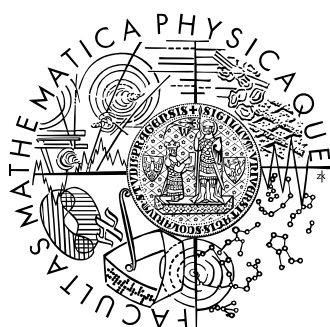


Univerzita Karlova v Praze  
Matematicko-fyzikální fakulta

# BAKALÁŘSKÁ PRÁCE



Adrián Andráš

## **William S. Gosset a jeho význam pro rozvoj teorie a praxe statistiky**

Katedra pravděpodobnosti a matematické statistiky  
Vedoucí bakalářské práce: Mgr. Michal Kulich, Ph.D.

Studijní program: Matematika  
Studijní obor: Obecná matematika

2009

Ďakujem Mgr. Michalovi Kulichovi, Ph.D. za odborné vedenie mojej bakalárskej práce, za cennú pomoc a čas, ktorý mi pri písaní venoval. Ďakujem tiež Mgr. Šárke Došlej, MSc. za lepšie porozumenie danej problematike a za pomoc pri numerických výpočtoch použitých v tejto práci. Veľká vďaka patrí aj PhDr. Milene Režnej a Veronike Stankovianskej, ktoré mi pomáhali s prekladom anglických odborných textov a tiež Nadke Langovej za užitočné rady pri estetickom spracovaní textu.

Prehlasujem, že som svoju bakalársku prácu napísal samostatne a výhradne s použitím citovaných prameňov. Súhlasím so zapožičaním práce a jej zverejňovaním.

V Prahe dňa 7. 8. 2009

Adrián Andráš

# OBSAH

<b>1</b>	<b>William Sealy Gosset .....</b>	<b>5</b>
1.1	Životopis .....	5
1.2	Prečo práve Student?.....	6
1.3	Studentova štatistická filozofia .....	6
<b>2</b>	<b>Pravdepodobná chyba priemeru .....</b>	<b>10</b>
2.1	Úvod.....	10
2.2	Rozdelenie výberovej smerodajnej odchýlky .....	10
2.3	Korelácia medzi $X - EX$ a $s$ výberu z normálneho rozdelenia .....	18
2.4	Hľadanie hustoty nezávislej na $\sigma$ .....	20
2.5	Vlastnosti hustoty z.....	24
2.6	Použitie tabuľky v praxi.....	26
<b>3</b>	<b>Objavené nepresnosti v Studentovom článku .....</b>	<b>28</b>
<b>4</b>	<b>Pravdepodobná chyba korelačného koeficientu.....</b>	<b>30</b>
<b>5</b>	<b>Štatistici, ktorých ovplyvnila Studentova tvorba .....</b>	<b>34</b>
<b>6</b>	<b>Návrat k súčasnosti .....</b>	<b>35</b>
6.1	Studentovo t-rozdelenie .....	35
6.2	Výberový korelačný koeficient .....	36
<b>7</b>	<b>Záver.....</b>	<b>37</b>
	<b>Literatúra.....</b>	<b>38</b>

**Názov práce:** William S. Gosset a jeho význam pro rozvoj teorie a praxe statistiky

**Autor:** Adrián Andráš

**Katedra:** Katedra pravděpodobnosti a matematické statistiky

**Vedúci bakalárskej práce:** Mgr. Michal Kulich, Ph.D.

**e-mail vedúceho:** kulich@karlin.mff.cuni.cz

**Abstrakt:** Táto práca je písaná na počesť stého výročia Gossetovho významného článku „Pravdepodobná chyba priemeru“. Úvod je venovaný životopisu Williama Sealyho Gosseta, no prevažná časť sa už venuje jeho článku. Približuje nám ho, hovorí o tom, ako Gossetovo (Studentovo)  $t$ -rozdelenie vzniklo, akých chýb sa Gosset dopustil pri výpočtoch a akej filozofie sa držal. Taktiež (stručne) rozoberá jeho druhý významný článok „Pravdepodobná chyba korelačného koeficientu“, ktorý sa venuje rozdeleniu výberového korelačného koeficientu za predpokladu náhodného výberu z dvojrozmerného normálneho rozdelenia. V závere sú spomenutí významní štatistici a matematici, ktorí venovali Gossetovmu článku značnú pozornosť, prípadne čerpali z jeho článkov niektoré závery. Posledná kapitola porovnáva  $t$ -rozdelenie v súčasnej štatistike s tou pred 100 rokmi.

**Kľúčové slová:** Gosset, Studentovo  $t$  – rozdelenie, pravdepodobná chyba, výberový korelačný koeficient

**Title:** William S. Gosset and his contribution to the advancement of statistical theory and practice

**Author:** Adrián Andráš

**Department:** Department of Probability and Mathematical Statistics

**Supervisor:** Mgr. Michal Kulich, Ph.D.

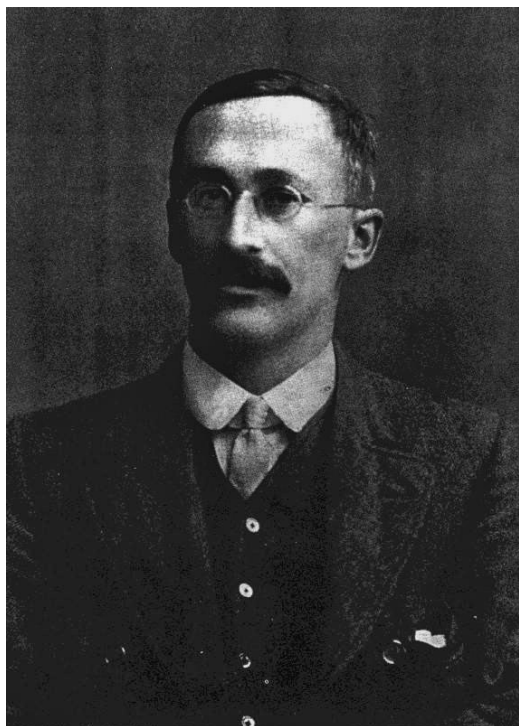
**Supervisor's e-mail address:** kulich@karlin.mff.cuni.cz

**Abstract:** This thesis is written in honour of the centenary of Gosset's seminal paper "The Probable Error of a Mean." While the introduction is devoted to William Sealy Gosset's biography, the rest focuses on Gosset's (Student's) article. The work makes the original text accessible as it explains the origins of Student's  $t$  – distribution and Gosset's errors in calculation in addition to the concepts he adopted. Rather briefly, the paper equally deals with another Gosset's influential article "The Probable Error of a Correlation Coefficient," addressing the distribution of the sample correlation coefficient in the case of a random sample from a bivariate normal distribution. A number of distinguished mathematicians and statisticians, either paying considerate attention to Gosset's work or basing their results on his findings, are mentioned. The conclusive chapter is concerned with the comparison of the  $t$  – distribution as perceived in contemporary statistics as opposed to the statistics of a century ago.

**Keywords:** Gosset, Student's  $t$  – distribution, the probable error, a sample correlation coefficient

# 1 William Sealy Gosset

## 1.1 Životopis



*William Sealy Gosset (Student)*

(1876 – 1937)

William Sealy Gosset sa narodil 13. júna 1876 v Caterbury, Veľká Británia (VB) a zomrel 16. októbra 1937, Beaconsfield, VB. Študoval na Winchester College a New College na Oxforde, VB. Počas štúdia bol ocenený niekoľkými titulmi za jeho vedomosti z matematiky („First in the Mathematical Moderations examination“) a z chémie („First-Class Degree in Chemistry“).

