



**MATEMATICKO-FYZIKÁLNÍ  
FAKULTA**  
Univerzita Karlova

## **BAKALÁRSKA PRÁCA**

Filip Kulla

# **Dvojvýberový Wilcoxonov test v prípade existencie zhôd**

Katedra pravděpodobnosti a matematické statistiky, MFF UK

Vedúci bakalárskej práce: doc. Ing. Omelka Marek, Ph.D.

Študijný program: Matematika

Študijný obor: Obecná matematika

Praha 2020

Prehlasujem, že som túto bakalársku prácu vypracoval samostatne a výhradne s použitím citovaných prameňov, literatúry a ďalších odborných zdrojov. Táto práca nebola využitá k získaniu iného alebo rovnakého titulu.

Beriem na vedomie, že sa na moju prácu vzťahujú práva a povinnosti vyplývajúce zo zákona č. 121/2000 Sb., autorského zákona v platnom znení, najmä skutočnosť, že Univerzita Karlova má právo na uzavrenie licenčnej zmluvy o použití tejto práce ako školského diela podľa §60 odst. 1 autorského zákona.

V ..... dňa .....

Podpis autora

Ďakujem môjmu školiteľovi, doc. Ing. Marekovi Omelkovi, Ph.D., za vedenie, užitočné rady a prejavenu ochotu. Rovnako ďakujem aj mojej rodine a priateľom.

Názov práce: Dvojjvýberový Wilcoxonov test v prípade existencie zhôd

Autor: Filip Kulla

Katedra: Katedra pravděpodobnosti a matematické statistiky, MFF UK

Vedúci bakalárskej práce: doc. Ing. Omelka Marek, Ph.D., Katedra pravděpodobnosti a matematické statistiky, MFF UK

Abstrakt: Táto práca sa venuje dvojjvýberovému Wilcoxonovmu testu v prípade existencie zhôd. V práci je s využitím známych výsledkov o U-štatistikách odvodené asymptotické rozdelenie Wilcoxonovej testovej štatistiky v prípade existencie zhôd. S pomocou tohto výsledku je navrhnutá korekcia testu pre dáta obsahujúce zhody. Ďalej sa v práci skúma súvislosť odvodenej korekcie a podmieneného rozptylu testovej štatistiky (a podmienenej strednej hodnoty) za  $H_0$ . Nakoniec je v práci pomocou simulácií preskúmané, aký vplyv má odvodená korekcia na skutočnú hladinu testu pri zvyšujúcom sa počte zhodných pozorovaní.

Kľúčové slová: dvojjvýberový Wilcoxonov test, Mannov-Whitneyov U test, zhody.

Title: Two-sample Wilcoxon test in the presence of ties

Author: Filip Kulla

Department: Department of Probability and Mathematical Statistics, MFF UK

Supervisor: doc. Ing. Omelka Marek, Ph.D., Department of Probability and Mathematical Statistics, MFF UK

Abstract: This work is devoted to the Wilcoxon rank-sum test in the presence of ties. In this work, asymptotic distribution of the Wilcoxon test statistic in the presence of ties is derived using well-known results about U-statistics. With the aid of this result, a corrected test statistic for data containing ties is proposed. Furthermore, the thesis examines the relation between the derived correction and conditional variance of the test statistic (and conditional expectation). Finally, by means of simulation, the effect of the derived correction on the actual significance level is examined as the number of tied observations increases.

Keywords: Wilcoxon rank-sum test, Mann–Whitney U test, ties.

# Obsah

Úvod	2
<b>1 Dvojjvýberový Wilcoxonov test bez zhôd</b>	<b>3</b>
1.1 Zavedenie testovej štatistiky . . . . .	3
1.2 Odvodenie asymptotického rozdelenia testovej štatistiky v prípade bez zhôd . . . . .	4
<b>2 Dvojjvýberový Wilcoxonov test so zhodami</b>	<b>7</b>
2.1 Zavedenie testovej štatistiky v prípade zhôd . . . . .	7
2.2 Odvodenie asymptotického rozdelenia testovej štatistiky v prípade zhôd . . . . .	8
<b>3 Súvislosť asymptotického a podmieneného rozptylu <math>W_X^*</math></b>	<b>16</b>
<b>4 Porovnanie testu bez korekcie a testu s korekciou pomocou simulácie</b>	<b>20</b>
<b>5 Dodatky</b>	<b>23</b>
5.1 U-štatistiky . . . . .	23
5.2 Odvodenie $\sigma_{10}^2$ z vety 5 . . . . .	23
<b>Záver</b>	<b>25</b>
<b>Zoznam použitej literatúry</b>	<b>26</b>

# Úvod

Táto práca sa venuje dvojvýberovému Wilcoxonovmu testu. Ide o neparametrický test založený na poradiach, pomocou ktorého porovnáваме 2 nezávislé náhodné výbery (napr. pacienti, ktorým bolo podané liečivo a pacienti, ktorí dostali placebo). Štandardná teória poradových testov predpokladá, že dáta pochádzajú zo spojitých rozdelení, čo zaručuje, že pravdepodobnosť zhodných pozorovaní je nulová. V praxi sa ale častokrát stáva, že predpoklad spojitosti nie je splnený a v dátach sa zhody nachádzajú (napr. nevieme rozhodnúť či je stav prvého pacienta lepší ako stav druhého pacienta alebo telesná teplota prvého pacienta je po zaokrúhlení na 1 desatinné miesto rovnaká ako zaokrúhlená teplota druhého pacienta).

V tejto práci budeme skúmať vplyv prítomnosti zhodných pozorovaní na dvojvýberový Wilcoxonov test. Naším cieľom bude odvodiť asymptotické rozdelenie Wilcoxonovej testovej štatistiky v prípade existencie zhôd, na základe dosiahnutého výsledku navrhnúť korekciu testu pre prípad, že naše dáta obsahujú zhody a nakoniec pomocou simulácií porovnať skutočnú hladinu testu s korekciou a testu bez korekcie v prítomnosti zhodných pozorovaní.

Hlavná časť práce je rozdelená do štyroch kapitol. Prvá kapitola je venovaná predstaveniu dvojvýberového Wilcoxonovho testu, pričom zatiaľ predpokladáme spojitost rozdelení, z ktorých dáta pochádzajú. V druhej kapitole upúšťame od predpokladu spojitosti. Jej podstatnú časť tvorí odvodenie asymptotického rozdelenia testovej štatistiky v prípade existencie zhôd, ktoré nás vedie k návrhu korekcie testu pre dáta so zhodami. Tretia kapitola ukazuje súvislosť medzi odvodenou korekciou a podmieneným rozptylom testovej štatistiky (a podmienenou strednou hodnotou) za  $H_0$ . Štvrtá kapitola obsahuje simulačnú štúdiu, v ktorej zistujeme skutočnú hladinu testu s korekciou a testu bez korekcie.

# 1. Dvojvýberový Wilcoxonov test bez zhôd

Dvojvýberový Wilcoxonov test je neparametrický test založený na poradiach. Uvažujme 2 nezávislé náhodné výbery  $\mathbf{X} = (X_1, \dots, X_n)$  a  $\mathbf{Y} = (Y_1, \dots, Y_m)$ . Nech  $\mathbf{X}$  je náhodný výber z rozdelenia s distribučnou funkciou  $F_X$  a  $\mathbf{Y}$  je náhodný výber z rozdelenia s distribučnou funkciou  $F_Y$ . Predpokladajme, že  $F_X$  a  $F_Y$  sú spojité. Tento predpoklad implikuje, že pravdepodobnosť zhodných pozorovaní je 0. Formálnejšie, ak zavedieme značenie  $N = n + m$  a  $\mathbf{Z} = (Z_1, \dots, Z_N) = (X_1, \dots, X_n, Y_1, \dots, Y_m)$ , tak

$$P(Z_i = Z_j, \text{ pre nejaké } i, j \in \{1, \dots, N\}) = 0. \quad (1.1)$$

Predpokladajme, že navyše

$$\exists \delta \in \mathbb{R} \forall x \in \mathbb{R} : F_X(x) = F_Y(x - \delta).$$

Takýto model sa niekedy nazýva model posunutia v polohe.

Zaujíma nás posunutie  $\delta$ . Presnejšie, chceme testovať hypotézu

$$H_0 : \delta = 0, \quad (1.2)$$

proti alternatíve

$$H_1 : \delta \neq 0.$$

V modeli posunutia v polohe teda dvojvýberový Wilcoxonov test testuje, či náhodné výbery  $\mathbf{X}$  a  $\mathbf{Y}$  pochádzajú z rovnakého rozdelenia.

## 1.1 Zavedenie testovej štatistiky

Zoradíme všetky náhodné veličiny v združenom výbere  $\mathbf{Z}$  od najmenej po najväčšiu:  $Z_{(1)} < Z_{(2)} < \dots < Z_{(N)}$ . Symbol  $Z_{(k)}$  označuje  $k$ -tu najmenšiu hodnotu medzi pozorovaniami zo  $\mathbf{Z}$ . Vďaka predpokladu (1.1), vyplývajúcej zo spojitosti  $F_X$  a  $F_Y$ , je usporiadanie pozorovaní zo  $\mathbf{Z}$  dobre definované a nemusíme uvažovať neostré nerovnosti.

**Definícia 1.** Vzostupne usporiadaný výber  $\mathbf{Z}$ ,  $Z_{(1)} < Z_{(2)} < \dots < Z_{(N)}$ , budeme nazývať usporiadaný výber a označovať  $\mathbf{Z}_{(\cdot)}$

**Definícia 2.** Poradím náhodnej veličiny  $Z_i$ ,  $i \in \{1, \dots, N\}$ , v združenom výbere  $\mathbf{Z}$  rozumieme prirodzené číslo  $R_i \in \{1, \dots, N\}$  také, že  $Z_i = Z_{(R_i)}$ .

Platí, že

$$R_i = \sum_{j=1}^N 1\{Z_j \leq Z_i\}. \quad (1.3)$$

Špeciálne, pre  $i \in \{1, \dots, n\}$  je  $R_i = \sum_{j=1}^n 1\{X_j \leq X_i\} + \sum_{j=1}^m 1\{Y_j \leq X_i\}$ .

Wilcoxonova testová štatistika je definovaná ako suma poradí náhodných veličín z  $\mathbf{X}$  v združenom výbere  $\mathbf{Z}$ , teda

$$W_X = \sum_{i=1}^n R_i.$$

Testová štatistika  $W_X$  môže nadobúdať hodnoty  $n(n+1)/2, \dots, nm + n(n+1)/2$ . Proti  $H_0$  zrejme svedčia príliš malé a zároveň príliš veľké hodnoty  $W_X$ .

