}$	148,41
6	$2,48 \cdot 10^{-3}$	403,43
7	$9,12 \cdot 10^{-4}$	1 096,63
8	$3,35 \cdot 10^{-4}$	2 980,96
9	$1,23 \cdot 10^{-4}$	8 103,08
10	$4,54 \cdot 10^{-5}$	22 026,47

Lze si například všimnout, že již pro $X_1 \sim \text{Po}(7)$ nabývá náhodná veličina X_1 hodnoty 0 přibližně v jednom z 1000 případů.

Obecně je tedy dobré používat testy založené na nulovém indexu spíše v případech, kdy chceme testovat hypotézu, že náhodný výběr X_1, X_2, \dots, X_n pochází z rozdělení $\text{Po}(\lambda)$ pro nějakou malou hodnotu parametru $\lambda > 0$. V situacích, kdy očekáváme, že konstanta λ bude větší číslo, by mohl být vhodnější jiný typ testu dobré shody (viz kapitola 3).

3. Další varianty testů dobré shody

Testy dobré shody založené na nulovém indexu, kterým jsme se podrobně věnovali v kapitole 2, jsou pouze jedním z mnoha způsobů, jak lze nulovou hypotézu (1.1) testovat. Z tohoto důvodu se hodí zmínit i některé další metody, které se v této situaci také dají použít. Oproti předchozí kapitole však uvedené testy popíšeme daleko stručněji, s tím, že detaily lze nalézt v uvedené literatuře.

3.1 χ^2 -testy dobré shody

Asi nejznámější test pro testování dobré shody s diskretním rozdělením je tzv. Pearsonův χ^2 -test dobré shody, kterému je věnována například velká část kapitoly 12 v knize Anděl (2007). Zde se pokusíme nastínit jeho základní princip.

Uvažujme nejprve situaci, kdy X_1, X_2, \dots, X_n je náhodný výběr z diskretního rozdělení takového, že

$$P(X_1 \in \{1, 2, \dots, K\}) = 1,$$

kde K je nějaké známé přirozené číslo. Zajímá nás test nulové hypotézy

$$P(X_1 = i) = p_i^0, \quad i = 1, \dots, K \tag{3.1}$$

pro nějaké předepsané pravděpodobnosti p_1^0, \dots, p_K^0 splňující $\sum_{j=1}^K p_j^0 = 1$, tj. chceme testovat dobrou shodu s nějakým předem určeným diskretním rozdělením, které má konečný nosič.

Označíme $Y_i := \sum_{j=1}^n \mathbf{1}_{[X_j=i]}$, $i = 1, \dots, K$. Potom testová statistika bude mít tvar

$$\chi^2 = \sum_{j=1}^K \frac{(Y_j - np_j^0)^2}{np_j^0}.$$

Výraz v čitateli zde nějakým způsobem reflektuje rozdíl mezi empirickými četnostmi výskytu hodnot $1, \dots, K$ (ty udávají veličiny Y_i) a jejich očekávanými četnostmi. Za nulové hypotézy by tyto dvě kvantily měly být podobné, tudíž proti (3.1) svědčí vysoké hodnoty testové statistiky.

Dá se ukázat, že za platnosti nulové hypotézy (3.1) má tato statistika asymptoticky χ^2 -rozdělení s $K - 1$ stupni volnosti (viz Anděl (2007), Věta 12.5). Pro získání asymptotického testu s hladinou významnosti α tedy budeme zamítat (3.1), pokud $\chi^2 \geq \chi_{K-1}^2(1-\alpha)$, kde $\chi_{K-1}^2(1-\alpha)$ značí $(1-\alpha)$ -tý kvantil rozdělení χ_{K-1}^2 .

K aplikování zmíněné metody na náš problém (testování (1.1)) brání skutečnost, že nosič Poissonova rozdělení je nekonečná množina. Tuto překážku je ovšem možné částečně odstranit tím, že si hodnoty, kterých veličina s Poissonovým rozdělením nabývá, rozdělíme do konečného počtu tříd.

Konkrétně si nejprve zvolíme konstanty $k \geq 0$ a $r \geq 0$. Dále definujeme veličiny

$$\begin{aligned} Y_1 &:= \sum_{j=1}^n \mathbf{1}_{[X_j \leq r]}, \\ Y_i &:= \sum_{j=1}^n \mathbf{1}_{[X_j = r+i-1]}, \quad i = 2, \dots, k-1, \\ Y_k &:= \sum_{j=1}^n \mathbf{1}_{[X_j \geq r+k-1]}. \end{aligned}$$

Naše pozorování jsme tak rozdělili do k tříd. První z nich je tvořena realizacemi $\leq r$, a poslední realizacemi, které jsou $\geq r+k-1$. Do zbylých $k-2$ jsme zařadili pozorování, která jsou rovna číslům $r+1, r+2, \dots, r+k-2$. Veličiny Y_1, \dots, Y_k budou podle své definice udávat počty realizací v jednotlivých třídách (takzvané empirické četnosti). Dále pravděpodobnost, že náhodná veličina X_1 spadne do i -té třídy ($i = 1, \dots, k$) označíme p_i . Snadno ověříme, že pro $X_1 \sim \text{Po}(\lambda)$ platí

$$p_1 = \sum_{j=0}^r q_j, \quad p_i = q_{r+i-1} \quad \text{pro } i = 2, \dots, k-1, \quad p_k = \sum_{j=r+k-1}^{\infty} q_j,$$

kde

$$q_i = \text{P}(X_1 = i) = e^{-\lambda} \cdot \frac{\lambda^i}{i!}, \quad i = 0, 1, 2, \dots$$