Po skončení školy ho zamestnal Arthur Guinness do svojej spoločnosti v Dubline, Írsku v roku 1899, ktorá sa zaoberala výrobou piva. Gosset si postupom času v práci začal

uvedomovať dôležitosť štatistických metód pri postupoch, ktoré sa vo firme prevádzali a obzvlášť tiež problém malých výberov. To Gosseta donútilo, aby sa začal hlbšie venovať práve týmto problémom (zhrnul to v nepublikovanom zázname „Aplikácia ‘Zákonu chyby’ na prácu pivovaru“ datovanom 3. novembra 1904). Keď nenašiel potrebné informácie v literatúre, dohodol si stretnutie s Karlom Pearsonom (významným štatistikom) v júli 1905. A neskôr v rokoch 1906-1907 s povolením Guinnessa, Gosset strávil dva semestre akademického roku v Pearsonovom biometrickom laboratóriu na univerzite v Londýne (University College London). Problémy, ktorými sa zaoberal, zhrnul vo vedeckom časopise *Biometrika*, ktorého spoluzakladateľom bol Pearson. Počas 25 ročného obdobia celkovo Gosset publikoval v *Biometrike* 14 článkov, ktoré napísal pod pseudonymom *Student*.

Podrobnejší Studentov životopis je možné nájsť v knihe [6].

## 1.2 Prečo práve Student?

Touto otázkou sa zaoberal článok [5] str. 189, v ktorom sa vysvetľuje Gossetovo tvrdenie, že jeho matematické a filozofické závery nebolo možné použiť na praktické využitie u konkurencie, no nakoniec ich mohol publikovať, ale len pod pseudonymom, aby sa zabránilo ťažkostiam so zvyškom personálu.

No taktiež sa stretávame s tvrdením, že anonymita bola vyžadovaná priamo Guinnessom, aby sa konkurencia nedozvedela, že je veľmi užitočné zamestnávať štatistikov. Toto tvrdenie je zase úplne protichodné Hotellingovmu (matematickému štatistikovi a významnému ekonomickému teoretikovi), podľa ktorého anonymita bola potrebná kvôli pivovarníkom „vo vnútri“ a nie tých „zvonka“.

Pseudonym Student bolo vymyslené Christopherom Diggesom La Toucheom, ktorý bol výkonným riaditeľom pivovaru.

## 1.3 Studentova štatistická filozofia

Student (odteraz mu budeme takto stále hovoriť, lebo pseudonym sa zachoval v štatistike až do súčasnosti) bol ovplyvnený spočiatku *Bayesovou filozofiou* aj preto, že čítal články a knihy od Karla Pearsona.

Pre priblíženie Bayesovej filozofie a základných princípov *bayesovských metód* uvedieme jednoduchý príklad, ktorý je možné nájsť v Andělovej knihe [2].

Máme veľkú sériu výrobkov. Nech  $N$  je počet kusov v tejto sérii a  $M$  značí počet vadných výrobkov. Číslo  $M$  nie je známe. Pravdepodobnosť, že vybraný kus bude chybný, je rovná  $p = M/N$  (tá tiež nie je známa).

Uskutočnime tzv. *výber s vrátením*. To znamená, že preskúmaný výrobok sa znova vráti do série a v ďalších ťahoch môže byť znova vytiahnutý. Tým sa docieli, že jednotlivé ťahy budú na sebe nezávislé.

Pravdepodobnosť, že z  $n$  vytiahnutých výrobkov bude práve  $m$  vadných je

$$r(m | p) = \binom{n}{m} p^m (1-p)^{n-m}, \quad 0 \leq m \leq n. \quad (1.1)$$

Ak vo výbere bolo nájdených práve  $m$  vadných výrobkov, je najlepším nestranným odhadom parametru  $p$  veličina  $p^* = m/n$ .

Odôvodnenie: Nech séria výrobkov pozostáva z náhodných rovnako rozdelených veličín  $X_1, \dots, X_n$ , kde  $X_i \sim \text{Alt}(p)$ .

Nech  $P(X_i = 1) = p$  je pravdepodobnosť, že vybraný výrobok bude chybný a  $P(X_i = 0) = 1 - p$  je pravdepodobnosť, že vybraný výrobok nebude chybný.

Ak  $m = \sum_{i=1}^n X_i$  označuje počet vadných výrobkov vo výbere a vieme, že  $X_i$  sú rovnako rozdelené, potom

$$E p^* = \frac{E \sum_{i=1}^n X_i}{n} = \frac{\sum_{i=1}^n E X_i}{n} = \frac{n \cdot E X_1}{n} = \frac{n p}{n} = p,$$

z čoho plynie nestrannosť  $p^*$ .

Rozptyl  $p^*$  je rovný  $\text{var } p^* = p(1-p)/n$ .

Ak sa stane, že výroba takých rovnako veľkých sérií prebieha už dlhšiu dobu a vedie sa dostatočne podrobná evidencia, je možné zistiť po čase presne, koľko bolo vyrobených vadných výrobkov v každej sérii, a teda sú známe aj podiely  $p_1, \dots, p_K$  vadných výrobkov v každej zo skôr vyrobených  $K$  sérií. Keďže podiely nie sú obvykle totožné, ľahko prijmeme názor, že podiel  $p$  vadných výrobkov v sérii je náhodná veličina, o ktorej rozdelení nás informuje náhodný výber  $p_1, \dots, p_K$ .

Keďže náhodná veličina  $p$  môže nadobúdať len hodnoty  $0, 1/N, \dots, N/N$ , je diskrétna. Ak je  $N$  veľké, je možné distribučnú funkciu tejto veličiny dostatočne presne aproximovať nejakou absolútne spojitou distribučnou funkciou.

V našom prípade budeme predpokladať nasledovnú aproximáciu náhodnej veličiny  $p$  B-rozdelením s hustotou

$$q(p) = \frac{1}{B(a,b)} p^{a-1} (1-p)^{b-1}, \quad 0 < p < 1, \quad (1.2)$$

kde  $a > 0$ ,  $b > 0$  sú známe parametre napr. dostatočne presne určené z hodnôt  $p_1, \dots, p_K$ , ktoré považujeme za výber z rozdelenia s hustotou (1.2). Hustota  $q$  je marginálna hustota veličiny  $p$  a nazýva sa *apriórna hustota*. Predpokladá sa totiž, že

je známa skôr, než sa začneme zaoberať danou skupinou výrobkov. Pravdepodobnosti (1.1) potom tvoria podmienenú hustotu náhodnej veličiny  $m$  pri danej hodnote  $p$ .

Aposteriórnou hustotou nazývame výraz  $\pi(p | m)$ , ktorý je podľa Bayesovej vety rovný

$$\begin{aligned}\pi(p | m) &= \frac{\frac{1}{B(a,b)} p^{a-1} (1-p)^{b-1} \binom{n}{m} p^m (1-p)^{n-m}}{\int_0^1 \frac{1}{B(a,b)} p^{a-1} (1-p)^{b-1} \binom{n}{m} p^m (1-p)^{n-m} dp} = \\ &= \frac{1}{B(a+m, b+n-m)} p^{a+m-1} (1-p)^{b+n-m-1}.\end{aligned}\quad (1.3)$$

Na základe aposteriórnej hustoty  $\pi(p | m)$  je možné zostrojiť intervalový odhad pre veličinu  $p$  nasledovne.