Pokiaľ nie sú  $n$  a  $m$  príliš veľké, je možné určiť presné rozdelenie  $W_X$  za  $H_0$ . Toto presné rozdelenie sa dá odvodiť z toho, že za  $H_0$  je ľubovoľné usporiadanie  $\mathbf{Z}$  rovnako pravdepodobné, a teda

$$P(R_1 = r_1, \dots, R_n = r_n) = \frac{m!}{N!},$$

pre všetky  $r_1, \dots, r_n \in \{1, \dots, N\}$  rôzne.

V prípade, že  $n$  a  $m$  sú veľké, využijeme na zvolenie kritického oboru tvrdenie 1 a vetu 3 zo sekcie 1.2. Aby sme dostali test s asymptotickou hladinou  $\alpha$ , zvolíme kritický obor tak, že

$$H_0 \text{ zamietame} \Leftrightarrow \frac{|W_X - \frac{n(N+1)}{2}|}{\sqrt{\frac{nm(N+1)}{12}}} \geq u_{1-\alpha/2},$$

kde  $u_{1-\alpha/2}$  značí  $(1 - \alpha/2)$ -kvantil rozdelenia  $N(0,1)$ .

**Tvrdenie 1.** *Za platnosti  $H_0$  je*

$$E W_X = \frac{n(n+m+1)}{2} \quad a \quad var W_X = \frac{nm(n+m+1)}{12}.$$

*Dôkaz.* Postup ako určiť strednú hodnotu a rozptyl sa dá nájsť v knihe Anděl (2007b, strana 102). □

Na odvodenie asymptotického rozdelenia  $W_X$  použijeme teóriu U-štatistík, viď sekcia 5.1 v kapitole Dodatky.

## 1.2 Odvodenie asymptotického rozdelenia testovej štatistiky v prípade bez zhôd

V tejto sekcii použijeme vetu 12 zo sekcie 5.1 na odvodenie asymptotického rozdelenia  $W_X$ . Definujme

$$W_{XY} = \sum_{i=1}^n \sum_{j=1}^m 1\{X_i < Y_j\}.$$

Štatistika  $W_{XY}$  sa často označuje ako Mannova-Whitneyova. S Wilcoxonovou štatistikou ju spája deterministický vzťah.

**Tvrdenie 2.** *Platí, že  $W_{XY} = nm + \frac{n(n+1)}{2} - W_X$ .*



*Dôkaz.* Toto tvrdenie sa dokáže analogicky ako tvrdenie 4 v sekcii 2.2. □

**Veta 3.** Ak platí  $H_0$  a  $n/N \rightarrow p \in (0,1)$ , keď  $n + m \rightarrow \infty$ , tak

$$\frac{W_X - E_{H_0} W_X}{\sqrt{\text{var}_{H_0} W_X}} \xrightarrow{D} N(0,1).$$

*Dôkaz.* Dôkaz bude aplikáciou vety 12. Definujme jadro  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$  predpisom

$$h(x,y) = \begin{cases} 1, & \text{ak } x < y, \\ 0, & \text{inak.} \end{cases}$$

Potom korešpondujúca U-štatistika pre jadro  $h$  je tvaru

$$U_{n,m} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m h(X_i, Y_j) = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m 1\{X_i < Y_j\}.$$

Zrejme platí, že

$$U_{n,m} = \frac{W_{XY}}{nm}. \quad (1.4)$$

Teraz spočítame  $\theta, \sigma_{10}^2, \sigma_{01}^2$  z vety 12 a overíme konečnosť  $\sigma_{11}^2$ .

Platí, že

$$\theta = E[h(X_i, Y_j)] = E[1\{X_i < Y_j\}] = P(X_i < Y_j) \stackrel{H_0}{=} \frac{1}{2},$$

pričom v kroku označenom  $\stackrel{H_0}{=}$  používame predpoklad, že platí  $H_0$ .

Pre  $\sigma_{10}^2$  a  $\sigma_{01}^2$  dostávame nasledujúce vzťahy:

$$\begin{aligned} \sigma_{10}^2 &= \text{cov}(1\{X_1 < Y_1\}, 1\{X_1 < Y_2\}) = E\left[1\{X_1 < Y_1\}1\{X_1 < Y_2\}\right] - \theta^2 \\ &= P(X_1 < Y_1, X_1 < Y_2) - \theta^2 \stackrel{H_0}{=} \frac{1}{3} - \frac{1}{4} = \frac{1}{12}, \end{aligned}$$

$$\begin{aligned} \sigma_{01}^2 &= \text{cov}(1\{X_1 < Y_1\}, 1\{X_2 < Y_1\}) = E\left[1\{X_1 < Y_1\}1\{X_2 < Y_1\}\right] - \theta^2 \\ &= P(X_1 < Y_1, X_2 < Y_1) - \theta^2 \stackrel{H_0}{=} \frac{1}{3} - \frac{1}{4} = \frac{1}{12}, \end{aligned}$$

pričom predposledná rovnosť v každom vzťahu plynie z toho, že máme 6 možností ako usporiadať  $X_1, Y_1$  a  $Y_2$  (respektíve  $X_1, X_2$  a  $Y_1$ ) a každá z nich je rovnako pravdepodobná.

Konečnosť  $\sigma_{11}^2$  je zrejماً, keďže ide o rozptyl  $1\{X_1 < Y_1\}$ .

Z vety 12 dostávame, že ak  $n/N \rightarrow p \in (0,1)$ , keď  $n + m \rightarrow \infty$ , tak

$$\sqrt{N} \left( U_{n,m} - \frac{1}{2} \right) \xrightarrow{D} N \left( 0, \frac{1}{12} \left( \frac{1}{p} + \frac{1}{1-p} \right) \right).$$

Veta o spojitej transformácii dá, že

$$\frac{\sqrt{N} \left( U_{n,m} - \frac{1}{2} \right)}{\sqrt{\frac{1}{12} \left( \frac{1}{p} + \frac{1}{1-p} \right)}} \xrightarrow{D} N(0, 1).$$

Vďaka Cramérovej-Sluckého vete môžeme výraz  $\frac{1}{p} + \frac{1}{1-p}$  v menovateli nahradiť výrazom  $\frac{N^2}{nm}$ , čím je konvergencia v distribúcii zachovaná, a teda

$$\frac{\sqrt{N} \left( U_{n,m} - \frac{1}{2} \right)}{\sqrt{\frac{N^2}{12nm}}} \xrightarrow{D} N(0, 1). \quad (1.5)$$

Výraz na ľavej strane (1.5) teraz prevedieme na tvar, v ktorom vystupuje  $W_X$ . Postupnosťou úprav s využitím tvrdenia 2 a vzťahu (1.4) dostávame, že

$$\begin{aligned} \frac{\sqrt{N} \left( U_{n,m} - \frac{1}{2} \right)}{\sqrt{\frac{N^2}{12nm}}} &= \frac{\sqrt{N} \left( \frac{W_{XY}}{nm} - \frac{1}{2} \right)}{\sqrt{\frac{N^2}{12nm}}} = \frac{\sqrt{N} \left( \frac{nm+n(n+1)-2W_X}{2nm} \right)}{\sqrt{\frac{N^2}{12nm}}} \\ &= -\frac{1}{\sqrt{N}} \frac{\left( \frac{W_X - \frac{n(n+m+1)}{2}}{nm} \right)}{\sqrt{\frac{1}{12nm}}} = -\frac{W_X - \frac{n(N+1)}{2}}{\sqrt{\frac{nmN}{12}}}. \end{aligned}$$

Z tvrdenia 1 a Cramérovej-Sluckého vety konečne dostávame, že ak  $n/N \rightarrow p \in (0,1)$ , keď  $n + m \rightarrow \infty$ , tak

$$\begin{aligned} \frac{W_X - \mathbf{E}_{H_0} W_X}{\sqrt{\text{var}_{H_0} W_X}} &= \frac{W_X - \mathbf{E}_{H_0} W_X}{\sqrt{\frac{nmN}{12}}} \frac{\sqrt{\frac{nmN}{12}}}{\sqrt{\text{var}_{H_0} W_X}} \\ &= \frac{W_X - \mathbf{E}_{H_0} W_X}{\sqrt{\frac{nmN}{12}}} \frac{\sqrt{\frac{nmN}{12}}}{\sqrt{\frac{nm(N+1)}{12}}} \xrightarrow{D} N(0,1). \end{aligned}$$

□

## 2. Dvojvýberový Wilcoxonov test so zhodami

V tejto kapitole opäť uvažujeme 2 nezávislé náhodné výbery  $\mathbf{X} = (X_1, \dots, X_n)$  a  $\mathbf{Y} = (Y_1, \dots, Y_m)$  z rozdelení s distribučnými funkciami  $F_X$  a  $F_Y$ , avšak už nepredpokladáme, že  $F_X$  a  $F_Y$  sú spojité. Už teda nemôžeme predpokladať, že pravdepodobnosť zhodných pozorovaní je 0, ako sme predpokladali v (1.1).

### 2.1 Zavedenie testovej štatistiky v prípade zhôd

V prípade existencie zhôd prestáva dávať definícia poradí z prvej kapitoly zmysel, pretože už nie je možné zoradiť náhodné veličiny v združenom výbere  $\mathbf{Z}$  od najmenej po najväčšiu tak, aby boli splnené ostré nerovnosti  $Z_{(1)} < Z_{(2)} < \dots < Z_{(N)}$ . Z tohto dôvodu prejdeme k definícii takzvaných priemerných poradí. Napriek nemožnosti usporiadať združený výber  $\mathbf{Z}$  tak, aby platili ostré nerovnosti, ho stále môžeme usporiadať tak, aby platili neostré nerovnosti  $Z_{(1)} \leq Z_{(2)} \leq \dots \leq Z_{(N)}$ . Uvažujme blok  $k$  zhodných pozorovaní:  $\dots < Z_{(i)} = Z_{(i+1)} = \dots = Z_{(i+k-1)} < \dots$ .  $\mathbf{Z}$  množiny  $1, \dots, N$  by mali tohto bloku zrejme prislúchať čísla  $i, i+1, \dots, i+k-1$ . Definície 2 však každej náhodnej veličine z tohto bloku priraduje všetky tieto čísla, zatiaľ čo vzťah (1.3) každej náhodnej veličine z tohto bloku priraduje číslo  $i+k-1$ . Priemerné poradie definujeme tak, aby priemerné poradie každej náhodnej veličiny z tohto bloku bolo aritmetickým priemerom čísel  $i, i+1, \dots, i+k-1$ , teda  $(2i+k-1)/2$ .