Principem testování bude opět porovnávání empirických četností Y_i s teoretickými četnostmi $n \cdot p_i$. Protože pravděpodobnosti p_i závisí na neznámém parametru λ , je potřeba tento parametr odhadnout. Podle sekce 12.3 z Anděl (2007) se odhad λ získá řešením rovnice

$$\sum_{i=1}^k \frac{Y_i}{p_i(\lambda)} \cdot \frac{\partial p_i(\lambda)}{\partial \lambda} = 0,$$

která se obvykle řeší numericky (iteračně). Můžeme si ještě všimnout, že

$$\begin{aligned} \frac{\partial p_1(\lambda)}{\partial \lambda} &= \sum_{j=0}^r \frac{\partial q_j(\lambda)}{\partial \lambda}, & \frac{\partial p_k(\lambda)}{\partial \lambda} &= \sum_{j=r+k-1}^{\infty} \frac{\partial q_j(\lambda)}{\partial \lambda}, \\ \frac{\partial p_i(\lambda)}{\partial \lambda} &= \frac{\partial q_{r+i-1}(\lambda)}{\partial \lambda}, & i &= 2, \dots, k-1, \end{aligned}$$

a platí

$$\frac{\partial q_i(\lambda)}{\partial \lambda} = -e^{-\lambda} \cdot \frac{\lambda^i}{i!} + e^{-\lambda} \cdot \frac{\lambda^{i-1}}{(i-1)!} = q_i(\lambda) \left(\frac{i}{\lambda} - 1 \right).$$

Označíme-li jako $\hat{\lambda}$ získaný odhad λ , použijeme testovou statistiku

$$\chi^2 = \sum_{i=1}^k \frac{(Y_i - np_i(\hat{\lambda}))^2}{np_i(\hat{\lambda})}.$$

Ta má za nulové hypotézy asymptoticky χ_{k-2}^2 rozdělení (odhadem λ ztratíme navíc jeden stupeň volnosti). Abychom dostali test s asymptotickou hladinou významnosti α (pro předem zvolené $\alpha \in (0, 1)$), budeme se rozhodovat podle pravidla:

$$\text{zamítneme } (H_0) \iff \chi^2 \geq \chi_{k-2}^2(1 - \alpha).$$

3.2 Testy založené na indexu disperze

Jak již bylo zmíněno v sekci 1.2, pro náhodný výběr X_1, X_2, \dots, X_n pocházející z rozdělení $\text{Po}(\lambda)$, pro nějaké $\lambda > 0$, platí

$$\mathbf{E}X_1 = \lambda = \mathbf{var}X_1.$$

Z toho plyne, že pokud jsou střední hodnota a rozptyl nějaké náhodné veličiny různé, daná veličina nemá Poissonovo rozdělení. Jeden ze způsobů, který vede k zamítnutí hypotézy (1.1), je tedy ukázat, že rozdělení našeho náhodného výběru má rozdílný první a druhý centrovaný moment.

Protože skutečné hodnoty střední hodnoty a rozptylu nejsou známé, budeme se opět muset spokojit s odhady těchto kvantit. K odhadnutí střední hodnoty náhodného výběru použijeme stejně jako v předchozí části výběrový průměr $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. Pro odhad rozptylu je obvyklou volbou výběrový rozptyl

$$S_n^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2,$$

jakožto nestranný a konzistentní odhad $\mathbf{var}X_1$ (viz (iii) ve větě 1).

Statistika používaná k testování nulové hypotézy bude v tomto případě nějaká funkce \bar{X}_n a S_n^2 , která dokáže odrážet, jak moc jsou jejich hodnoty odlišné. Konkrétně se bude jednat o takzvaný index disperze, jehož definici jsme převzali z článku Weiß a kol. (2019).

Definice 3 (Index disperze). *Nechť funkce $f_d : H \rightarrow \mathbb{R}$, $(0, \infty)^2 \subseteq H$, jejíž předpis nezávisí na neznámých parametrech, splňuje*

- (i) $f_d(\mu, \mu) = 1$ pro $\forall (\mu, \mu)^\top \in H$;
- (ii) $f_d(\mu, \sigma^2) \neq 1$, pokud $(\mu, \sigma^2)^\top \in H$ a $\mu \neq \sigma^2$.