Nech čísla  $D$  a  $H$  ( $0 \leq D < H \leq 1$ ) splňujú podmienky

$$\int_D^H \pi(p | m) dp = 1 - \alpha, \quad \int_0^D \pi(p | m) dp = \alpha/2, \quad \int_H^1 \pi(p | m) dp = \alpha/2,$$

kde  $\alpha \in (0,1)$ . Potom platí

$$P(D < p < H | m) = 1 - \alpha. \quad (1.4)$$

Interval  $(D, H)$  je *bayesovský interval spoľahlivosti* pre  $p$ .

Ak sú čísla  $D$  a  $H$  zvolené tak, aby splňali podmienku (1.4), môžeme pristúpiť k testovaniu hypotéz pre parameter  $p$ .

Predpokladáme test hypotézy  $H_0 : p \in A$ , kde  $A$  je ľubovoľný interval na  $(0, 1)$ . Ak celý interval  $A$  padne mimo  $(D, H)$ , hypotézu  $H_0$  zamietame. V tomto prípade je  $P(A | m) \leq \alpha$ .

Nedá sa však založiť test len na hodnote  $P(A | m)$  bez porovnania s intervalom  $(D, H)$ . Pretože, ak by interval  $A$  bol príliš krátky, tak  $P(A | m)$  by bola veľmi malá a to by spôsobilo zamietnutie  $H_0$ , napriek tomu, že poloha  $A$  by vzhľadom k hustote (1.3) nebola nijako „podozrivá“.

Ak sa vrátíme k nášmu príkladu, parameter  $p$  (binomického rozdelenia) nebol považovaný za konštantu, ale za náhodnú veličinu. O tejto možnosti prvýkrát hovoril Thomas Bayes, a preto sa jeho metódy nazývajú *bayesovské*.



Najväčším kladom bayesovských postupov je ten, že sa neobmedzujú len na samotný výber, ale berú do úvahy aj doplňujúce informácie. Táto výhoda sa však stráca, ak doplňujúce informácie nemáme k dispozícii.

Konkrétne využitie Bayesovej filozofie Student sčasti uplatnil napr. vo svojom článku [9], ku ktorému sa neskôr dostaneme.

V ďalšej kapitole sa však pozrieme najskôr na Studentov významnejší článok [10].

## 2 Pravdepodobná chyba priemeru

### 2.1 Úvod

V celom článku sa predpokladá náhodný výber  $n$  ( $n \geq 2$ ) jedincov  $Y_1, \dots, Y_n$  z populácie s normálnym rozdelením, čiže  $Y_i \sim N(\mu, \sigma^2)$ , kde  $\mu$  značí strednú hodnotu náhodnej veličiny  $Y_i$  a  $\sigma$  označuje jej smerodajnú odchýlku.

Ďalej vieme, že výberový priemer  $\bar{Y}$  má normálne rozdelenie (kde  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ ) so strednou hodnotou  $\mu$  a smerodajnou odchýlkou  $\frac{\sigma}{\sqrt{n}}$ , čiže  $\bar{Y} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ .

Označíme  $s$  ako výberovú smerodajnú odchýlku náhodného výberu  $X_1, \dots, X_n$ , kde

$$X_i = Y_i - \mu \quad \text{a} \quad s = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2} \quad \text{pre } i = 1, \dots, n.$$

Cieľom Studentovho článku, bolo určiť rozdelenie náhodnej veličiny

$$z = \frac{\bar{Y} - \mu}{s}$$

a pre túto veličinu zostaviť štatistické tabuľky distribučnej funkcie  $F(z)$ .

### 2.2 Rozdelenie výberovej smerodajnej odchýlky

Ak uvažujeme výberovú smerodajnú odchýlku  $s$  náhodného výberu  $X_1, \dots, X_n$  a náhodnú veličinu  $X_i$  ako v úvode (kap. 2.1), potom

$$s^2 = \frac{\sum_{i=1}^n X_i^2}{n} - \left( \frac{\sum_{i=1}^n X_i}{n} \right)^2 = \frac{\sum_{i=1}^n X_i^2}{n} - \frac{\sum_{i=1}^n X_i^2}{n^2} - \frac{2 \sum_{\substack{i=1 \\ i \neq j}}^n \sum_{j=1}^n X_i X_j}{n^2}.$$

Strednú hodnotu veličiny  $s^2$  spočítame ako

$$\begin{aligned}
 Es^2 &= \frac{E \sum_{i=1}^n X_i^2}{n} - E \left( \frac{\sum_{i=1}^n X_i}{n} \right)^2 = \frac{E \sum_{i=1}^n X_i^2}{n} - \frac{E \sum_{i=1}^n X_i^2}{n^2} - \frac{2 \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n EX_i X_j}{n^2} = \\
 &= \frac{\sum_{i=1}^n EX_i^2}{n} - \frac{\sum_{i=1}^n EX_i^2}{n^2} - \frac{2 \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n EX_i X_j}{n^2} = \\
 &= \frac{nEX_i^2}{n} - \frac{nEX_i^2}{n^2} - \frac{2 \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n EX_i X_j}{n^2} = EX_i^2 - \frac{EX_i^2}{n} - \frac{2 \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n EX_i X_j}{n^2}. \quad (2.1)
 \end{aligned}$$

Keďže

$$EX_i = EY_i - E\mu = \mu - \mu = 0,$$

$$\text{var } X_i = EX_i^2 - (EX_i)^2 = EX_i^2 = E \left( Y_i - \underbrace{\mu}_{=EY_i} \right)^2 = \sigma^2,$$

dostávame, že

$$X_i \sim N(0, \sigma^2). \quad (2.2)$$

Vieme, že  $X_1, \dots, X_n$  sú nekorelované náhodné veličiny (lebo sa jedná o náhodný výber z normálneho rozdelenia a odčítanie  $\mu$  tomu nijak neublíži). Potom môžeme písať, že

$$\text{corr}(X_i, X_j) = 0, \text{ čiže } \text{cov}(X_i, X_j) = 0 = EX_i X_j - \underbrace{EX_i EX_j}_{=0}, \text{ a teda } EX_i X_j = 0.$$

Teda posledný člen  $\frac{2 \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n EX_i X_j}{n^2}$  v (2.1) vypadne a po zavedení  $\mu_k = EX^k$ , dostaneme

$$Es^2 = \mu_2 \frac{(n-1)}{n},$$

čo je v podstate 1. moment (označíme ho  $M_1'$ ) veličiny  $s^2$ .

Student sa pokúsil spočítať aj ďalšie momenty veličiny  $s^2$ .

Výpočet 2. momentu ( $M_2'$ )  $s^2$ :

$$\begin{aligned}
 s^4 &= \left\{ \frac{\sum_{i=1}^n X_i^2}{n} - \left( \frac{\sum_{i=1}^n X_i}{n} \right)^2 \right\}^2 = \left( \frac{\sum_{i=1}^n X_i^2}{n} \right)^2 - \frac{2 \sum_{i=1}^n X_i^2}{n} \left( \frac{\sum_{i=1}^n X_i}{n} \right) + \left( \frac{\sum_{i=1}^n X_i}{n} \right)^4 = \quad (2.3) \\
 &= \frac{\sum_{i=1}^n X_i^4}{n^2} + \frac{2 \sum_{\substack{i=1 \\ i \neq j}}^n \sum_{j=1}^n X_i^2 X_j^2}{n^2} - \frac{2 \sum_{i=1}^n X_i^4}{n^3} - \frac{4 \sum_{\substack{i=1 \\ i \neq j}}^n \sum_{j=1}^n X_i^2 X_j^2}{n^3} + \frac{\sum_{i=1}^n X_i^4}{n^4} + \frac{6 \sum_{\substack{i=1 \\ i \neq j}}^n \sum_{j=1}^n X_i^2 X_j^2}{n^4} \\
 &\quad + \text{ďalšie výrazy, pri ktorých je } X_i \text{ umocnené na nepárnu mocninu.}
 \end{aligned}$$