**Definícia 3.** Priemerné poradie náhodnej veličiny  $Z_i$ ,  $i \in \{1, \dots, N\}$ , v združenom výbere  $\mathbf{Z}$  definujeme ako

$$R_i^* = 1 + \sum_{j=1}^N 1\{Z_j < Z_i\} + \frac{1}{2} \sum_{\substack{j=1 \\ j \neq i}}^N 1\{Z_j = Z_i\}.$$

*Poznámka.* Usporiadánym výberom  $\mathbf{Z}_{(\cdot)}$  budeme stále rozumieť vzostupne usporiadaný výber  $\mathbf{Z}$ , len s neostrými nerovnosťami:  $Z_{(1)} \leq Z_{(2)} \leq \dots \leq Z_{(N)}$ .

Wilcoxonova testová štatistika je v prípade existencie zhôd definovaná ako suma priemerných poradí náhodných veličín z  $\mathbf{X}$  v združenom výbere  $\mathbf{Z}$ , teda

$$W_X^* = \sum_{i=1}^n R_i^*.$$

Rovnako ako v prípade bez zhôd svedčia proti  $H_0$  (1.2) príliš malé a zároveň príliš veľké hodnoty testovej štatistiky  $W_X^*$ .

Pokiaľ nie sú  $n$  a  $m$  príliš veľké, je možné určiť presné rozdelenie  $W_X^*$  za  $H_0$  pri danom  $\mathbf{Z}_{(\cdot)}$ . Toto presné rozdelenie sa dá odvodiť z toho, že za  $H_0$  je ľubovoľné usporiadanie  $\mathbf{Z}$  rovnako pravdepodobné. Nasledujúci príklad ilustruje, ako by sme postupovali, pre  $n$  a  $m$  dostatočne malé.

*Príklad.* Predpokladajme, že platí  $H_0$ , že  $n = m = 2$ , a že  $\mathbf{Z}_{(\cdot)} = (2, 7, 7, 10)$ . Vektor priemerných poradí je potom  $(1, \frac{5}{2}, \frac{5}{2}, 4)$ . Za tejto informácie je 12 možností

ako usporiadať  $(X_1, X_2, Y_1, Y_2)$ , pretože máme  $\binom{4}{2}$  možností ako vybrať 2 náhodné veličiny, ktoré sa budú rovnať a v každej z týchto možností máme 2 možnosti, ktorá zo zvyšných náhodných veličín bude najmenšia. Potom

$$\begin{aligned} P(R_1^* = 1, R_2^* = \frac{5}{2}) &= \frac{2}{12}, & P(R_1^* = \frac{5}{2}, R_2^* = 1) &= \frac{2}{12}, & P(R_1^* = 4, R_2^* = 1) &= \frac{1}{12}, \\ P(R_1^* = 1, R_2^* = 4) &= \frac{1}{12}, & P(R_1^* = \frac{5}{2}, R_2^* = \frac{5}{2}) &= \frac{2}{12}, & P(R_1^* = 4, R_2^* = \frac{5}{2}) &= \frac{2}{12}, \\ & & P(R_1^* = \frac{5}{2}, R_2^* = 4) &= \frac{2}{12}. \end{aligned}$$

Vidíme, že  $W_X^*$  môže nadobúdať hodnoty  $\frac{7}{2}$ , 5 a  $\frac{13}{2}$ . Každú z nich s pravdepodobnosťou  $\frac{1}{3}$ .

V prípade, že  $n$  a  $m$  sú veľké, využijeme na zvolenie kritického oboru vetu 8 zo sekcie 2.2. Aby sme dostali test s asymptotickou hladinou  $\alpha$ , zvolíme kritický obor tak, že

$$H_0 \text{ zamietame} \Leftrightarrow \frac{|W_X^* - \frac{n(N+1)}{2}|}{\sqrt{\frac{nm(N+1)}{12} - \frac{nm\kappa}{(N-1)N}}} \geq u_{1-\alpha/2},$$

kde

$$\kappa = \sum_z \frac{k_z^3 - k_z}{12},$$

pričom  $k_z$  značí počet, koľkokrát sa medzi hodnotami  $Z_1, \dots, Z_N$  vyskytla hodnota  $z$ . Suma  $\sum_z$  značí sčítanie cez všetky rôzne hodnoty  $Z_1, \dots, Z_N$ .

Pre odvodenie asymptotického rozdelenia testovej štatistiky  $W_X^*$  budeme postupovať analogicky ako v prvej kapitole.

## 2.2 Odvodenie asymptotického rozdelenia testovej štatistiky v prípade zhôd

V tejto sekcii opäť použijeme vetu 12 zo sekcie 5.1 na odvodenie asymptotického rozdelenia  $W_X^*$ . Definujme

$$W_{XY}^* = \sum_{i=1}^n \sum_{j=1}^m 1\{X_i < Y_j\} + \frac{1}{2} 1\{X_i = Y_j\}.$$

Podobne ako v prípade bez zhôd spája Mannovu-Whitneyovu štatistiku  $W_{XY}^*$  s Wilcoxonovou štatistikou deterministický vzťah.

**Tvrdenie 4.** Platí, že  $W_{XY}^* = nm + \frac{n(n+1)}{2} - W_X^*$ .

*Dôkaz.* Podľa definície 3, pre  $i \in \{1, \dots, n\}$  máme, že

$$\begin{aligned} R_i^* &= 1 + \sum_{j=1}^n 1\{X_j < X_i\} + \sum_{j=1}^m 1\{Y_j < X_i\} \\ &\quad + \frac{1}{2} \sum_{\substack{j=1 \\ j \neq i}}^n 1\{X_j = X_i\} + \frac{1}{2} \sum_{j=1}^m 1\{Y_j = X_i\}. \end{aligned} \tag{2.1}$$

Označme  $S_i^*$  priemerné poradie náhodnej veličiny  $X_i$  v  $\mathbf{X}$ . S využitím (2.1) je možné upraviť súčet  $W_{XY}^*$  a  $W_X^*$  nasledovne:

$$\begin{aligned}
W_{XY}^* + W_X^* &= \sum_{i=1}^n \sum_{j=1}^m 1\{X_i < Y_j\} + \frac{1}{2} 1\{X_i = Y_j\} + \sum_{i=1}^n R_i^* \\
&= \sum_{i=1}^n \sum_{j=1}^m 1\{X_i < Y_j\} + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m 1\{X_i = Y_j\} \\
&\quad + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m 1\{Y_j = X_i\} + \sum_{i=1}^n \sum_{j=1}^m 1\{Y_j < X_i\} \\
&\quad + \sum_{i=1}^n 1 + \sum_{i=1}^n \sum_{j=1}^n 1\{X_j < X_i\} + \frac{1}{2} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n 1\{X_j = X_i\} \\
&= nm + \sum_{i=1}^n \left( 1 + \sum_{j=1}^n 1\{X_j < X_i\} + \frac{1}{2} \sum_{\substack{j=1 \\ j \neq i}}^n 1\{X_j = X_i\} \right) \\
&= nm + \sum_{i=1}^n S_i^* = nm + \frac{n(n+1)}{2}.
\end{aligned}$$

□

**Veta 5.** Ak platí  $H_0$  a  $n/N \rightarrow p \in (0,1)$ , keď  $n + m \rightarrow \infty$ , tak

$$\sqrt{N} \left( \frac{W_{XY}^*}{nm} - \frac{1}{2} \right) \xrightarrow{D} N \left( 0, \xi \frac{1}{p(1-p)} \right),$$

kde  $\xi = \text{var}(F(Z_i) + F(Z_i-)) / 4$ , pričom  $F$  značí spoločnú distribučnú funkciu pre náhodné veličiny z  $\mathbf{X}$  aj z  $\mathbf{Y}$ .

*Dôkaz.* K dôkazu použijeme vetu 12, rovnako ako sme ju použili v dôkaze vety 3. Definujme jadro  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$  predpisom

$$h(x,y) = \begin{cases} 1, & \text{ak } x < y, \\ 1/2, & \text{ak } x = y, \\ 0, & \text{inak.} \end{cases}$$

Potom korešpondujúca U-štatistika pre jadro  $h$  je tvaru

$$U_{n,m} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m h(X_i, Y_j) = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m 1\{X_i < Y_j\} + \frac{1}{2} 1\{X_i = Y_j\}.$$

Zrejme platí, že

$$U_{n,m} = \frac{W_{XY}^*}{nm}.$$

Teraz spočítame  $\theta$  z vety 12 a overíme konečnosť  $\sigma_{11}^2$ .

Platí, že

$$\begin{aligned}
\theta &= \mathbb{E}[h(X_i, Y_j)] = \mathbb{E}\left[1\{X_i < Y_j\} + \frac{1}{2} 1\{X_i = Y_j\}\right] \\
&= \mathbb{P}(X_i < Y_j) + \frac{1}{2} \mathbb{P}(X_i = Y_j) \stackrel{H_0}{=} \frac{1}{2},
\end{aligned}$$

pričom v kroku označenom  $\stackrel{H_0}{=}$  používame predpoklad, že platí  $H_0$ , a teda náhodné veličiny  $X_i$  a  $Y_j$  pochádzajú z toho istého rozdelenia.