Nechť \bar{X}_n a S_n^2 jsou definované jako výše. Potom statistika

$$f_d(\bar{X}_n, S_n^2),$$

*je-li definovaná, se nazývá **index disperze**.*

Platí-li (i) a (ii), potom opět pro $\forall (\mu, \sigma^2)^\top \in H$:

$$f_d(\mu, \sigma^2) = 1 \quad \Leftrightarrow \quad \mu = \sigma^2.$$

Použitím věty o spojitě transformaci dostaneme, že pokud je X_1, X_2, \dots, X_n náhodný výběr z rozdělení s konečným druhým momentem a f_d je spojitá v bodě $(\mathbf{E}X_1, \mathbf{var}X_1)^\top \in H$, platí

$$f_d(\bar{X}_n, S_n^2) \xrightarrow[n \rightarrow \infty]{\text{P}} f_d(\mathbf{E}X_1, \mathbf{var}X_1). \quad (3.2)$$

Dosazením $\mu = \mathbf{E}X_1$ a $\sigma^2 = \mathbf{var}X_1$ do definice indexu disperze potom díky (3.2) dostaneme vztah

$$f_d(\bar{X}_n, S_n^2) \xrightarrow[n \rightarrow \infty]{\text{P}} 1 \quad \Leftrightarrow \quad \mathbf{E}X_1 = \mathbf{var}X_1.$$

Za nulové hypotézy a pro vysoká n tedy bude hodnota indexu disperze $f_d(\bar{X}_n, S_n^2)$ blízká jedné. Naopak v situaci, kdy je zmíněná hodnota významně vzdálená

od čísla 1, můžeme očekávat, že $\mathbf{E}X_1 \neq \mathbf{var}X_1$, a tedy uvažovaný náhodný výběr nepochází z Poissonova rozdělení.

Jednou z nejintuitivnějších možností je zvolit funkci f_d jako

$$\begin{aligned} f_d : (0, \infty) \times \langle 0, \infty \rangle &\rightarrow \langle 0, \infty \rangle, \\ (x, y) &\mapsto \frac{y}{x}. \end{aligned}$$

Snadno ověříme, že podmínky (i) a (ii) jsou splněny, tudíž

$$\hat{\mathbf{I}}_{d,n} := f_d(\bar{X}_n, S_n^2) = \frac{S_n^2}{\bar{X}_n}$$

je podle definice 3 index disperze. Náhodný výběr, jehož nosičem je podmnožina nezáporných celých čísel, navíc v případě $\bar{X}_n = 0$ má i nulový výběrový rozptyl. Použijeme-li konvenci $\frac{0}{0} := 1$, můžeme uvažovat index disperze $\hat{\mathbf{I}}_{d,n}$ i pro $\bar{X}_n = 0$. Z 3.2 máme, že pokud $\mathbf{E}X_1 > 0$, platí

$$\hat{\mathbf{I}}_{d,n} = \frac{S_n^2}{\bar{X}_n} \xrightarrow[n \rightarrow \infty]{\mathbf{P}} \frac{\mathbf{var}X_1}{\mathbf{E}X_1} \stackrel{(1.1)}{=} \frac{\lambda}{\lambda} = 1.$$

Dá se navíc ukázat (viz Semjonov (2020), Věta 7), že za nulové hypotézy (1.1)

$$\sqrt{n}(\hat{\mathbf{I}}_{d,n} - 1) \xrightarrow[n \rightarrow \infty]{\mathbf{D}} \mathbf{N}(0, 2),$$

což je díky vlastnostem normálního rozdělení ekvivalentní s

$$\sqrt{\frac{n}{2}}(\hat{\mathbf{I}}_{d,n} - 1) \xrightarrow[n \rightarrow \infty]{\mathbf{D}} \mathbf{N}(0, 1).$$

Analogicky jako v důkazu věty 9 dostaneme, že za (\mathbf{H}_0) splňuje testová statistika $D_n := \sqrt{\frac{n}{2}}(\hat{\mathbf{I}}_{d,n} - 1)$ rovnost

$$\lim_{n \rightarrow \infty} \mathbf{P}(|D_n| \geq u_{1-\frac{\alpha}{2}}) = \alpha.$$

Z toho plyne, že test s rozhodovacím pravidlem

$$\text{zamítáme } (\mathbf{H}_0) \iff |D_n| \geq u_{1-\frac{\alpha}{2}}$$

má asymptotickou hladinu významnosti α , kde $\alpha \in (0, 1)$ je nějaké předem zvolené číslo.

Podobně jako testy založené na nulovém indexu, je výše popsáný test rovněž asymptotický test, který testuje pouze jednu charakteristiku Poissonova rozdělení. Nastane-li situace, kdy $\mathbf{E}X_1 = \mathbf{var}X_1$, test dobré shody založený na indexu disperze obvykle potenciálně neplatnou hypotézu nezamítne.

4. Simulační studie

V dosavadní části práce jsme teoreticky popsali řadu metod, pomocí nichž lze testovat dobrou shodu s Poissonovým rozdělením. Konkrétně se jednalo o testy založené na nulovém indexu, test využívající index disperze a χ^2 -testy dobré shody. Nyní se pokusíme zmíněné způsoby testování porovnat pomocí dat nasimulovaných v programu R Core Team (2021). Všechny testy si budeme implementovat sami (nebudeme používat žádné dodatečné balíčky).

4.1 Hladina testů

Jako první budeme zkoumat, jak moc dané testy dodržují předepsanou hladinu významnosti. Bude nás tedy zajímat, v jaké míře bude v jednotlivých případech docházet k zamítání platné nulové hypotézy. Z předchozích kapitol víme, že všechny uvedené testy jsou pouze asymptotické. Můžeme tak očekávat, že zejména pro nízké hodnoty n (rozsahu náhodného výběru) stanovená hladina příliš dodržována nebude, pro různé testy se však míra jejího dodržování může lišit.