Vieme, že platí  $X_i \sim N(0, \sigma^2)$  a pre  $k \in \mathbf{N}$ ,

$$\begin{aligned}
 E(X_i - EX_i)^k &\stackrel{EX_i=0}{=} EX_i^k = \int_{-\infty}^{+\infty} x^k f(x) dx = \int_{-\infty}^{+\infty} x^k \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2}} dx = \quad (2.4) \\
 &= \begin{cases} 0 & \text{pre } k \text{ nepárne,} \\ (k-1)(k-3) \dots \cdot 1 \cdot \sigma^k & \text{pre } k \text{ párne.} \end{cases}
 \end{aligned}$$

Keďže sú  $X_1, \dots, X_n$  nekorelované náhodné veličiny, dostávame pre  $l \in \mathbf{N}$  a  $i \neq j$

$$EX_i^k X_j^l = \begin{cases} 0 \text{ (plynie z (2.4))} & \text{pre } k \text{ alebo } l \text{ nepárne,} \\ EX_i^k EX_j^l & \text{pre } k \text{ aj } l \text{ párne,} \end{cases} \quad (2.5)$$

Teda všetky výrazy, v ktorých sa vyskytuje  $X_i$  umocnené na nepárnu mocninu pri

výpočte  $Es^4$  v (2.3) zmiznú a keďže  $\sum_{i=1}^n X_i^4$  má  $n$  členov a  $\sum_{\substack{i=1 \\ i \neq j}}^n \sum_{j=1}^n X_i^2 X_j^2$  má

$\frac{n(n-1)}{2}$  členov, ostane nám

$$\begin{aligned}
 M_2' = Es^4 &= \frac{\mu_4}{n} + \mu_2^2 \frac{(n-1)}{n} - \frac{2\mu_4}{n^2} - 2\mu_2^2 \frac{(n-1)}{n^2} + \frac{\mu_4}{n^3} + 3\mu_2^2 \frac{(n-1)}{n^4} = \quad (2.6) \\
 &= \frac{\mu_4}{n^3} \{n^2 - 2n + 1\} + \frac{\mu_2^2}{n^3} (n-1) \{n^2 - 2n + 3\}
 \end{aligned}$$

Z (2.4) plynie, že

$$\begin{aligned}
 \mu_2 &= EX^2 = \sigma^2, \\
 \mu_4 &= EX^4 = 3\sigma^4,
 \end{aligned}$$

čiže

$$\mu_4 = 3\mu_2^2.$$

Teda po dosadení do (2.6) dostaneme

$$M_2' = \mu_2^2 \frac{(n-1)}{n^3} \{3n - 3 + n^2 - 2n + 3\} = \mu_2^2 \frac{(n-1)(n+1)}{n^2}.$$

Analogicky sa spočíta 3. a 4. moment veličiny  $s^2$  a dostaneme

$$M_3' = \mu_2^3 \frac{(n-1)(n+1)(n+3)}{n^3},$$

$$M_4' = \mu_2^4 \frac{(n-1)(n+1)(n+3)(n+5)}{n^4}.$$

Čiže sa črtá pravidlo formovania jednotlivých momentov  $s^2$ .

Označme teraz  $M_R$  ako  $R$  - tý *centrálny moment* veličiny  $s^2$ . Potom

$$\begin{aligned} M_2 &= E(s^2 - Es^2)^2 = \underbrace{E(s^2)^2}_{M_2'} - (E(s^2))^2 = \mu_2^2 \frac{(n-1)(n+1)}{n^2} - \left( \mu_2 \frac{(n-1)}{n} \right)^2 \\ &= 2 \cdot \frac{\mu_2^2 (n-1)}{n^2} \end{aligned}$$

a analogicky sa spočítajú ďalšie momenty (podrobnejší výpočet je možné nájsť v Studentovom článku [10]).

$$M_3 = E(s^2 - Es^2)^3 = 8\mu_2^3 \frac{(n-1)}{n^3},$$

$$M_4 = E(s^2 - Es^2)^4 = 12\mu_2^4 \frac{(n-1)(n+3)}{n^4}.$$

Ďalej Student zistil, že ak spočíta šikmosť náhodnej veličiny  $s^2$

$$\beta_1 = \left[ \frac{E(s^2 - Es^2)^3}{\left( \sqrt{E(s^2 - Es^2)^2} \right)^3} \right]^2 = \frac{M_3^2}{M_2^3} = \frac{8}{n-1},$$

a následne špicatosť

$$\beta_2 = \frac{E(s^2 - Es^2)^4}{\left( \sqrt{E(s^2 - Es^2)^2} \right)^4} = \frac{M_4}{M_2^2} = \frac{3(n+3)}{n-1},$$

vyjde mu

$$2\beta_2 - 3\beta_1 - 6 = 2 \cdot \frac{3(n+3)}{n-1} - 3 \cdot \frac{8}{n-1} - 6 = 0$$

a to odpovedá *Pearsonovému rozdeleniu typu III*, ktoré v súčasnosti poznáme pod názvom *gama rozdelenie*.

Teda hustota náhodnej veličiny  $s^2$  bude vyzerat' ako

$$f(x) = Cx^p e^{-\gamma x},$$

kde

$$\gamma = 2 \frac{M_2}{M_3} = \frac{4\mu_2^2(n-1)n^3}{8n^2\mu_2^3(n-1)} = \frac{n}{2\mu_2},$$

$$p = \frac{4}{\beta_1} - 1 = \frac{n-1}{2} - 1 = \frac{n-3}{2},$$

z čoho plynie hustota

$$f(x) = Cx^{\frac{n-3}{2}} e^{-\frac{nx}{2\mu_2}}, \quad \text{kde } x > 0.$$

**Poznámka:** Presné zdôvodnenie *Pearsonového rozdelenia typu III*. je možné nájsť v Pearsonovom článku [7].

Ak si označíme

$$I = C \int_0^{+\infty} x^{\frac{n-3}{2}} e^{-\frac{nx}{2\mu_2}} dx,$$

potom

$$1 = \frac{C \int_0^{+\infty} x^{\frac{n-3}{2}} e^{-\frac{nx}{2\mu_2}} dx}{I},$$

z čoho plynie, že  $\frac{C x^{\frac{n-3}{2}} e^{-\frac{nx}{2\mu_2}}}{I}$  je hustota  $s^2$ .