Ďalej,

$$\begin{aligned}\sigma_{11}^2 &= \text{cov} \left( 1\{X_1 < Y_1\} + \frac{1}{2}1\{X_1 = Y_1\}, 1\{X_1 < Y_1\} + \frac{1}{2}1\{X_1 = Y_1\} \right) \\ &= \text{E} \left[ 1\{X_1 < Y_1\} + \frac{1}{4}1\{X_1 = Y_1\} \right] - \theta^2 < \infty.\end{aligned}$$

Pre  $\sigma_{10}^2$  a  $\sigma_{01}^2$  dostávame nasledujúce vzťahy:

$$\begin{aligned}\sigma_{10}^2 &= \text{cov} \left( 1\{X_1 < Y_1\} + \frac{1}{2}1\{X_1 = Y_1\}, 1\{X_1 < Y_2\} + \frac{1}{2}1\{X_1 = Y_2\} \right) \\ &= \text{E} \left[ 1\{X_1 < Y_1\}1\{X_1 < Y_2\} + \frac{1}{2}1\{X_1 < Y_1\}1\{X_1 = Y_2\} \right. \\ &\quad \left. + \frac{1}{2}1\{X_1 < Y_2\}1\{X_1 = Y_1\} + \frac{1}{4}1\{X_1 = Y_1\}1\{X_1 = Y_2\} \right] - \theta^2,\end{aligned}\tag{2.2}$$

$$\begin{aligned}\sigma_{01}^2 &= \text{cov} \left( 1\{X_1 < Y_1\} + \frac{1}{2}1\{X_1 = Y_1\}, 1\{X_2 < Y_1\} + \frac{1}{2}1\{X_2 = Y_1\} \right) \\ &= \text{E} \left[ 1\{X_1 < Y_1\}1\{X_2 < Y_1\} + \frac{1}{2}1\{X_1 < Y_1\}1\{X_2 = Y_1\} \right. \\ &\quad \left. + \frac{1}{2}1\{X_2 < Y_1\}1\{X_1 = Y_1\} + \frac{1}{4}1\{X_1 = Y_1\}1\{X_2 = Y_1\} \right] - \theta^2.\end{aligned}\tag{2.3}$$

Pokračujme ďalej s úpravou  $\sigma_{01}^2$  s využitím podmienenej strednej hodnoty. Špeciálne, zamerajme sa na člen  $\text{E} \left[ 1\{X_1 < Y_1\}1\{X_2 < Y_1\} \right]$ . Postupnosťou úprav dostávame, že

$$\begin{aligned}\text{E} \left[ 1\{X_1 < Y_1\}1\{X_2 < Y_1\} \right] &= \text{E} \left[ \text{E} \left[ 1\{X_1 < Y_1\}1\{X_2 < Y_1\} \mid Y_1 \right] \right] \\ &= \int_{\mathbb{R}} \text{E} \left[ 1\{X_1 < Y_1\}1\{X_2 < Y_1\} \mid Y_1 = y \right] d\mathbf{P}_{Y_1}(y) \\ &= \int_{\mathbb{R}} \text{E} \left[ 1\{X_1 < y\}1\{X_2 < y\} \right] d\mathbf{P}_{Y_1}(y) \\ &= \int_{\mathbb{R}} \text{E} \left[ 1\{X_1 < y\} \right] \text{E} \left[ 1\{X_2 < y\} \right] d\mathbf{P}_{Y_1}(y) \\ &= \int_{\mathbb{R}} F_X(y-)F_X(y-) d\mathbf{P}_{Y_1}(y) = \text{E} \left[ F_X(Y_1-)F_X(Y_1-) \right].\end{aligned}\tag{2.4}$$

Analogicky, pre zvyšné členy vystupujúce v  $\sigma_{01}^2$  dostávame, že

$$\begin{aligned}\text{E} \left[ 1\{X_1 < Y_1\}1\{X_2 = Y_1\} \right] &= \int_{\mathbb{R}} \text{E} \left[ 1\{X_1 < y\} \right] \text{E} \left[ 1\{X_2 = y\} \right] d\mathbf{P}_{Y_1}(y) \\ &= \int_{\mathbb{R}} F_X(y-) \left( F_X(y) - F_X(y-) \right) d\mathbf{P}_{Y_1}(y) \\ &= \text{E} \left[ F_X(Y_1-) \left( F_X(Y_1) - F_X(Y_1-) \right) \right],\end{aligned}\tag{2.5}$$

$$\begin{aligned}
\mathbb{E} \left[ 1\{X_1 = Y_1\} 1\{X_2 = Y_1\} \right] &= \int_{\mathbb{R}} \mathbb{E} \left[ 1\{X_1 = y\} \right] \mathbb{E} \left[ 1\{X_2 = y\} \right] d\mathbf{P}_{Y_1}(y) \\
&= \int_{\mathbb{R}} \left( F_X(y) - F_X(y-) \right)^2 d\mathbf{P}_{Y_1}(y) \\
&= \mathbb{E} \left[ \left( F_X(Y_1) - F_X(Y_1-) \right)^2 \right]. \tag{2.6}
\end{aligned}$$

S využitím podmienenej strednej hodnoty navyše rovnako (podmienением náhodnou veličinou  $Y_1$ ) upravíme  $\theta$ :

$$\begin{aligned}
\theta &= \mathbb{E} \left[ 1\{X_1 < Y_1\} \right] + \frac{1}{2} \mathbb{E} \left[ 1\{X_1 = Y_1\} \right] \\
&= \int_{\mathbb{R}} \mathbb{E} \left[ 1\{X_1 < y\} \right] d\mathbf{P}_{Y_1}(y) + \frac{1}{2} \int_{\mathbb{R}} \mathbb{E} \left[ 1\{X_1 = y\} \right] d\mathbf{P}_{Y_1}(y) \\
&= \int_{\mathbb{R}} F_X(y-) d\mathbf{P}_{Y_1}(y) + \frac{1}{2} \int_{\mathbb{R}} F_X(y) - F_X(y-) d\mathbf{P}_{Y_1}(y) \\
&= \mathbb{E} \left[ F_X(Y_1-) \right] + \frac{1}{2} \mathbb{E} \left[ F_X(Y_1) - F_X(Y_1-) \right] \\
&= \frac{1}{2} \mathbb{E} \left[ F_X(Y_1) + F_X(Y_1-) \right]. \tag{2.7}
\end{aligned}$$

Konečne, dosadením rovností (2.4), (2.5), (2.6) a (2.7) do (2.3) dostávame, že

$$\begin{aligned}
\sigma_{01}^2 &= \mathbb{E} \left[ F_X(Y_1-) F_X(Y_1-) \right] + \mathbb{E} \left[ F_X(Y_1-) \left( F_X(Y_1) - F_X(Y_1-) \right) \right] \\
&\quad + \frac{1}{4} \mathbb{E} \left[ \left( F_X(Y_1) - F_X(Y_1-) \right)^2 \right] - \left( \frac{1}{2} \mathbb{E} \left[ F_X(Y_1) + F_X(Y_1-) \right] \right)^2 \\
&= \frac{1}{4} \mathbb{E} \left[ \left( F_X(Y_1) + F_X(Y_1-) \right)^2 \right] - \frac{1}{4} \left( \mathbb{E} \left[ F_X(Y_1) + F_X(Y_1-) \right] \right)^2 \\
&= \frac{1}{4} \text{var} \left( F_X(Y_1) + F_X(Y_1-) \right).
\end{aligned}$$

Analogicky, podmienением náhodnou veličinou  $X_1$  dostávame, že

$$\sigma_{10}^2 = \frac{1}{4} \text{var} \left( F_Y(X_1) + F_Y(X_1-) \right).$$

Táto rovnosť je podrobne odvodená v kapitole Dodatky.

Predpokladáme, že platí  $H_0$ . Potom platí, že  $\mathbf{Z}$  je náhodný výber z rozdelenia s distribučnou funkciou  $F$ , kde  $F = F_X = F_Y$ . Ďalej platí, že  $\sigma_{10}^2 = \sigma_{01}^2$ . Označme túto spoločnú hodnotu  $\xi$ .

□

Problémom však zatiaľ ostáva, že nepoznáme spoločnú distribučnú funkciu  $F$  z tvrdenia vety 5, a teda ani  $\xi$ . Z tohto dôvodu prejdeme k odhadu  $\xi$  pomocou výberového rozptylu a empirickej distribučnej funkcie. Zavedme empirickú distribučnú funkciu.

**Definícia 4.** *Funkciu*

$$\hat{F}(x) = \frac{1}{N} \sum_{i=1}^N 1\{Z_i \leq x\} = \frac{1}{N} \left( \sum_{i=1}^n 1\{X_i \leq x\} + \sum_{j=1}^m 1\{Y_j \leq x\} \right)$$

*budeme nazývať empirickou distribučnou funkciou  $\mathbf{Z}$ .*

Teraz s využitím  $\widehat{F}$  zavedme odhad  $\xi$  ako

$$\widehat{\xi}_N = \frac{1}{4} \frac{1}{N-1} \sum_{i=1}^N (U_i - \overline{U}_N)^2, \text{ kde } U_i = \widehat{F}(Z_i) + \widehat{F}(Z_{i-}).$$

Vidíme, že  $\widehat{\xi}_N$  sme skonštruovali ako zloženie výberového rozptylu a empirickej distribučnej funkcie, teda dvoch konzistentných odhadov.

**Tvrdenie 6.** *Ak platí  $H_0$  a  $N = n + m \rightarrow \infty$ , tak  $\widehat{\xi}_N \xrightarrow{P} \xi$ .*

*Dôkaz.* K dôkazu tohto tvrdenia použijeme Glivenkovu-Cantelliho vetu, ktorá hovorí, že empirická distribučná funkcia konverguje rovnomerne v pravdepodobnosti k skutočnej distribučnej funkcii, z ktorej pochádza náhodný výber, teda, že

$$\sup_{t \in \mathbb{R}} |\widehat{F}(t) - F(t)| \xrightarrow{P} 0.$$

Dôkaz Glivenkovej-Cantelliho vety je uvedený v knihe van der Vaart (1998, strana 266).

Štandardná úprava výberového rozptylu dá, že

$$4\widehat{\xi}_N = \frac{1}{N-1} \sum_{i=1}^N (U_i - \overline{U}_N)^2 = \frac{1}{N-1} \sum_{i=1}^N U_i^2 - \frac{N}{N-1} \overline{U}_N^2. \quad (2.8)$$

Najprv vyšetříme konvergenciu člena  $\frac{N}{N-1} \overline{U}_N^2$ . Platí, že

$$\frac{1}{N} \sum_{i=1}^N \widehat{F}(Z_i) = \frac{1}{N} \sum_{i=1}^N \left( \widehat{F}(Z_i) - F(Z_i) \right) + \frac{1}{N} \sum_{i=1}^N F(Z_i) \xrightarrow{P} \mathbf{E} F(Z_i),$$

pretože

$$\left| \frac{1}{N} \sum_{i=1}^N \left( \widehat{F}(Z_i) - F(Z_i) \right) \right| \leq \sup_{t \in \mathbb{R}} |\widehat{F}(t) - F(t)| \xrightarrow{P} 0$$

a podľa zákona veľkých čísel

$$\frac{1}{N} \sum_{i=1}^N F(Z_i) \xrightarrow{P} \mathbf{E} F(Z_i).$$

Potom pre  $\frac{N}{N-1} \overline{U}_N^2$  dostávame, že

$$\begin{aligned} \frac{N}{N-1} \overline{U}_N^2 &= \frac{N}{N-1} \left( \frac{1}{N} \sum_{i=1}^N \widehat{F}(Z_i) + \frac{1}{N} \sum_{i=1}^N \widehat{F}(Z_{i-}) \right)^2 \\ &\xrightarrow{P} \left( \mathbf{E} \left[ F(Z_i) + F(Z_{i-}) \right] \right)^2. \end{aligned} \quad (2.9)$$

Teraz vyšetříme konvergenciu člena  $\frac{1}{N-1} \sum_{i=1}^N U_i^2$ .