Náš postup bude nyní spočívat v tom, že si nejprve vygenerujeme realizace náhodného výběru X_1, X_2, \dots, X_n z rozdělení $\text{Po}(\lambda)$, a následně provedeme několik dílčích testů dobré shody, přičemž si vždy zaznamenáme, pro které z nich jsme zamítli (v této situaci platnou) hypotézu (1.1). Tento proces 1000krát zopakujeme. Pro každý test tak dostaneme číslo N_i udávající počet případů, kdy došlo k zamítnutí (1.1). Pomocí něj spočteme odhad hladiny daného testu jako

$$\frac{N_i}{1000},$$

tj. poměr případů zamítnutí platné hypotézy ku celkovému počtu pokusů. Navíc rozlišíme několik různých hodnot parametru λ a rozsahu výběru n .

Mezi zkoumané testy budou patřit v první řadě testy založené na nulových indexech ze sekce 2.2 popsané ve větě 10, dále test s testovou statistikou D_n založený na indexu disperze (viz sekce 3.2) a χ^2 -testy dobré shody ze sekce 3.1. U χ^2 -testů budeme uvažovat případy pro 4 a 6 kategorií ($k = 4$ a $k = 6$), přičemž u obou zvolíme $r = 0$. Pro první volbu parametrů tak získáme statistiku

$$\chi_{1,n}^2 := \sum_{i=1}^4 \frac{(Y_i - np_i(\hat{\lambda}))^2}{np_i(\hat{\lambda})},$$

kde

$$Y_1 := \sum_{j=1}^n \mathbf{1}_{[X_j=0]}, \quad Y_2 := \sum_{j=1}^n \mathbf{1}_{[X_j=1]}, \quad Y_3 := \sum_{j=1}^n \mathbf{1}_{[X_j=2]}, \quad Y_4 := \sum_{j=1}^n \mathbf{1}_{[X_j \geq 3]},$$

a $\hat{\lambda}$ je odhad λ získaný řešením příslušné rovnice (viz sekce 3.1). Hypotézu (1.1) potom budeme zamítat právě tehdy, když $\chi_{1,n}^2 \geq \chi_2^2(1 - \alpha)$. Testovou statistiku příslušející volbě $k = 6$ označíme $\chi_{2,n}^2$. U testů, kde se v kritickém oboru vyskytuje kvantil normovaného normálního rozdělení (tj. testů založených na nulových

indexech a indexu disperze) navíc vyzkoušíme i variantu, kde tento kvantil nahradíme kvantilem t-rozdělení s $n - 1$ stupni volnosti tak, jak bylo diskutováno v sekci 2.4.

Všechny zmíněné testy budeme vždy provádět na hladině $\alpha = 0,05$.

V následujících tabulkách jsou v jednotlivých sloupcích vypsány získané odhady hladiny pro test založený na dané testové statistice. Pokud je testová statistika označena hvězdičkou, provádíme test s kvantilem t-rozdělení namísto $N(0, 1)$ rozdělení. Pro přehlednost u těchto případů uvádíme pouze ty hodnoty, které se od těch původních liší o více než 0,005.

Tabulka 4.1: Odhady hladiny jednotlivých testů pro případ $\lambda = 1$

n	$T_{1,n}$	$T_{1,n}^*$	$T_{2,n}$	$T_{2,n}^*$	$T_{3,n}$	$T_{3,n}^*$	$T_{4,n}$	$T_{4,n}^*$
20	0,080	0,058	0,056	0,045	0,105	0,092	0,057	0,036
50	0,071	...	0,055	...	0,083	0,077	0,057	0,050
100	0,052	0,043	0,051	...	0,061	...	0,050	...
500	0,056	...	0,053	...	0,056	...	0,054	...
2000	0,047	...	0,044	...	0,047	...	0,042	...

n	D_n	D_n^*	$\chi_{1,n}^2$	$\chi_{2,n}^2$
20	0,046	0,034	0,048	0,056
50	0,051	...	0,046	0,058
100	0,054	...	0,051	0,047
500	0,049	...	0,046	0,052
2000	0,048	...	0,048	0,054

Tabulka 4.2: Odhady hladiny jednotlivých testů pro případ $\lambda = 2$

n	$T_{1,n}$	$T_{1,n}^*$	$T_{2,n}$	$T_{2,n}^*$	$T_{3,n}$	$T_{3,n}^*$	$T_{4,n}$	$T_{4,n}^*$
20	0,121	0,111	0,051	0,033	0,189	0,176	0,100	0,087
50	0,089	0,082	0,051	...	0,141	...	0,061	0,053
100	0,077	...	0,055	...	0,091	...	0,060	...
500	0,061	...	0,051	...	0,059	...	0,052	...
2000	0,051	...	0,046	...	0,056	...	0,046	...

n	D_n	D_n^*	$\chi_{1,n}^2$	$\chi_{2,n}^2$
20	0,056	0,047	0,047	0,036
50	0,040	...	0,041	0,046
100	0,039	...	0,059	0,053
500	0,053	...	0,044	0,052
2000	0,043	...	0,043	0,047

Z tabulek 4.1 a 4.2 můžeme vidět, že zatímco testy založené na indexu disperze a χ^2 -testy dodržují stanovenou hladinu $\alpha = 0,05$ poměrně dobře i pro nízké rozsahy výběrů, u testů založených na nulových indexech (s výjimkou testů se statistikou $T_{2,n}$) je pro malá n získaný odhad hladiny mnohdy významně vyšší. Podobný závěr si můžeme udělat i z tabulky 4.3 níže.