Potom 1. moment náhodnej veličiny  $s^2$  je (keďže *gama funkcia* nebola dobre známa, Student použil metódu *per partes*)

$$Es^2 = \frac{C \int_0^{+\infty} \overbrace{x \cdot x^{\frac{n-3}{2}}}^{\frac{n-1}{2}} e^{-\frac{nx}{2\mu_2}} dx}{I} \stackrel{\text{per partes}}{=} \left| \begin{array}{l} u = x^{\frac{n-1}{2}} \\ v' = e^{-\frac{nx}{2\mu_2}} \end{array} \right. \left. \begin{array}{l} u' = \left( \frac{n-1}{2} \right) x^{\frac{n-3}{2}} \\ v = \frac{e^{-\frac{nx}{2\mu_2}}}{-\frac{n}{2\mu_2}} \end{array} \right| =$$

$$\begin{aligned}
&= \frac{C}{I} \left\{ \left[ -\frac{2\mu_2}{n} x^{\frac{n-1}{2}} e^{-\frac{nx}{2\mu_2}} \right]_{x=0}^{+\infty} + \frac{2\mu_2}{n} \left( \frac{n-1}{2} \right) \underbrace{\int_0^{+\infty} x^{\frac{n-3}{2}} e^{-\frac{nx}{2\mu_2}} dx}_{\frac{I}{C}} \right\} = \\
&= \left( -\frac{2\mu_2}{n} \cdot \lim_{x \rightarrow \infty} \frac{x^{\frac{n-1}{2}}}{e^{\frac{nx}{2\mu_2}}} - 0 \right) + \frac{C}{I} \mu_2 \frac{(n-1)}{n} \frac{I}{C} \stackrel{(\#)}{=} 0 + \frac{C}{I} \mu_2 \frac{(n-1)}{n} \frac{I}{C} = \mu_2 \frac{(n-1)}{n}.
\end{aligned}$$

Analogicky sa spočítajú druhy, tretí, štvrtý moment, ... a dostaneme

$$\begin{aligned}
M_2' &= \mu_2^2 \frac{(n-1)(n+1)}{n^2}, \\
M_3' &= \mu_2^3 \frac{(n-1)(n+1)(n+3)}{n^3}, \\
M_4' &= \mu_2^4 \frac{(n-1)(n+1)(n+3)(n+5)}{n^4}.
\end{aligned}$$

Teda rozdelenie  $s^2$  má všetky momenty rovnaké ako Pearsonovo rozdelenie typu III a Student prehlásil, že rozdelenie  $s^2$  je teda práve Pearsonovým rozdelením typu III.

dané hustotou  $f(x) = Cx^{\frac{n-3}{2}} e^{-\frac{nx}{2\mu_2}}$ . (V súčasnosti je známe, že neplatí tvrdenie: Rovnosť momentov dvoch rozdelení implikuje aj rovnosť rozdelení. Vid' príklad v Andělovej knihe [3] str. 16. Teda Studentov postup tu nie je úplne ošetrový.)

Keďže už poznáme rozdelenie náhodnej veličiny  $s^2$ , nie je problém konečne nájsť hľadané rozdelenie náhodnej veličiny  $s$ . K tomu použijeme *vetu o transformácií náhodnej veličiny*.

**Veta (O transformácií náhodnej veličiny)** *Nech  $X$  má spojitú distribučnú funkciu  $F$ . Nech  $F'(x) = f(x)$  existuje všade až na konečne mnoho bodov. Nech je  $t$  rýdzo monotónna funkcia, ktorá má všade nenulovú deriváciu. Položme  $Y = t(X)$ . Označme  $\tau$  inverznú funkciu k  $t$ . Potom  $Y$  má hustotu*

$$g(y) = f[\tau(y)] |\tau'(y)|.$$

*Dôkaz.* Vid' Andělovu knihu [3] str. 48.

V našom prípade môžeme označiť

---

(#) Na výpočet limity použijeme  $(n-1)/2$  krát L'Hospitalovo pravidlo pre prípad " $\frac{\infty}{\infty}$ "

$$X = s^2, \quad \text{kde } \phi(s^2) = Cx^{\frac{n-3}{2}} e^{-\frac{nx}{2\mu_2}} \text{ je hustota } X,$$

$$Y = s = \sqrt{X}, \quad \text{kde } \psi(s) \text{ je hustota } Y \text{ a } x = y^2.$$

Najskôr však potrebujeme overiť predpoklady vety, aby sme zistili, či vetu je možné aplikovať na náš prípad. Čiže máme zobrazenie

$$t: X \rightarrow \sqrt{X}.$$

Deriváciu  $t$  spočítame ako

$$t'(x) = \frac{1}{2\sqrt{x}} \neq 0, \quad \forall x \in \mathfrak{R}^+ \setminus \{0\}.$$

A keďže

$$t'(x) = \frac{1}{2\sqrt{x}} > 0, \quad \forall x \in \mathfrak{R}^+ \setminus \{0\},$$

dostávame rýdzo-monotónnosť  $t$ .

Keďže všetky predpoklady sú splnené, môžeme použiť vetu nasledujúcim spôsobom:

$$\psi(y) = \phi[\tau(y)]|\tau'(y)|,$$

kde

$$\tau(y) = x = y^2,$$

čiže po dosadení do  $\tau(y)$  a  $\tau'(y)$  dostaneme

$$\psi(y) = \phi(y^2) \cdot 2y.$$

Keď nahradíme  $y$  premennou  $s$ , môžeme písať

$$\psi(s) = 2 \cdot C \cdot s \cdot (s^2)^{\frac{n-3}{2}} e^{-\frac{ns^2}{2\mu_2}} = 2 \cdot C \cdot s^{n-2} e^{-\frac{ns^2}{2\mu_2}}, \quad \text{kde } s > 0,$$

čiže sme nakoniec dostali hustotu  $s$ .

Teda ak  $A = 2C$ , potom

$$f(s) = A \cdot s^{n-2} e^{-\frac{ns^2}{2\mu_2}} \tag{2.7}$$

je hustota výberovej smerodajnej odchýlky  $s$  výberu z populácie, ktorého náhodné veličiny majú normálne rozdelenie so smerodajnou odchýlkou  $\sigma$  ( $X_i \sim N(0, \sigma^2)$ ).

Chceme teraz spočítať konštantu  $A$ . Tú vypočítame z hustoty smerodajnej odchýlky  $s$  a hustota musí spĺňať



$$(\text{Plocha pod krivkou } f(s)) = A \int_0^{+\infty} s^{n-2} e^{-\frac{ns^2}{2\mu_2}} ds. \quad (2.8)$$

Nech

$$I_p = \int_0^{+\infty} s^p e^{-\frac{ns^2}{2\mu_2}} ds$$

a z (2.2) plynie, že

$$\mu_2 = \sigma^2.$$

Potom

$$\begin{aligned} I_p &= \frac{\sigma^2}{n} \int_0^{+\infty} s^{p-1} \frac{d}{dx} \left( -e^{-\frac{ns^2}{2\sigma^2}} \right) ds = \\ &= \frac{\sigma^2}{n} \left[ s^{p-1} \left( -e^{-\frac{ns^2}{2\sigma^2}} \right) \right]_{s=0}^{+\infty} + \frac{\sigma^2}{n} (p-1) \int_0^{+\infty} s^{p-2} e^{-\frac{ns^2}{2\sigma^2}} ds = \\ &= \frac{\sigma^2}{n} (p-1) I_{p-2}. \end{aligned}$$

Ak budeme ďalej pokračovať týmto rekurentným postupom, dostaneme

$$\begin{aligned} I_{n-2} &= \left( \frac{\sigma^2}{n} \right)^{\frac{n-2}{2}} (n-3)(n-5)\dots 3 \cdot 1 I_0 \quad \text{pre } n \text{ párne,} \\ &= \left( \frac{\sigma^2}{n} \right)^{\frac{n-3}{2}} (n-3)(n-5)\dots 4 \cdot 2 I_1 \quad \text{pre } n \text{ nepárne,} \end{aligned} \quad (2.9)$$

kde

$$I_0 = \int_0^{+\infty} e^{-\frac{ns^2}{2\sigma^2}} ds = \sqrt{\frac{\pi}{2n}} \sigma, \quad (2.10)$$

$$I_1 = \int_0^{+\infty} s e^{-\frac{ns^2}{2\sigma^2}} ds = \frac{\sigma^2}{n}. \quad (2.11)$$