$$\frac{1}{N-1} \sum_{i=1}^N U_i^2 = \frac{1}{N-1} \sum_{i=1}^N \widehat{F}(Z_i)^2 + \frac{2}{N-1} \sum_{i=1}^N \widehat{F}(Z_i) \widehat{F}(Z_{i-}) + \frac{1}{N-1} \sum_{i=1}^N \widehat{F}(Z_{i-})^2$$



Platí, že

$$\begin{aligned} \frac{1}{N-1} \sum_{i=1}^N \hat{F}(Z_i)^2 &= \frac{1}{N-1} \sum_{i=1}^N \left( \hat{F}(Z_i) - F(Z_i) \right)^2 \\ &\quad + \frac{2}{N-1} \sum_{i=1}^N F(Z_i) \left( \hat{F}(Z_i) - F(Z_i) \right) \\ &\quad + \frac{1}{N-1} \sum_{i=1}^N F(Z_i)^2 \xrightarrow{P} \mathbf{E} F(Z_i)^2, \end{aligned}$$

pretože

$$\begin{aligned} \left| \frac{1}{N-1} \sum_{i=1}^N \left( \hat{F}(Z_i) - F(Z_i) \right)^2 \right| &\leq \frac{N}{N-1} \left( \sup_{t \in \mathbb{R}} |\hat{F}(t) - F(t)| \right)^2 \xrightarrow{P} 0, \\ \left| \frac{2}{N-1} \sum_{i=1}^N F(Z_i) \left( \hat{F}(Z_i) - F(Z_i) \right) \right| &\leq \frac{2}{N-1} \sum_{i=1}^N F(Z_i) \sup_{t \in \mathbb{R}} |\hat{F}(t) - F(t)| \\ &\leq 2 \frac{N}{N-1} \sup_{t \in \mathbb{R}} |\hat{F}(t) - F(t)| \xrightarrow{P} 0 \end{aligned}$$

a podľa zákona veľkých čísel

$$\frac{1}{N-1} \sum_{i=1}^N F(Z_i)^2 \xrightarrow{P} \mathbf{E} F(Z_i)^2.$$

Analogicky,

$$\frac{1}{N-1} \sum_{i=1}^N \hat{F}(Z_{i-})^2 \xrightarrow{P} \mathbf{E} F(Z_{i-})^2.$$

Zostáva sa pozrieť na člen  $\frac{2}{N-1} \sum_{i=1}^N \hat{F}(Z_i) \hat{F}(Z_{i-})$ , pre ktorý platí, že

$$\begin{aligned} \frac{2}{N-1} \sum_{i=1}^N \hat{F}(Z_i) \hat{F}(Z_{i-}) &= \frac{2}{N-1} \sum_{i=1}^N \left( \hat{F}(Z_i) \hat{F}(Z_{i-}) - \hat{F}(Z_i) F(Z_{i-}) \right) \\ &\quad + \frac{2}{N-1} \sum_{i=1}^N \left( \hat{F}(Z_i) F(Z_{i-}) - F(Z_i) F(Z_{i-}) \right) \\ &\quad + \frac{2}{N-1} \sum_{i=1}^N F(Z_i) F(Z_{i-}) \xrightarrow{P} 2 \mathbf{E} F(Z_i) F(Z_{i-}), \end{aligned}$$

pretože

$$\begin{aligned} &\left| \frac{2}{N-1} \sum_{i=1}^N \left( \hat{F}(Z_i) \hat{F}(Z_{i-}) - \hat{F}(Z_i) F(Z_{i-}) \right) \right| \\ &\leq \frac{2}{N-1} \sum_{i=1}^N \hat{F}(Z_i) \sup_{t \in \mathbb{R}} |\hat{F}(t-) - F(t-)| \xrightarrow{P} 0, \\ &\left| \frac{2}{N-1} \sum_{i=1}^N \left( \hat{F}(Z_i) F(Z_{i-}) - F(Z_i) F(Z_{i-}) \right) \right| \\ &\leq \frac{2}{N-1} \sum_{i=1}^N F(Z_{i-}) \sup_{t \in \mathbb{R}} |\hat{F}(t) - F(t)| \xrightarrow{P} 0 \end{aligned}$$

a podľa zákona veľkých čísel

$$\frac{2}{N-1} \sum_{i=1}^N F(Z_i)F(Z_{i-}) \xrightarrow{P} 2 \mathbf{E} F(Z_i)F(Z_{i-}).$$

Pre  $\frac{1}{N-1} \sum_{i=1}^N U_i^2$  teda dostávame, že

$$\begin{aligned} \frac{1}{N-1} \sum_{i=1}^N U_i^2 &\xrightarrow{P} \mathbf{E} F(Z_i)^2 + 2 \mathbf{E} F(Z_i)F(Z_{i-}) + \mathbf{E} F(Z_{i-})^2 \\ &= \mathbf{E} \left( F(Z_i) + F(Z_{i-}) \right)^2 \end{aligned} \quad (2.10)$$

Konečne, z rozpisu (2.8) pre  $4\widehat{\xi}_N$  dostávame pomocou (2.9) a (2.10), že

$$4\widehat{\xi}_N \xrightarrow{P} \mathbf{E} \left( F(Z_i) + F(Z_{i-}) \right)^2 - \left( \mathbf{E} \left[ F(Z_i) + F(Z_{i-}) \right] \right)^2 = \text{var} \left( F(Z_i) + F(Z_{i-}) \right),$$

a teda  $\widehat{\xi}_N \xrightarrow{P} \xi$ . □

**Tvrdenie 7.** Platí, že  $U_i = \frac{1}{N}(2R_i^* - 1)$  a  $\overline{U}_N = 1$ .

*Dôkaz.* Nech  $l$  je počet náhodných veličín zo  $\mathbf{Z}$  menších ako  $Z_i$  a  $k$  je počet náhodných veličín zo  $\mathbf{Z}$  menších alebo rovných  $Z_i$ . Náhodná veličina  $Z_i$  je teda rovná  $k - 1$  ďalším náhodným veličinám zo  $\mathbf{Z}$ . Potom

$$\begin{aligned} U_i &= \frac{1}{N} \sum_{j=1}^N 1\{Z_j \leq Z_i\} + \frac{1}{N} \sum_{j=1}^N 1\{Z_j < Z_i\} = \frac{1}{N}(l + k) + \frac{1}{N}l \\ &= \frac{1}{N} \left( 2(1 + l + \frac{1}{2}(k - 1)) - 1 \right) = \frac{1}{N}(2R_i^* - 1). \end{aligned}$$

Pre výberový priemer  $U_i$  potom platí, že

$$\overline{U}_N = \frac{1}{N} \sum_{i=1}^N \frac{1}{N}(2R_i^* - 1) = \frac{1}{N^2} \left( 2 \frac{N(N+1)}{2} - N \right) = 1. \quad \square$$

**Veta 8.** Ak platí  $H_0$  a  $n/N \rightarrow p \in (0,1)$ , keď  $n + m \rightarrow \infty$ , tak

$$\frac{W_X^* - \frac{n(N+1)}{2}}{\sqrt{\frac{nm(N+1)}{12} - \frac{nm\kappa}{(N-1)N}}} \xrightarrow{D} N(0,1),$$

kde

$$\kappa = \sum_z \frac{k_z^3 - k_z}{12},$$

pričom  $k_z$  značí počet, koľkokrát sa medzi hodnotami  $Z_1, \dots, Z_N$  vyskytla hodnota  $z$ . Suma  $\sum_z$  značí sčítanie cez všetky rôzne hodnoty  $Z_1, \dots, Z_N$ .

*Dôkaz.* Vďaka Cramérovej-Sluckého vete a tvrdeniu 6 dostávame z vety 5, že

$$\frac{\sqrt{N} \left( \frac{W_{XY}^*}{nm} - \frac{1}{2} \right)}{\sqrt{\widehat{\xi}_N \frac{N^2}{nm}}} \xrightarrow{D} N(0, 1). \quad (2.11)$$

Z lemy 10 z kapitoly 3 vieme, že

$$\sum_{i=1}^N R_i^{*2} = \frac{N(N+1)(2N+1)}{6} - \kappa.$$

S využitím tejto rovnosti a tvrdenia 7 sa  $\widehat{\xi}_N$  zjednoduší nasledovne:

$$\begin{aligned} \widehat{\xi}_N &= \frac{1}{4} \frac{1}{N-1} \sum_{i=1}^N (U_i - \overline{U_N})^2 \\ &= \frac{1}{4} \frac{1}{N-1} \left( \sum_{i=1}^N \frac{1}{N^2} (2R_i^* - 1)^2 - 2 \sum_{i=1}^N \frac{1}{N} (2R_i^* - 1) + N \right) \\ &= \frac{1}{4} \frac{1}{N-1} \left( \sum_{i=1}^N \frac{1}{N^2} (2R_i^* - 1)^2 - \frac{2}{N} N^2 + N \right) \\ &= \frac{1}{4} \frac{1}{N-1} \left( \frac{4}{N^2} \sum_{i=1}^N R_i^{*2} - \frac{4}{N^2} \sum_{i=1}^N R_i^* + \frac{1}{N} - N \right) \\ &= \frac{1}{4} \frac{1}{N-1} \left( \frac{4}{N^2} \left( \frac{N(N+1)(2N+1)}{6} - \kappa \right) - \frac{4}{N^2} \frac{N(N+1)}{2} + \frac{1}{N} - N \right) \\ &= \frac{1}{4} \frac{1}{N-1} \left( \frac{4(N+1)(2N+1) - 12(N+1) + 6 - 6N^2}{6N} - \frac{4}{N^2} \kappa \right) \\ &= \frac{1}{4} \frac{1}{N-1} \frac{2N^2 - 2}{6N} - \frac{1}{N-1} \frac{1}{N^2} \kappa \\ &= \frac{1}{12} \frac{N+1}{N} - \frac{1}{N-1} \frac{1}{N^2} \kappa. \end{aligned} \quad (2.12)$$