Tabulka 4.3: Odhady hladiny jednotlivých testů pro případ $\lambda = 3$

n	$T_{1,n}$	$T_{1,n}^*$	$T_{2,n}$	$T_{2,n}^*$	$T_{3,n}$	$T_{3,n}^*$	$T_{4,n}$	$T_{4,n}^*$
20	0,364	...	0,027	0,019	0,371	...	0,375	...
50	0,112	...	0,033	...	0,213	...	0,099	...
100	0,095	...	0,050	...	0,151	0,145	0,077	...
500	0,064	...	0,060	...	0,073	...	0,063	...
2000	0,058	...	0,051	...	0,058	...	0,052	...

n	D_n	D_n^*	$\chi_{1,n}^2$	$\chi_{2,n}^2$
20	0,041	0,030	0,040	0,047
50	0,046	...	0,045	0,048
100	0,052	...	0,050	0,050
500	0,044	...	0,064	0,052
2000	0,041	...	0,051	0,055

Lze si také všimnout, že u většiny zkoumaných testů založených na nulových indexech platí, že čím vyšší je parametr (skutečného) Poissonova rozdělení, tím více dat je potřeba k tomu, aby test dodržoval požadovanou hladinu. To může souviset s tím, že s rostoucím λ klesá poměr nulových pozorování a častěji dochází k případům, kdy jsou dokonce všechna data nenulová (viz sekce 2.4).

Z definice jednotlivých nulových indexů a následných testových statistik vidíme, že pokud $\hat{p}_0 = 0$, dostaneme bez ohledu na bližší podobu dat

$$T_{1,n} = -\infty \quad \text{a} \quad T_{3,n} = \infty,$$

tedy nulová hypotéza bude pro testy s danými testovými statistikami (a k nim náležející testy využívající kvantily t-rozdělení) zamítnuta. V situaci, kdy generujeme data z rozdělení $\text{Po}(3)$ o rozsahu 20, se přitom přibližně ve třetině případů nevyskytne ani jedna nulová hodnota.

U testů založených na nulových indexech $\hat{I}_{2,n}$ a $\hat{I}_{4,n}$ bude v případě absence nulových pozorování platit $\hat{I}_{2,n} = -1$ a $\hat{I}_{4,n} = 1$, a tedy

$$T_{2,n} = \frac{-1}{\sqrt{\frac{e^{\bar{X}_n} - \bar{X}_n - 1}{n}}}, \quad T_{4,n} = \frac{1}{\sqrt{\frac{e^{\bar{X}_n} - \bar{X}_n - 1}{n(\bar{X}_n)^2}}}.$$

Absolutní hodnoty daných testových statistik tak v této situaci závisejí pouze na velikosti odhadu asymptotického rozptylu. Toto může být jeden z důvodů, proč pro vyšší hodnoty parametru λ vycházejí některé odhady hladiny naopak výrazně nižší než stanovených 0,05.

Rovněž je vidět, že používání kvantilů Studentova t-rozdělení dává významně rozdílné výsledky pouze u nízkých rozsahů výběru (v našem případě $n = 20$). Při větším počtu dat je tak nahrazování kvantilu normovaného normálního rozdělení tímto kvantilem vcelku zbytečné.

Na základě získaných odhadů se jako nejlepší z hlediska dodržování hladiny jeví test založený na nulovém indexu $\hat{I}_{2,n}$, test založený na indexu disperze se statistikou D_n a oba χ^2 -testy.

4.2 Síla testů

Nyní nás bude zajímat síla zkoumaných testů proti určitým alternativám. Připomeneme, že síla testu při dané konkrétní alternativě je rovna pravděpodobnosti, se kterou zamítneme (v této situaci neplatnou) nulovou hypotézu, pochází-li naše data z rozdělení příslušného této alternativě.

Při odhadování síly uvažovaných testů budeme postupovat analogicky jako v předchozí části s jediným rozdílem – data budeme generovat z rozdělení různého od Poissonova.

Jako alternativu budeme uvažovat takzvané negativně binomické rozdělení (pro různé hodnoty parametrů). Pro připomenutí, pro $X \sim \text{NB}(\mu, r)$, kde $\mu > 0$ a $r > 0$, platí

$$\mathbb{P}(X = k) = \binom{r+k-1}{k} \left(\frac{r}{r+\mu}\right)^r \left(\frac{\mu}{r+\mu}\right)^k, \quad k \in \mathbb{N}_0.$$

Zároveň

$$\mathbb{E}X = \mu \quad \text{a} \quad \text{var}X = \mu + \frac{\mu^2}{r}.$$

Jedná se tedy o diskrétní rozdělení se stejným nosičem jako má Poissonovo rozdělení.

Poznámka. Toto je alternativní parametrizace negativně binomického rozdělení. Obvykle definujeme negativně binomické rozdělení pomocí definice uvedené např. v sekci 1.2.3 v knize Anděl (2007), tj. jako rozdělení s parametry $r > 0$ a $p \in (0, 1)$ splňující

$$\mathbb{P}(X = k) = \binom{r+k-1}{k} p^r (1-p)^k, \quad k \in \mathbb{N}_0.$$

Platí $\mathbb{E}X = \frac{r(1-p)}{p}$ a $\text{var}X = \frac{r(1-p)}{p^2}$. Výše uvedený způsob parametrizace tak získáme, pokud položíme $\mu := \mathbb{E}X = \frac{r(1-p)}{p}$. Není těžké ověřit, že potom $p = \frac{r}{r+\mu}$.