Potom z (2.8), (2.9), (2.10) a (2.11) plynie

$$\begin{aligned} A &= \frac{(\text{Plocha pod krivkou } f(s))}{(n-3)(n-5)\dots 3 \cdot 1 \cdot \sqrt{\frac{\pi}{2}} \left( \frac{\sigma^2}{n} \right)^{\frac{n-1}{2}}} \quad \text{pre } n \text{ párne,} \\ &= \frac{(\text{Plocha pod krivkou } f(s))}{(n-3)(n-5)\dots 4 \cdot 2 \cdot \left( \frac{\sigma^2}{n} \right)^{\frac{n-1}{2}}} \quad \text{pre } n \text{ nepárne.} \end{aligned}$$

Dosadením do (2.7) dostaneme

$$f(s) = \frac{N}{(n-3)(n-5)\dots 3 \cdot 1} \cdot \sqrt{\frac{\pi}{2}} \left(\frac{n}{\sigma^2}\right)^{\frac{n-1}{2}} s^{n-2} e^{-\frac{ns^2}{2\sigma^2}} \quad \text{pre } n \text{ párne,}$$

$$= \frac{N}{(n-3)(n-5)\dots 4 \cdot 2} \cdot \left(\frac{n}{\sigma^2}\right)^{\frac{n-1}{2}} s^{n-2} e^{-\frac{ns^2}{2\sigma^2}} \quad \text{pre } n \text{ nepárne,}$$
(2.12)

kde  $N$  značí plochu pod krivkou  $f(s)$ . (V súčasnosti je známe, že  $N = 1$ , pretože vo vzorci (2.8) sa  $A$  volí tak, aby plocha pod krivkou  $f(s)$  bola rovná 1).

(2.12) je už konečný tvar hustoty  $s$  výberu z normálneho rozdelenia.

## 2.3 Korelácia medzi $\bar{X} - EX$ a $s$ výberu z normálneho rozdelenia

Predpokladajme opäť náhodný výber  $X_1, \dots, X_n$  z normálneho rozdelenia  $X_i \sim N(0, \sigma^2)$ . Vid' kap. 2.1.

Najskôr Student ukázal, že náhodné veličiny  $\bar{X} - EX \stackrel{EX=0}{=} \bar{X}$  a  $s$  sú nekorelované, argumentoval to tým, že je rovnako pravdepodobná pozícia výberového priemeru ( $\bar{X}$ ) naľavo aj napravo od počiatku 0. (t.j. že platí  $P(\bar{X} > 0) = P(\bar{X} < 0) = \frac{1}{2}$  resp. že sa jedná o symetrické rozdelenie.) V dnešnej dobe je známe, že Studentov argument nebol správny a neplatí. Vieme však, že  $\bar{X}$  a  $s^2$  sú nezávislé (vid' Andělovu knihu [3] str. 70) a z toho plynie nezávislosť  $\bar{X}$  a  $s$  resp. aj nekorelovanosť  $\bar{X}$  a  $s$ .

Student ešte zisťoval koreláciu medzi  $(\bar{X})^2$  a  $s^2$ , aby si bol istý výpočtami ďalej.

Označme

$$u^2 := (\bar{X})^2, \text{ kde } X = Y_i - \mu \quad \text{pre } i = 1, \dots, n.$$

Ďalej zavedieme premenné  $m_1$  a  $M_1$  nasledujúcim spôsobom ako

$$m_1 = Eu^2 = E \left( \frac{\sum_{i=1}^n X_i}{n} \right)^2 = \frac{E \sum_{i=1}^n X_i^2}{n^2} + \underbrace{\frac{2 \sum_{i=1}^n \sum_{j=1, j \neq i}^n EX_i X_j}{n^2}}_{\substack{z(2.5) \\ = 0}} = \frac{\mu_2}{n},$$

$$M_1 = Es^2 = \mu_2 \frac{(n-1)}{n}.$$

Spočítame

$$u^2 s^2 = \left( \frac{\sum_{i=1}^n X_i}{n} \right)^2 \left[ \frac{\sum_{i=1}^n X_i^2}{n} - \left( \frac{\sum_{i=1}^n X_i}{n} \right)^2 \right] =$$

$$= \left( \frac{\sum_{i=1}^n X_i^2}{n} \right)^2 + 2 \frac{\sum_{i=1}^n \sum_{j=1, j \neq i}^n X_i X_j \cdot \sum_{i=1}^n X_i^2}{n^2} - \frac{\sum_{i=1}^n X_i^4}{n^4} - \frac{6 \cdot \sum_{i=1}^n \sum_{j=1, j \neq i}^n X_i X_j}{n^4} - (\text{zvyšné})$$

členy, ktoré pri strednej hodnote zmiznú),

potom

$$Eu^2 s^2 = E \left( \frac{\sum_{i=1}^n X_i^2}{n} \right)^2 + 2E \frac{\sum_{i=1}^n \sum_{j=1, j \neq i}^n X_i X_j \cdot \sum_{i=1}^n X_i^2}{n^2} - E \frac{\sum_{i=1}^n X_i^4}{n^4} - E \frac{6 \cdot \sum_{i=1}^n \sum_{j=1, j \neq i}^n X_i X_j}{n^4} =$$

$$= \frac{\mu_4}{n^2} + \mu_2^2 \frac{(n-1)}{n^2} - \frac{\mu_4}{n^3} - 3\mu_2^2 \frac{(n-1)}{n^3}.$$

Keďže pre koreláciu medzi náhodnými veličinami  $u^2$  a  $s^2$  platí vzťah

$$\text{corr}(u^2, s^2) = \frac{\text{cov}(u^2, s^2)}{\sqrt{\text{var } u^2} \sqrt{\text{var } s^2}} = \frac{Eu^2 s^2 - Eu^2 Es^2}{\sqrt{\text{var } u^2} \sqrt{\text{var } s^2}} = \frac{Eu^2 s^2 - m_1 M_1}{\sqrt{\text{var } u^2} \sqrt{\text{var } s^2}},$$

dostávame, že

$$\text{corr}(u^2, s^2) \sqrt{\text{var } u^2} \sqrt{\text{var } s^2} + m_1 M_1 = Eu^2 s^2. \quad (2.13)$$

Po dosadení za  $m_1 M_1$  a  $Eu^2 s^2$  vo vzorci (2.13) s použitím (2.4) dostaneme

$$\text{corr}(u^2, s^2) \sqrt{\text{var } u^2} \sqrt{\text{var } s^2} + \frac{\mu_2}{n} \mu_2 \frac{(n-1)}{n} = \frac{\mu_4}{n^2} + \mu_2^2 \frac{(n-1)}{n^2} - \frac{\mu_4}{n^3} - 3\mu_2^2 \frac{(n-1)}{n^3},$$

$$\text{corr}(u^2, s^2) \sqrt{\text{var } u^2} \sqrt{\text{var } s^2} + \mu_2^2 \frac{(n-1)}{n^2} = \frac{3\mu_2^2}{n^2} + \mu_2^2 \frac{(n-1)}{n^2} - \frac{3\mu_2^2}{n^3} - 3\mu_2^2 \frac{(n-1)}{n^3},$$

$$\text{corr}(u^2, s^2) \sqrt{\text{var } u^2} \sqrt{\text{var } s^2} + \mu_2^2 \frac{(n-1)}{n^2} = \mu_2^2 \frac{(n-1)}{n^2}.$$

Teda

$$\text{corr}(u^2, s^2) \sqrt{\text{var } u^2} \sqrt{\text{var } s^2} = 0, \text{ čiže } \text{corr}(u^2, s^2) = 0$$

a vidíme, že veličiny  $u^2$  a  $s^2$  sú nekorelované.