Výraz na ľavej strane (2.11) už len prevedieme na tvar, v ktorom vystupuje  $W_X^*$ . Postupnosťou úprav s využitím tvrdenia 4 a vzťahu (2.12) dostávame, že

$$\begin{aligned} \frac{\sqrt{N} \left( \frac{W_{XY}^*}{nm} - \frac{1}{2} \right)}{\sqrt{\widehat{\xi}_N \frac{N^2}{nm}}} &= \frac{\sqrt{N} \left( \frac{W_{XY}^*}{nm} - \frac{1}{2} \right)}{N \sqrt{\frac{1}{nm} \sqrt{\frac{N+1}{12N}} - \frac{1}{(N-1)N^2} \kappa}} = \frac{\sqrt{N} \left( \frac{nm + \frac{n(n+1)}{2} - W_X^*}{nm} - \frac{1}{2} \right)}{N \sqrt{\frac{1}{nm} \sqrt{\frac{N+1}{12N}} - \frac{1}{(N-1)N^2} \kappa}} \\ &= \frac{\sqrt{N} \frac{1}{nm} \left( \frac{n(n+m+1)}{2} - W_X^* \right)}{N \sqrt{\frac{1}{nm} \sqrt{\frac{N+1}{12N}} - \frac{1}{(N-1)N^2} \kappa}} = - \frac{W_X^* - \frac{n(n+m+1)}{2}}{\sqrt{nm} \sqrt{\frac{N+1}{12} - \frac{1}{(N-1)N} \kappa}} \\ &= - \frac{W_X^* - \frac{n(N+1)}{2}}{\sqrt{\frac{nm(N+1)}{12} - \frac{nm\kappa}{(N-1)N}}}. \end{aligned}$$

Alternatívny dôkaz je uvedený v knihe Lehmann a D'Abbrera (1975, strana 355).  $\square$

### 3. Súvislosť asymptotického a podmieneného rozptylu $W_X^*$

V tejto kapitole odvodíme

$$E [W_X^* | \mathbf{Z}_{(\cdot)}] \quad \text{a} \quad \text{var} [W_X^* | \mathbf{Z}_{(\cdot)}],$$

za platnosti  $H_0$  (1.2), čím ukážeme súvislosť s asymptotickým výsledkom z vety 8. K tomu sa nám budú hodiť takzvané lineárne poradové štatistiky a náhodné poradia.

Nech  $\mathbf{X} = (X_1, \dots, X_n)$  a  $\mathbf{Y} = (Y_1, \dots, Y_m)$  sú nezávislé náhodné výbery z rozdelení s distribučnými funkciami  $F_X$  a  $F_Y$ . Zatiaľ predpokladajme, že  $F_X$  a  $F_Y$  sú spojité. Opäť označme  $\mathbf{Z} = (Z_1, \dots, Z_N) = (X_1, \dots, X_n, Y_1, \dots, Y_m)$ .

Nech  $R_i$  je poradie náhodnej veličiny  $Z_i$  v  $\mathbf{Z}$ . Ďalej, nech pre  $i \in \{1, \dots, N\}$  je  $c_i \in \mathbb{R}$  reálny koeficient a  $a : \{1, \dots, N\} \rightarrow \mathbb{R}$  je takzvaná skórová funkcia. Položme

$$S = \sum_{i=1}^N c_i a(R_i)$$

a označme

$$\begin{aligned} \bar{a} &= \frac{1}{N} \sum_{i=1}^N a(i), & \sigma_a^2 &= \frac{1}{N} \sum_{i=1}^N [a(i) - \bar{a}]^2, \\ \bar{c} &= \frac{1}{N} \sum_{i=1}^N c_i, & \sigma_c^2 &= \frac{1}{N} \sum_{i=1}^N [c_i - \bar{c}]^2. \end{aligned}$$

**Veta 9.** *Nech  $\mathbf{X}$  a  $\mathbf{Y}$  sú nezávislé náhodné výbery zo spojitého rozdelenia s distribučnou funkciou  $F = F_X = F_Y$ . Potom platí, že*

$$E S = N\bar{a}\bar{c} \quad \text{a} \quad \text{var} S = \frac{N^2}{N-1} \sigma_a^2 \sigma_c^2.$$

*Dôkaz.* Dôkaz plynie z toho, že náhodné veličiny  $R_i$  nadobúdajú každú z hodnôt  $1, \dots, N$  s pravdepodobnosťou  $1/N$ . Podrobný dôkaz je uvedený v knihe Anděl (2007a, strana 236). □

V tejto chvíli upustíme od predpokladu spojitosti  $F_X$  a  $F_Y$ . Už teda nemôžeme predpokladať, že pravdepodobnosť zhodných pozorovaní je 0, ako sme predpokladali v (1.1). Vieme, že v prípade existencie zhôd nedáva definícia poradí z prvej kapitoly zmysel, pretože už nie je možné zoradiť náhodné veličiny v združenom výbere  $\mathbf{Z}$  od najmensej po najväčšiu tak, aby boli splnené ostré nerovnosti  $Z_{(1)} < Z_{(2)} < \dots < Z_{(N)}$ . Z tohto dôvodu prejdeme k takzvaných náhodným poradiam. Napriek nemožnosti usporiadať združený výber  $\mathbf{Z}$  tak, aby platili ostré nerovnosti, ho stále môžeme usporiadať tak, aby platili neostré nerovnosti  $Z_{(1)} \leq Z_{(2)} \leq \dots \leq Z_{(N)}$ . Uvažujme blok  $k$  zhodných pozorovaní:  $\dots < Z_{(i)} = Z_{(i+1)} = \dots = Z_{(i+k-1)} < \dots$ . Náhodné poradia náhodných veličín z tohto bloku budú dané náhodne vybranou permutáciou množiny

$\{i, i+1, \dots, i+k-1\}$ . Za predpokladu, že  $F_X = F_Y$  náhodné poradie  $(\widetilde{R}_i)$  stále splňa, že je rovnomerne rozdelené na množine  $\{1, \dots, N\}$ , a teda veta 9 platí aj s náhodnými poradiami, a to aj v prípade, že spojitost  $F_X$  a  $F_Y$  nepredpokladáme. Využijeme ju aby sme spočítali strednú hodnotu a rozptyl štatistiky  $W_X^*$  za podmienky, že poznáme  $\mathbf{Z}_{(\cdot)}$ .

Položme

$$c_i = \begin{cases} 1, & \text{ak } i = 1, \dots, n, \\ 0, & \text{ak } i = n+1, \dots, N. \end{cases}$$

Označme  $z_{(1)}, \dots, z_{(N)}$  vzostupne usporiadanú realizáciu združeného výberu  $\mathbf{Z}$  a definujme

$$a(i) = 1 + \sum_{j=1}^N 1[z_{(j)} < z_{(i)}] + \frac{1}{2} \sum_{\substack{j=1 \\ j \neq i}}^N 1[z_{(j)} = z_{(i)}].$$

Pre takto zvolené konštanty  $c_i$  a skórovú funkciu  $a$  platí, že  $W_X^* = \sum_{i=1}^N c_i a(\widetilde{R}_i)$ .

Zrejme platí, že

$$\bar{a} = \frac{N+1}{2} \quad \text{a} \quad \bar{c} = \frac{n}{N}. \quad (3.1)$$

Ďalej,

$$\begin{aligned} \sigma_c^2 &= \frac{1}{N} \left( \sum_{i=1}^n \left(1 - \frac{n}{N}\right)^2 + \sum_{i=n+1}^N \left(-\frac{n}{N}\right)^2 \right) \\ &= \frac{1}{N} \left( n \left(1 - \frac{2n}{N} + \frac{n^2}{N^2}\right) + m \frac{n^2}{N^2} \right) \\ &= \frac{1}{N} \left( n - \frac{2n^2}{N} + \frac{n^2}{N} \right) = \frac{nm}{N^2}. \end{aligned} \quad (3.2)$$

Pozrime sa na  $\sigma_a^2$ :

$$\sigma_a^2 = \frac{1}{N} \sum_{i=1}^N [a(i) - \bar{a}]^2 = \frac{1}{N} \sum_{i=1}^N (a(i))^2 - \bar{a}^2. \quad (3.3)$$

Vidíme, že sa stačí zamerať na  $\sum_{i=1}^N (a(i))^2$ . Uvedomme si, že ide vlastne o sumu druhých mocnín prvkov postupnosti, ktorá vznikne z postupnosti prirodzených čísel  $1, \dots, N$  tak, že túto postupnosť rozdelíme do  $r$  blokov, kde  $r$  je počet rôznych hodnôt medzi hodnotami  $z_1, \dots, z_N$ , a každé číslo v danom bloku nahradíme aritmetickým priemerom daného bloku. S touto sumou nám pomôže nasledujúca lema.

**Lema 10.** *Nech  $(a_1^{(1)}, \dots, a_{s_1}^{(1)}, \dots, a_1^{(r)}, \dots, a_{s_r}^{(r)}) = (1, \dots, N)$ , kde  $N \in \mathbb{N}$  a  $r \in \{1, \dots, N\}$ . Pre  $i \in \{1, \dots, r\}$  a  $j \in \{1, \dots, s_i\}$  položme*

$$b_j^{(i)} = \frac{1}{s_i} \sum_{k=1}^{s_i} a_k^{(i)}.$$

Potom

$$\sum_{i=1}^r \sum_{j=1}^{s_i} (b_j^{(i)})^2 = \frac{N(N+1)(2N+1)}{6} - \kappa,$$

kde  $\kappa = \sum_{i=1}^r \frac{s_i^3 - s_i}{12}$ .

*Dôkaz.* Pre  $i \in \{1, \dots, r\}$  platí, že

$$\sum_{k=1}^{s_i} (a_k^{(i)})^2 = \sum_{k=0}^{s_i-1} (a_1^{(i)} + k)^2 = s_i (a_1^{(i)})^2 + a_1^{(i)} (s_i - 1) s_i + \frac{(s_i - 1) s_i (2s_i - 1)}{6} \quad (3.4)$$

a

$$s_i \left( \frac{2a_1^{(i)} + s_i - 1}{2} \right)^2 = s_i (a_1^{(i)})^2 + a_1^{(i)} (s_i - 1) s_i + \frac{(s_i - 1)^2 s_i}{4}. \quad (3.5)$$

Z (3.4) a (3.5) dostávame, že

$$s_i \left( \frac{2a_1^{(i)} + s_i - 1}{2} \right)^2 = \sum_{k=1}^{s_i} (a_k^{(i)})^2 - \frac{(s_i - 1) s_i (2s_i - 1)}{6} + \frac{(s_i - 1)^2 s_i}{4},$$

vďaka čomu už môžeme dôkaz priamočiara dokončiť:

$$\begin{aligned} \sum_{i=1}^r \sum_{j=1}^{s_i} (b_j^{(i)})^2 &= \sum_{i=1}^r \sum_{j=1}^{s_i} \left( \frac{1}{s_i} \sum_{k=1}^{s_i} a_k^{(i)} \right)^2 = \sum_{i=1}^r \frac{1}{s_i} \left( \sum_{k=1}^{s_i} a_k^{(i)} \right)^2 \\ &= \sum_{i=1}^r \frac{1}{s_i} \left( s_i \frac{2a_1^{(i)} + s_i - 1}{2} \right)^2 \\ &= \sum_{i=1}^r \left( \sum_{k=1}^{s_i} (a_k^{(i)})^2 - \frac{(s_i - 1) s_i (2s_i - 1)}{6} + \frac{(s_i - 1)^2 s_i}{4} \right) \\ &= \sum_{i=1}^N i^2 - \sum_{i=1}^r \frac{s_i^3 - s_i}{12} = \frac{N(N+1)(2N+1)}{6} - \sum_{i=1}^r \frac{s_i^3 - s_i}{12}. \end{aligned}$$