Z definice negativně binomického rozdělení dostaneme, že pro náhodný výběr X_1, X_2, \dots, X_n pocházející z rozdělení $\text{NB}(\mu, r)$ máme

$$\mathbb{P}(X_1 = 0) = \left(\frac{r}{r+\mu}\right)^r, \quad (4.1)$$

tedy obecně neplatí $p_0 = e^{-\mathbb{E}X_1}$. Můžeme si ovšem všimnout, že

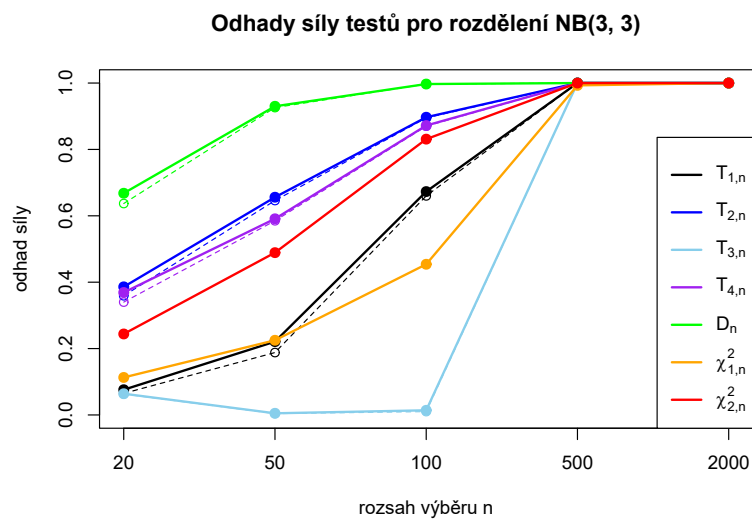
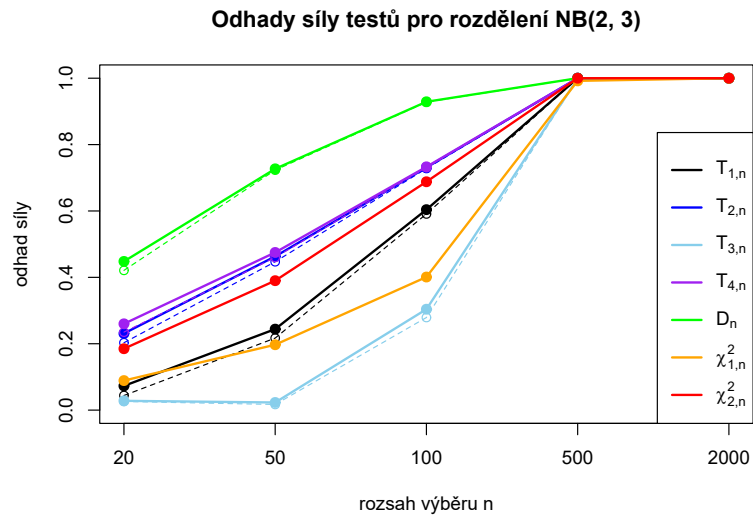
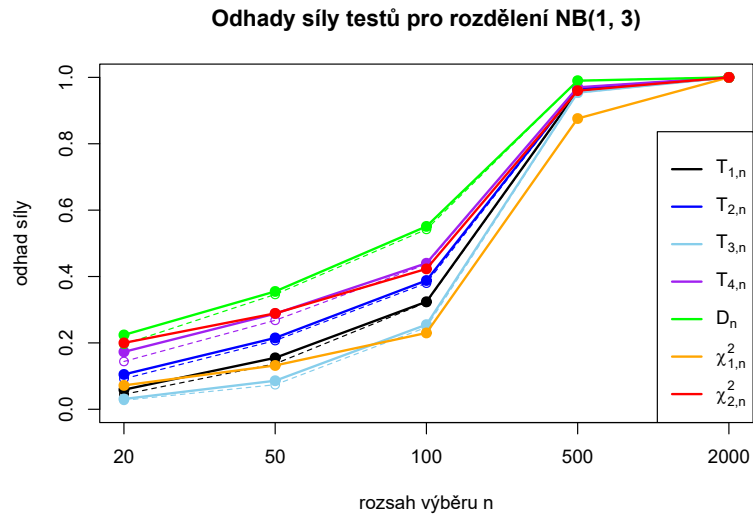
$$\lim_{r \rightarrow \infty} \mathbb{P}(X_1 = 0) = e^{-\mu} = e^{-\mathbb{E}X_1}.$$

Je možné ukázat (viz Šír (2020), 2.4) i silnější vztah:

$$\lim_{r \rightarrow \infty} \mathbb{P}(X_1 = k) = e^{-\mu} \cdot \frac{\mu^k}{k!}, \quad k \in \mathbb{N}_0,$$