**Poznámka** Student vo svojej práci [10] použil značenie  $R_{u^2, s^2}$  pre  $\text{corr}(u^2, s^2)$ ,

$$\sigma_{u^2} \text{ pre } \sqrt{\text{var } u^2},$$

$$\sigma_{s^2} \text{ pre } \sqrt{\text{var } s^2}.$$

## 2.4 Hľadanie hustoty nezávislej na $\sigma$

Vzorec  $f(s) = \frac{C}{\sigma^{n-1}} s^{n-2} e^{-\frac{ns^2}{2\sigma^2}},$

kde

$$C = \frac{N}{(n-3)(n-5)\dots 3 \cdot 1} \cdot \sqrt{\frac{2}{\pi}} \cdot n^{\frac{n-1}{2}} \quad \text{pre } n \text{ párne,}$$

$$= \frac{N}{(n-3)(n-5)\dots 4 \cdot 2} \cdot n^{\frac{n-1}{2}} \quad \text{pre } n \text{ nepárne,}$$

označuje hustotu smerodajnej odchýlky  $s$  náhodného výberu  $n$  jedincov populácie z normálneho rozdelenia ( $N(0, \sigma^2)$ ).

Rovnica  $f(x) = \frac{\sqrt{nN}}{\sqrt{2\pi}\sigma} e^{-\frac{nx^2}{2\sigma^2}}$  označuje hustotu aritmetického priemeru  $\bar{X}$  vyš-

šie zmieneného náhodného výberu. (Vzt'ah bol získaný z Airyho knihy [1].)

Nech

$$z = \frac{\bar{X}}{s},$$

kde

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n (Y_i - \mu) = \bar{Y} - \mu$$

a  $\mu$  označuje strednú hodnotu náhodnej veličiny  $Y_i$ .

Chceme zistiť rozdelenie náhodnej veličiny  $z$ . Na to použijeme vetu o transformácii náhodného vektoru.

**Veta (O transformácii náhodného vektoru)** *Nech náhodný vektor  $\mathbf{X} = (X_1, \dots, X_n)'$  má hustotu  $p$  vzhľadom k Lebesgueovej miere v  $\mathfrak{R}^n$ . Nech  $t$  je zobrazenie z  $\mathfrak{R}^n$  do  $\mathfrak{R}^n$ , ktoré je regulárne a prosté na takej otvorenej množine  $G$ , pre ktorú platí  $\int_G p(x)dx = 1$ . Označme  $\tau$  inverzné zobrazenie k  $t : G \rightarrow t(G)$ . Potom náhodný vektor  $\mathbf{Y} = t(\mathbf{X})$  má hustotu vzhľadom k Lebesgueovej miere a táto hustota je rovná*

$$q(\mathbf{y}) = \begin{cases} p[\tau(\mathbf{y})] |D_\tau(\mathbf{y})| & \text{pre } \mathbf{y} \in t(G), \\ 0 & \text{pre } \mathbf{y} \notin t(G). \end{cases}$$

*Dôkaz.* Vid' Andělova kniha [3] str. 49.

Vrátíme sa späť k nášmu problému, a tým je zistiť rozdelenie  $z = \frac{\bar{X}}{s}$ , kde  $\bar{X}$  budeme značiť ako  $x$ .

Pri počítaní združenej hustoty  $f(s, x)$  sa Student dopustil nepresnosti. Pretože podľa jeho výpočtov  $f(s, x) = f(s)f(x)$ , lenže to by platilo, pokiaľ by náhodné veličiny  $\bar{X}$  a  $s$  boli nezávislé. No Student dokázal len nekorelovanosť  $\bar{X}$  a  $s$  a potom nekorelovanosť  $\bar{X}^2$  a  $s^2$  (vid' kap. 2.3), a to neimplikuje nezávislosť náhodných veličín  $\bar{X}$  a  $s$  (vid' [1], str. 34). Teda podľa Studenta združená hustota náhodných veličín  $s$  a  $x$  je rovná

$$f(s, x) = \frac{C}{\sigma^{n-1}} s^{n-2} e^{-\frac{ns^2}{2\sigma^2}} \cdot \frac{\sqrt{nN}}{\sqrt{2\pi}\sigma} e^{-\frac{nx^2}{2\sigma^2}} = \frac{CN}{\sigma^n} \sqrt{\frac{n}{2\pi}} \cdot s^{n-2} e^{-\frac{n(s^2+x^2)}{2\sigma^2}}.$$

Ideme použiť vetu o transformácii náhodného vektoru.

Nech

$$t : \begin{pmatrix} s \\ x \end{pmatrix} \rightarrow \begin{pmatrix} z \\ w \end{pmatrix} = \begin{pmatrix} \frac{x}{s} \\ s \end{pmatrix}.$$

Keďže

$$z = \frac{x}{s}, \quad x = z \cdot w, \\ w = s,$$

pre inverzné zobrazenie  $\tau$  platí, že

$$\tau : \begin{pmatrix} z \\ w \end{pmatrix} \rightarrow \begin{pmatrix} w \\ z \cdot w \end{pmatrix},$$

kde  $s \in (0, \infty)$ ,  $x \in \mathfrak{R}$ , teda aj  $w \in (0, \infty)$ ,  $z \in \mathfrak{R}$ .

Najskôr overíme predpoklady vety.

Keďže

$$D_t = \det \begin{pmatrix} -\frac{x}{s^2} & \frac{1}{s} \\ 1 & 0 \end{pmatrix} = -\frac{1}{s} \neq 0, \quad s \neq 0,$$

$x \in \mathfrak{R}$ ,  $s \in (0, \infty) \Rightarrow G \subset \mathfrak{R} \times (0, \infty)$  je otvorená

a

$$\left. \begin{array}{l} \frac{\partial t_1}{\partial s} = -\frac{x}{s^2} \quad \frac{\partial t_1}{\partial x} = \frac{1}{s} \\ \frac{\partial t_2}{\partial s} = 1 \quad \frac{\partial t_2}{\partial x} = 0 \end{array} \right\} \text{spojité na } G,$$

dostávame, že zobrazenie  $t$  je regulárne na množine  $G \subset \mathfrak{R} \times (0, \infty)$ .

Teraz sa pozrieme na podmienku prostoty.

Keďže

$$z = \frac{x}{s} \quad \text{a} \quad w = s,$$

potom

$\forall (s, x) \in G, \forall (\tilde{s}, \tilde{x}) \in G$  platí, že ak  $[(s, x) \neq (\tilde{s}, \tilde{x})]$ , potom  $[t(s, x) \neq t(\tilde{s}, \tilde{x})]$ ,

teda zobrazenie  $t$  je prosté na množine  $G \subset \mathfrak{R} \times (0, \infty)$ .

Keďže predpoklady sú splnené, použijeme danú vetu.