□

Označme

$$\kappa = \sum_z \frac{k_z^3 - k_z}{12},$$

kde  $k_z$  značí počet, kôľkokrát sa medzi hodnotami  $z_1, \dots, z_N$  vyskytla hodnota  $z$ . Suma  $\sum_z$  značí sčítanie cez všetky rôzne hodnoty  $z_1, \dots, z_N$ . Podľa lemy 10 teda platí, že

$$\sum_{i=1}^N (a^{(i)})^2 = \frac{N(N+1)(2N+1)}{6} - \kappa.$$

Dosadením do (3.3) dostaneme, že

$$\begin{aligned} \sigma_a^2 &= \frac{(N+1)(2N+1)}{6} - \frac{\kappa}{N} - \bar{a}^2 \\ &= \frac{(N+1)(2N+1)}{6} - \frac{\kappa}{N} - \frac{(N+1)^2}{4} \\ &= \frac{(N+1)(N-1)}{12} - \frac{\kappa}{N}. \end{aligned} \quad (3.6)$$

**Veta 11.** Ak platí  $H_0$ , tak pre štatistiku  $W_X^*$  platí, že

$$E [W_X^* | \mathbf{Z}_{(\cdot)}] = \frac{n(N+1)}{2} \quad a \quad var [W_X^* | \mathbf{Z}_{(\cdot)}] = \frac{nm(N+1)}{12} - \frac{nm\kappa}{(N-1)N}.$$

*Dôkaz.* Tvrdenie plynie z dosadenia (3.1), (3.2) a (3.6) do vety 9.  
Pre rozptyl dostávame, že

$$\begin{aligned}\text{var} [W_X^* | \mathbf{Z}_{(\cdot)}] &= \frac{N^2}{N-1} \sigma_a^2 \sigma_c^2 = \frac{nm}{N-1} \sigma_a^2 = \frac{nm}{N-1} \left( \frac{(N+1)(N-1)}{12} - \frac{\kappa}{N} \right) \\ &= \frac{nm(N+1)}{12} - \frac{nm\kappa}{(N-1)N}.\end{aligned}$$

Pre strednú hodnotu dostávame, že

$$\mathbb{E} [W_X^* | \mathbf{Z}_{(\cdot)}] = N\bar{a}\bar{c} = \frac{n(N+1)}{2}.$$

□

Vidíme, že menovateľ výrazu z tvrdenia vety 8 je presne  $\sqrt{\text{var} [W_X^* | \mathbf{Z}_{(\cdot)}]}$ .  
Ešte si všimnime, že

$$\frac{\text{var} [W_X^* | \mathbf{Z}_{(\cdot)}]}{nmN} = \widehat{\xi}_N,$$

a teda takto znormovaný podmienený rozptyl konverguje podľa tvrdenia 6 za  $H_0$  s  $N = n + m \rightarrow \infty$  v pravdepodobnosti ku  $\xi$ .

## 4. Porovnanie testu bez korekcie a testu s korekciou pomocou simulácie

V tejto kapitole sa pokúsime pomocou simulácií nahliadnúť na to, aký je vplyv odvodenej korekcie na dodržovanie predpísanej hladiny  $\alpha$ . Zaujímá nás, ako sa správa hladina Wilcoxonovho testu s narastajúcim počtom zhodných pozorovaní v prípade, že sme v priebehu výpočtu zabudli zohľadniť korekciu. V priebehu celej kapitoly budeme predpokladať, že  $\alpha = 0,05$ .

Uvažujme 2 nezávislé náhodné výbery  $\mathbf{X}$  a  $\mathbf{Y}$  z toho istého rozdelenia. Keďže obidva náhodné výbery sú z toho istého rozdelenia, je splnená  $H_0$  (1.2). Chceme odhadnúť skutočnú hladinu testu s korekciou, respektíve testu bez korekcie.

Pri simulácii budeme postupovať tak, že 10 000-krát vygenerujeme 2 nezávislé náhodné výbery z  $N(100, \sigma^2)$ , vygenerované dáta zaokrúhlime na desiatky, aby nám vznikli zhodné pozorovania a zakaždým určíme, či jednotlivé testy (s korekciou alebo bez korekcie)  $H_0$  zamietli. Podiel počtu zamietnutí a 10 000 nám odhaduje skutočnú hladinu testu. Tento postup simuluje častú príčinu, prečo vznikajú zhody v reálnych dátach - zaokrúhľovanie.

So zmenšujúcou sa smerodajnou odchýlkou ( $\sigma$ ) klesá variabilita výberu, a teda počet zhodných pozorovaní stúpa. Pre  $\sigma = 15$  (relatívne malý počet zhodných pozorovaní) si test s korekciou drží predpísanú 5-percentnú hladinu, zatiaľ čo hladina testu bez korekcie už predpísaných 5 percent nedosahuje (viď Tabuľka 4.1). Pre  $\sigma = 10$  nie je na teste s korekciou badateľná žiadna zmena, stále si drží predpísanú hladinu, avšak hladina testu bez korekcie je už výrazne nižšia než požadovaných 5 percent (viď Tabuľka 4.2). Pre  $\sigma = 5$  (veľmi veľký počet zhodných pozorovaní) test s korekciou stále funguje správne. Hladina testu bez korekcie už naopak nedosahuje ani 2 percentá (viď Tabuľka 4.3).

Simulácie ukazujú, že hladina testu bez korekcie s klesajúcou smerodajnou odchýlkou systematicky klesá. Simulujme toto správanie ešte jedným grafom. Nech rozsah náhodného výberu  $\mathbf{X}$  je 50 a rozsah náhodného výberu  $\mathbf{Y}$  je takisto 50. Budeme simulovať hladinu testu s korekciou a hladinu testu bez korekcie pre smerodajnú odchýlku z intervalu  $[3, 17]$ . Vezmime delenie tohto intervalu také, že vzdialenosť medzi deliacimi bodmi je 0,1. Pre každý deliaci bod nasimulujeme hladinu testu s korekciou a hladinu testu bez korekcie pre smerodajnú odchýlku rovnú tomuto bodu. Výsledky vynesieme do grafu a vidíme, že test s korekciou si stabilne drží predpísanú hladinu, zatiaľ čo hladina testu bez korekcie s klesajúcou smerodajnou odchýlkou klesá k 0 (viď Obr. 4.1).



<b>Rozsah <i>X</i></b>	<b>Rozsah <i>Y</i></b>	<b>Podiel zamietnutí s korekciou [%]</b>	<b>Podiel zamietnutí bez korekcie [%]</b>
10	90	4,90	4,39
90	10	5,01	4,37
30	70	4,75	4,17
70	30	5,06	4,54
50	50	4,86	4,31

Tabuľka 4.1: Porovnanie hladiny jednotlivých testov pre  $\sigma = 15$ .

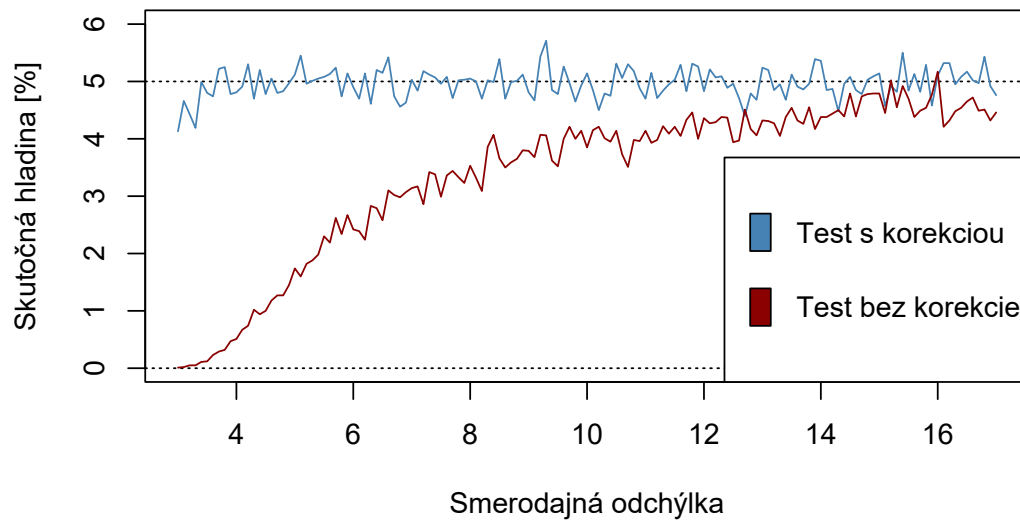
<b>Rozsah <i>X</i></b>	<b>Rozsah <i>Y</i></b>	<b>Podiel zamietnutí s korekciou [%]</b>	<b>Podiel zamietnutí bez korekcie [%]</b>
10	90	4,80	3,68
90	10	4,75	3,71
30	70	5,00	3,91
70	30	4,84	3,82
50	50	4,97	3,98

Tabuľka 4.2: Porovnanie hladiny jednotlivých testov pre  $\sigma = 10$ .

<b>Rozsah <i>X</i></b>	<b>Rozsah <i>Y</i></b>	<b>Podiel zamietnutí s korekciou [%]</b>	<b>Podiel zamietnutí bez korekcie [%]</b>
10	90	4,96	1,73
90	10	5,13	1,46
30	70	5,05	1,81
70	30	4,68	1,31
50	50	4,70	1,42

Tabuľka 4.3: Porovnanie hladiny jednotlivých testov pre  $\sigma = 5$ .

### Simulácia hladiny



Obr. 4.1: Porovnanie hladiny jednotlivých testov v závislosti na  $\sigma$ .