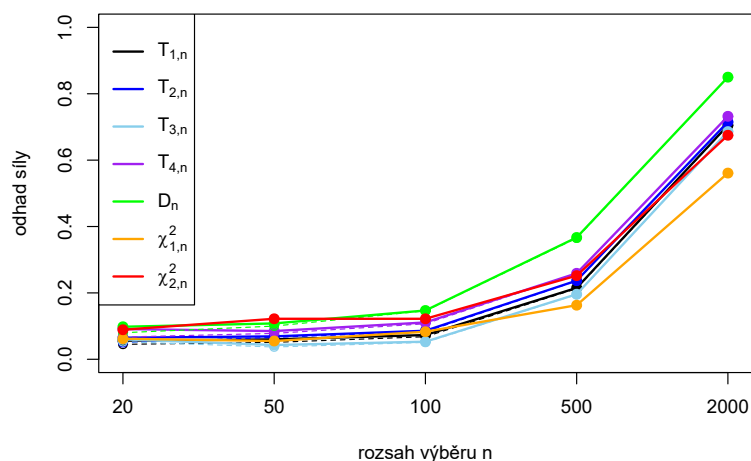
což nám dává $\lim_{r \rightarrow \infty} \mathbb{P}(X_1 = k) = \mathbb{P}(Y = k)$ pro $Y \sim \text{Po}(\mu)$. Platí rovněž $\lim_{r \rightarrow \infty} \text{var}X_1 = \mu = \mathbb{E}X_1$.

Můžeme tak očekávat, že pro testy dobré shody založené na nulovém indexu, χ^2 -testy i test založený na indexu disperze bude s rostoucím parametrem r (a pevným μ) jejich síla vůči alternativě $\text{NB}(\mu, r)$ klesat, neboť rozdíly mezi rozdělením $\text{NB}(\mu, r)$ a Poissonovým rozdělením s parametrem μ se budou zmenšovat.

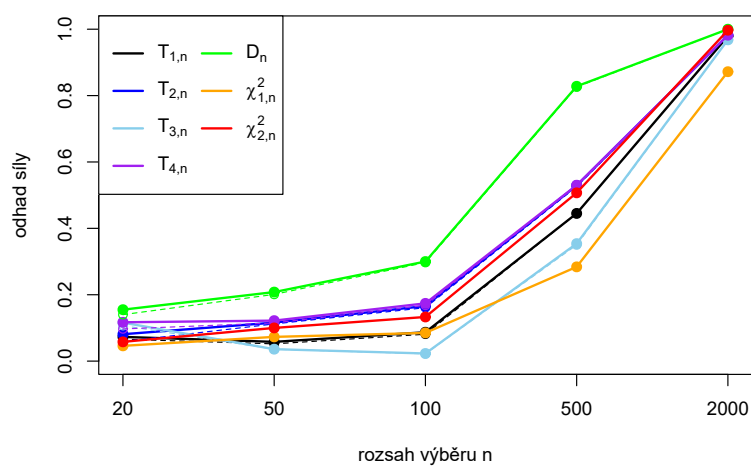


Obrázek 4.1: Odhady síly jednotlivých testů vůči alternativě $NB(\mu, r)$ pro $r = 3$ a různé hodnoty μ (přerušované čáry značí odhady pro testy, kde jsou kvantily rozdělení $N(0, 1)$ nahrazeny kvantily t-rozdělení)

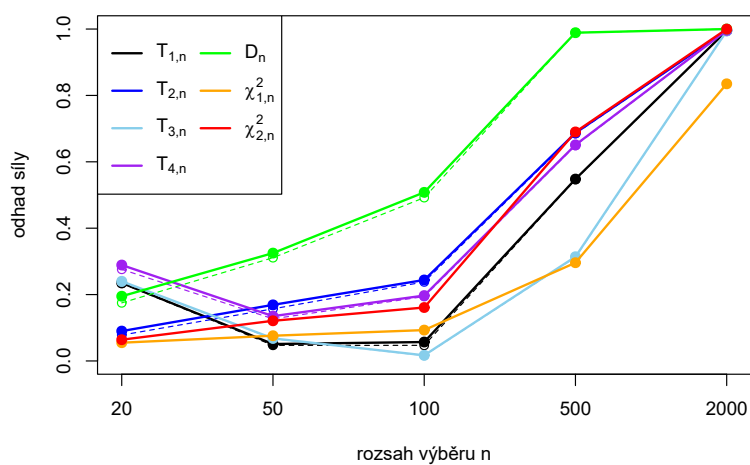
Odhady síly testů pro rozdělení NB(1, 10)



Odhady síly testů pro rozdělení NB(2, 10)



Odhady síly testů pro rozdělení NB(3, 10)



Obrázek 4.2: Odhady síly jednotlivých testů vůči alternativě $NB(\mu, r)$ pro $r = 10$ a různé hodnoty μ (přerušované čáry značí odhady pro testy, kde jsou kvantily rozdělení $N(0, 1)$ nahrazeny kvantily t-rozdělení)

Na obrázcích 4.1 a 4.2 jsme mohli vidět, že u každého testu se odhad jeho síly proti určité konkrétní alternativě pro $n \rightarrow \infty$ blíží k 1. Testy tak můžeme považovat za konzistentní (vůči daným alternativám).

Porovnáním 4.1 a 4.2 si můžeme všimnout, že odhadnuté síly testů proti alternativě $\text{NB}(\mu, 3)$ obvykle rostou při zvětšujícím se rozsahu výběru rychleji než odhadnuté síly proti alternativě $\text{NB}(\mu, 10)$ (je-li μ v obou případech stejné). To odpovídá poznatkům diskutovaným v předchozí části, tedy že pro vyšší hodnotu r je rozdělení $\text{NB}(\mu, r)$ „bližší“ Poissonovu rozdělení.

Trochu neobvyklá situace nastává u některých testů založených na nulových indexech, kde proti určitým alternativám (např. $\text{NB}(3, 10)$) dochází při zvyšování rozsahu výběru k dočasnému poklesu jejich odhadnuté síly. Konkrétně to lze vidět především u testů s testovými statistikami $T_{1,n}$ a $T_{3,n}$. Jak ovšem víme z předchozí sekce, tyto testy mají pro nízký počet pozorování (v našem případě $n = 20$) a vyšší λ (střední hodnotu skutečného rozdělení) i velmi vysokou hladinu významnosti, která zřejmě souvisí s absencí nulových pozorování. Pomocí vztahu (4.1) můžeme nahlédnout, že pro vyšší hodnoty parametrů μ a r bude pravděpodobnost $P(X_1 = 0)$, kde $X_1 \sim \text{NB}(\mu, r)$, rovněž velmi malá.

Z dostupných odhadů síly jednotlivých testů proti alternativám $\text{NB}(\mu, r)$ (pro dané hodnoty parametrů) lze jako nejsilnější test (vůči zmíněným alternativám) ze všech zkoumaných testů dobré shody považovat test založený na indexu disperze, tj. test s testovou statistikou D_n . Jako poměrně silné se jeví i testy založené na nulových indexech $\hat{I}_{2,n}$ a $\hat{I}_{4,n}$ a χ^2 -test dobré shody rozlišující 6 kategorií (test s testovou statistikou $\chi_{2,n}^2$). Naopak nejhůře v tomto ohledu dopadl test, jehož testová statistika má tvar $T_{3,n}$.

Závěr

K testování, zda rozdělení určitého náhodného výběru odpovídá Poissonovu rozdělení, existuje celá řada metod. Jedna z nich – testy založené na tzv. nulovém indexu – byla hlavním předmětem této práce.

Nulovým indexům a jejich využití při testování jsme věnovali celou druhou kapitolu. Nejprve jsme obecně zavedli Poissonův nulový index jakožto určitou funkci výběrového průměru a relativní četnosti nul v realizacích náhodného výběru. Poté jsme pomocí Δ -věty odvodili asymptotické rozdělení těchto indexů za platnosti hypotézy, že daný náhodný výběr pochází z Poissonova rozdělení, a využili jej ke konstrukci asymptotických testů dobré shody. Pro lepší představu jsme vše provedli i pro čtyři konkrétní příklady nulových indexů. Ve třetí kapitole jsme se podívali i na jiné způsoby testování dobré shody – variantu χ^2 -testů pro Poissonovo rozdělení a testy založené na indexu disperze. Všechny zmíněné typy testů jsme potom porovnali pomocí simulací v kapitole 4.

Mezi hlavní výhody testů založených na nulových indexech patří jejich výpočetní jednoduchost (obzvláště ve srovnání s χ^2 -testy). Na druhou stranu však tyto testy nemají tak všestranné využití, neboť nejsou konzistentní vůči všem alternativám. Také se příliš nehodí používat v situacích, kdy testujeme dobrou shodu s rozdělením $Po(\lambda)$ pro nějakou větší hodnotu parametru λ . V takovém případě je lepší volbou například test založený na indexu disperze se statistikou D_n či χ^2 -testy dobré shody s vhodně zvolenými kategoriemi.

Ze všech zkoumaných testů lze na základě simulační studie doporučit výše zmíněný test s testovou statistikou D_n a test založený na nulovém indexu se statistikou $T_{2,n}$, neboť tyto testy poměrně dobře dodržují stanovenou hladinu významnosti a jeví se i jako nejsilnější vůči zkoumaným alternativám (negativně binomickému rozdělení pro různé hodnoty parametrů). Naopak nejméně vhodný se zdá být test se statistikou $T_{3,n}$ (založený na nulovém indexu), který má obzvláště pro nižší rozsahy výběru (cca $n \leq 100$) velmi vysokou hladinu a zároveň malou sílu proti daným alternativám $NB(\mu, r)$.

Seznam použité literatury

- ANDĚL, J. (1978). *Matematická statistika*. SNTL/ALFA, Praha.
- ANDĚL, J. (2007). *Základy matematické statistiky*. Druhé opravené vydání. Matfyzpress, Praha. ISBN 80-7378-001-1.
- DUPAČ, V. a HUŠKOVÁ, M. (2013). *Pravděpodobnost a matematická statistika*. Druhé upravené vydání. Nakladatelství Karolinum, Praha. ISBN 978-80-246-2208-8.
- R CORE TEAM (2021). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- SEMJONOV, V. (2020). Index disperze pro diskrétní rozdělení. Bakalářská práce, Matematicko-fyzikální fakulta, Univerzita Karlova.
- VAN DER VAART, A. W. (1998). *Asymptotic Statistics*. Cambridge University Press, Cambridge, United Kingdom.
- WEISS, C. H., HOMBURG, A. a PUIG, P. (2019). Testing for zero inflation and overdispersion in INAR(1) models. *Statist. Papers* 60, **26**, 823–848.
- ŠÍR, D. (2020). Dvourozměrné negativně binomické rozdělení. Bakalářská práce, Matematicko-fyzikální fakulta, Univerzita Karlova.