# 5. Dodatky

## 5.1 U-štatistiky

**Definícia 5.** Pre 2 nezávislé náhodné výbery  $(X_1, \dots, X_n)$  a  $(Y_1, \dots, Y_m)$  a jadro  $h : \mathbb{R}^{r+s} \rightarrow \mathbb{R}$  symetrické v bloku prvých  $r$  premenných a v bloku posledných  $s$  premenných je príslušná dvojvýberová U-štatistika definovaná ako

$$U_{n,m} = \frac{1}{\binom{n}{r} \binom{m}{s}} \sum_{\alpha} \sum_{\beta} h(X_{\alpha_1}, \dots, X_{\alpha_r}, Y_{\beta_1}, \dots, Y_{\beta_s}),$$

kde  $\alpha = (\alpha_1, \dots, \alpha_r)$  značí sčítanie cez všetky  $r$ -prvkové podmnožiny  $\{1, \dots, n\}$  a  $\beta = (\beta_1, \dots, \beta_s)$  značí sčítanie cez všetky  $s$ -prvkové podmnožiny  $\{1, \dots, m\}$ .

Označme

$$\sigma_{ij}^2 = \text{cov} \left( h(X_1, \dots, X_i, X_{i+1}, \dots, X_r, Y_1, \dots, Y_j, Y_{j+1}, \dots, Y_s), \right. \\ \left. h(X_1, \dots, X_i, X_{i+1}^*, \dots, X_r^*, Y_1, \dots, Y_j, Y_{j+1}^*, \dots, Y_s^*) \right),$$

kde  $(X_1, \dots, X_r, X_{i+1}^*, \dots, X_r^*)$  a  $(Y_1, \dots, Y_s, Y_{j+1}^*, \dots, Y_s^*)$  sú nezávislé náhodné výbery. Ďalej označme

$$\theta = \mathbb{E} h(X_1, \dots, X_r, Y_1, \dots, Y_s).$$

**Veta 12.** Ak platí, že  $n/N \rightarrow p \in (0,1)$ , keď  $n+m \rightarrow \infty$  a  $\sigma_{11}^2$  je navyše konečné, tak

$$\sqrt{N} (U_{n,m} - \theta) \xrightarrow{D} N(0, \sigma^2), \text{ kde } \sigma^2 = \frac{\sigma_{10}^2}{p} + \frac{\sigma_{01}^2}{1-p}.$$

*Dôkaz.* Dôkaz je uvedený v knihe van der Vaart (1998, strana 166). □

## 5.2 Odvodenie $\sigma_{10}^2$ z vety 5

Analogicky ako pre  $\sigma_{01}^2$  pokračujme s úpravou  $\sigma_{10}^2$  s využitím podmienenej strednej hodnoty. Špeciálne, zamerajme sa na člen  $\mathbb{E} \left[ 1\{X_1 < Y_1\} 1\{X_1 < Y_2\} \right]$ . Postupnosťou úprav dostávame, že

$$\begin{aligned} \mathbb{E} \left[ 1\{X_1 < Y_1\} 1\{X_1 < Y_2\} \right] &= \mathbb{E} \left[ \mathbb{E} \left[ 1\{X_1 < Y_1\} 1\{X_1 < Y_2\} \mid X_1 \right] \right] \\ &= \int_{\mathbb{R}} \mathbb{E} \left[ 1\{X_1 < Y_1\} 1\{X_1 < Y_2\} \mid X_1 = x \right] d\mathbf{P}_{X_1}(x) \\ &= \int_{\mathbb{R}} \mathbb{E} \left[ 1\{x < Y_1\} 1\{x < Y_2\} \right] d\mathbf{P}_{X_1}(x) \\ &= \int_{\mathbb{R}} \mathbb{E} \left[ 1\{x < Y_1\} \right] \mathbb{E} \left[ 1\{x < Y_2\} \right] d\mathbf{P}_{X_1}(x) \\ &= \int_{\mathbb{R}} \left( 1 - F_Y(x) \right) \left( 1 - F_Y(x) \right) d\mathbf{P}_{X_1}(x) \\ &= \mathbb{E} \left[ \left( 1 - F_Y(X_1) \right) \left( 1 - F_Y(X_1) \right) \right]. \end{aligned} \tag{5.1}$$

Analogicky, pre zvyšné členy vystupujúce v  $\sigma_{10}^2$  dostávame, že

$$\begin{aligned}\mathbb{E} \left[ 1\{X_1 < Y_1\} 1\{X_1 = Y_2\} \right] &= \int_{\mathbb{R}} \mathbb{E} \left[ 1\{x < Y_1\} \right] \mathbb{E} \left[ 1\{x = Y_2\} \right] d\mathbf{P}_{X_1}(x) \\ &= \int_{\mathbb{R}} \left( 1 - F_Y(x) \right) \left( F_Y(x) - F_Y(x-) \right) d\mathbf{P}_{X_1}(x) \\ &= \mathbb{E} \left[ \left( 1 - F_Y(X_1) \right) \left( F_Y(X_1) - F_Y(X_1-) \right) \right],\end{aligned}\quad (5.2)$$

$$\begin{aligned}\mathbb{E} \left[ 1\{X_1 = Y_1\} 1\{X_1 = Y_2\} \right] &= \int_{\mathbb{R}} \mathbb{E} \left[ 1\{x = Y_1\} \right] \mathbb{E} \left[ 1\{x = Y_2\} \right] d\mathbf{P}_{X_1}(x) \\ &= \int_{\mathbb{R}} \left( F_Y(x) - F_Y(x-) \right) \left( F_Y(x) - F_Y(x-) \right) d\mathbf{P}_{X_1}(x) \\ &= \mathbb{E} \left[ \left( F_Y(X_1) - F_Y(X_1-) \right) \left( F_Y(X_1) - F_Y(X_1-) \right) \right].\end{aligned}\quad (5.3)$$

S využitím podmienenej strednej hodnoty navyiac rovnako (podmienení náhodnou veličinou  $X_1$ ) upravíme  $\theta$ :

$$\begin{aligned}\theta &= \mathbb{E} \left[ 1\{X_1 < Y_1\} \right] + \frac{1}{2} \mathbb{E} \left[ 1\{X_1 = Y_1\} \right] \\ &= \int_{\mathbb{R}} \mathbb{E} \left[ 1\{x < Y_1\} \right] d\mathbf{P}_{X_1}(x) + \frac{1}{2} \int_{\mathbb{R}} \mathbb{E} \left[ 1\{x = Y_1\} \right] d\mathbf{P}_{X_1}(x) \\ &= \int_{\mathbb{R}} 1 - F_Y(x) d\mathbf{P}_{X_1}(x) + \frac{1}{2} \int_{\mathbb{R}} F_Y(x) - F_Y(x-) d\mathbf{P}_{X_1}(x) \\ &= \mathbb{E} \left[ 1 - F_Y(X_1) \right] + \frac{1}{2} \mathbb{E} \left[ F_Y(X_1) - F_Y(X_1-) \right] \\ &= 1 - \frac{1}{2} \mathbb{E} \left[ F_Y(X_1) + F_Y(X_1-) \right].\end{aligned}\quad (5.4)$$

Konečne, dosadením (5.1), (5.2), (5.3) a (5.4) do (2.2) dostávame, že

$$\begin{aligned}\sigma_{10}^2 &= \mathbb{E} \left[ \left( 1 - F_Y(X_1) \right) \left( 1 - F_Y(X_1) \right) \right] \\ &\quad + \mathbb{E} \left[ \left( 1 - F_Y(X_1) \right) \left( F_Y(X_1) - F_Y(X_1-) \right) \right] \\ &\quad + \frac{1}{4} \mathbb{E} \left[ \left( F_Y(X_1) - F_Y(X_1-) \right) \left( F_Y(X_1) - F_Y(X_1-) \right) \right] \\ &\quad - \left( 1 - \frac{1}{2} \mathbb{E} \left[ F_Y(X_1) + F_Y(X_1-) \right] \right)^2 \\ &= \mathbb{E} \left[ F_Y(X_1) F_Y(X_1-) \right] + \frac{1}{4} \mathbb{E} \left[ F_Y(X_1)^2 - 2F_Y(X_1) F_Y(X_1-) + F_Y(X_1-)^2 \right] \\ &\quad - \frac{1}{4} \left( \mathbb{E} \left[ F_Y(X_1) + F_Y(X_1-) \right] \right)^2 \\ &= \frac{1}{4} \mathbb{E} \left[ \left( F_Y(X_1) + F_Y(X_1-) \right)^2 \right] - \frac{1}{4} \left( \mathbb{E} \left[ F_Y(X_1) + F_Y(X_1-) \right] \right)^2 \\ &= \frac{1}{4} \text{var} \left( F_Y(X_1) + F_Y(X_1-) \right).\end{aligned}$$

# Záver

Táto práca sa venovala dvojvýberovému Wilcoxonovmu testu. Najprv sme sa zaoberali jednoduchším prípadom, kedy sme predpokladali, že dáta pochádzajú zo spojitých rozdelení. Následne sme sa začali venovať zložitejšiemu prípadu, kedy sme sa vzdali predpokladu spojitosti a začali predpokladať, že v dátach sa môžu vyskytovať zhody.

S využitím známych výsledkov o U-štatistikách sme odvodili asymptotické rozdelenie Wilcoxonovej testovej štatistiky v prípade existencie zhôd, čo nás viedlo k návrhu korekcie testu pre dáta obsahujúce zhody. Ďalej sme poukázali na súvislosť odvodennej korekcie a podmieneného rozptylu testovej štatistiky (a podmienenej strednej hodnoty).

Hlavným prínosom teoretickej časti práce teda bola aplikácia známych výsledkov o U-štatistikách na odvodenie asymptotického rozdelenia Wilcoxonovej testovej štatistiky v prípade zhôd a výpočet podmieneného rozptylu tejto štatistiky v tretej kapitole.

Nakoniec sme pomocou simulácií sledovali aký vplyv má odvodená korekcia na skutočnú hladinu testu pri zvyšujúcom sa počte zhodných pozorovaní. Ukázali sme, že zatiaľ čo test s korekciou si pri rastúcom počte zhodných pozorovaní stabilne držal predpísanú hladinu, hladina testu bez korekcie systematicky klesala až k nule.

Test s korekciou stanovenú hladinu dodržiava len asymptoticky. V práci sme uviedli príklad ako by sme postupovali, keby sme chceli test založený na presnom rozdelení. Predmetom ďalšej štúdie by teda mohla byť otázka, ako súvisia rozsahy výberov (kedy sú už dostatočne veľké) a skutočná hladina testu s korekciou založeného na asymptotickom rozdelení.

# Zoznam použitej literatúry

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 978-80-7378-003-6.

LEHMANN, E. L. a D'ABRERA, H. J. M. (1975). *Nonparametrics: Statistical Methods Based on Ranks*. Holden-day series in probability and statistics. Holden-Day, San Francisco. ISBN 0-8162-4994-6.

VAN DER VAART, A. W. (1998). *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge. ISBN 0-521-78450-6.