Združená hustota náhodných veličín  $Z$  a  $W$  je

$$g_{ZW}(z, w) = f(\boldsymbol{\tau}(z, w)) |D_{\boldsymbol{\tau}}|, \quad (2.14)$$

kde

$$D_{\boldsymbol{\tau}} = \begin{pmatrix} \frac{\partial \tau_1}{\partial z} & \frac{\partial \tau_1}{\partial w} \\ \frac{\partial \tau_2}{\partial z} & \frac{\partial \tau_2}{\partial w} \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ w & z \end{pmatrix} = -w,$$

čiže po dosadení do (2.14) dostaneme

$$\begin{aligned} g_{ZW}(z, w) &= f(w, zw) |-w| = \frac{CN}{\sigma^n} \sqrt{\frac{n}{2\pi}} \cdot w^{n-2} e^{-\frac{n(zw^2+w^2)}{2\sigma^2}} w \\ &= \frac{CN}{\sigma^n} \sqrt{\frac{n}{2\pi}} \cdot w^{n-1} e^{-\frac{nw^2(z^2+1)}{2\sigma^2}} \end{aligned}$$

a hustotu  $g_Z(z)$  s použitím substitúcie

$$\begin{aligned} y &= \frac{nw^2(z^2+1)}{2\sigma^2}, \\ dy &= \frac{n(z^2+1)}{\sigma^2} w dw, \\ w &= \sqrt{\frac{2\sigma^2 y}{n(z^2+1)}} \end{aligned}$$

a gama funkcie

$$\Gamma(a) \stackrel{\text{def.}}{=} \int_0^{+\infty} x^{a-1} e^{-x} dx, \quad a > 0 \quad \text{a} \quad \Gamma(a+1) = a!, \quad \Gamma\left(\frac{1}{2}\right) \stackrel{\text{def.}}{=} \sqrt{\pi},$$

spočítame vzťahom

$$\begin{aligned} g_Z(z) &= \frac{CN}{\sigma^n} \sqrt{\frac{n}{2\pi}} \cdot \int_0^{+\infty} w^{n-1} e^{-\frac{nw^2(z^2+1)}{2\sigma^2}} dw = \\ &= \frac{CN}{\sigma^n} \sqrt{\frac{n}{2\pi}} \cdot \frac{\sigma^2}{n(z^2+1)} \int_0^{+\infty} \left[ \frac{2\sigma^2}{n(z^2+1)} \right]^{\frac{n-2}{2}} y^{\frac{n-1}{2}} e^{-y} dy \\ &= \frac{C}{\sqrt{2\pi}} N \cdot \frac{\sigma^{2+n-2}}{\sigma^n} \cdot \frac{\sqrt{n}}{n \cdot n^{\frac{n-2}{2}}} \cdot 2^{\frac{n-2}{2}} (z^2+1)^{-\left[\frac{n-2}{2}+1\right]} \Gamma\left(\frac{n-2}{2}+1\right) = \\ &= \frac{CN}{\sqrt{2\pi}} \cdot 2^{\frac{n-2}{2}} \cdot \frac{\Gamma\left(\frac{n}{2}\right)}{n^{\frac{n-1}{2}}} \cdot (z^2+1)^{-\frac{n}{2}}, \quad z \in \mathfrak{R}. \end{aligned} \quad (2.15)$$

C môžeme prepísať cez gama funkciu (integračným cvičením zo vzťahu (2.8) ) ako

$$C = \frac{N \cdot n^{\frac{n-1}{2}}}{2^{\frac{n-3}{2}}} \cdot \frac{1}{\Gamma\left(\frac{n-1}{2}\right)}.$$

Po dosadení do (2.15), dostaneme

$$g_z(z) = \frac{N^2 n^{\frac{n-1}{2}}}{\sqrt{2\pi} \cdot 2^{\frac{n-3}{2}} \Gamma\left(\frac{n-1}{2}\right)} \frac{2^{\frac{n-2}{2}} \Gamma\left(\frac{n}{2}\right)}{n^{\frac{n-1}{2}}} (1+z^2)^{-\frac{n}{2}} =$$

$$= \frac{N-1}{2} \frac{1}{n-3} \frac{n-2}{n-5} \dots \frac{5}{4} \cdot \frac{3}{2} (1+z^2)^{-\frac{n}{2}} \quad \text{pre } n \text{ nepárne,}$$

$$= \frac{N-1}{\pi} \frac{1}{n-3} \frac{n-2}{n-5} \dots \frac{4}{3} \cdot \frac{2}{1} (1+z^2)^{-\frac{n}{2}} \quad \text{pre } n \text{ párne.}$$

Pretože táto rovnica nám dáva nezávislosť na  $\sigma$ , dostali sme rozdelenie náhodnej veličiny  $z = \frac{\bar{X}}{s}$  pre akýkoľvek náhodný výber  $(X_1, \dots, X_n)$  z normálneho rozdelenia.

## 2.5 Vlastnosti hustoty $z$

Student chcel vytvoriť štatistickú tabuľku, z ktorej je možné vyčítať pravdepodobnosť, že výberový priemer  $\bar{X}$  vydelený  $s$   $\left(z = \frac{\bar{X}}{s}\right)$ , bude ležať medzi  $-\infty$  a  $z$  (čiže  $P(-\infty < \frac{\bar{X}}{s} < z) = F(z)$ ). Teda potreboval vypočítať distribučnú funkciu  $F(z)$ .

Tá sa spočíta ako

$$F(z) = K \cdot \int_{-\infty}^z \frac{1}{(1+x^2)^{\frac{n}{2}}} dx,$$

kde

$$K := \frac{1}{2} \frac{1}{n-3} \frac{n-2}{n-5} \dots \frac{5}{4} \cdot \frac{3}{2} \quad \text{pre } n \text{ nepárne,}$$

$$:= \frac{1}{\pi} \frac{1}{n-3} \frac{n-2}{n-5} \dots \frac{4}{3} \cdot \frac{2}{1} \quad \text{pre } n \text{ párne.}$$



S použitím substitúcie

$$x = \tan \theta, \quad \text{kde } \theta \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right],$$

$$dx = \frac{1}{\cos^2 \theta} d\theta,$$

kde medze spočítame ako

$$\theta = \arctan z,$$

$$\lim_{z \rightarrow -\infty} \arctan z = -\frac{\pi}{2},$$

$$\text{čiže } \theta \in \left[-\frac{\pi}{2}, \arctan z\right],$$

dostaneme

$$\begin{aligned} F(z) &= K \cdot \int_{-\frac{\pi}{2}}^{\arctan z} \frac{1}{(1 + \tan^2 \theta)^{\frac{n}{2}}} \cdot \frac{1}{\cos^2 \theta} d\theta = \\ &= K \cdot \int_{-\frac{\pi}{2}}^{\arctan z} \frac{1}{\left[\frac{1}{\cos^2 \theta}\right]^{\frac{n}{2}}} \cdot \frac{1}{\cos^2 \theta} d\theta = \\ &= K \cdot \int_{-\frac{\pi}{2}}^{\arctan z} [\cos^2 \theta]^{\frac{n}{2}-1} d\theta = \\ &= K \cdot \int_{-\frac{\pi}{2}}^{\arctan z} \cos^{n-2} \theta d\theta \end{aligned}$$

a to je možné dopočítať už pre akékoľvek  $z \in \mathfrak{R}$ .

Na základe toho Student zostavil *tabuľky distribučných funkcií*  $F(z)$ , kde

$$F(z) = \frac{n-2}{n-3} \frac{n-4}{n-5} \cdots \begin{pmatrix} \frac{3}{2} \cdot \frac{1}{2} & \text{pre } n \text{ nepárne} \\ \frac{2}{1} \cdot \frac{1}{\pi} & \text{pre } n \text{ párne} \end{pmatrix} \int_{-\frac{\pi}{2}}^{\arctan z} \cos^{n-2} \theta d\theta. \quad (2.16)$$

Pre porovnanie si vzal distribučnú funkciu  $F(x)$  normálneho rozdelenia so smerodajnou odchýlkou  $s = \sqrt{\frac{1}{7}}$  a  $n = 10$ :

$$F(x) = \frac{\sqrt{7}}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{7x^2}{2}} dx.$$

Čiže tabuľka vyzerá nasledovne.

2.5.1 Tabuľky distribučnej funkcie  $F(z)$  získané výpočtom zo vzťahu (2.16) pre hodnoty  $n = 10$  v porovnaní s distribučnou funkciou  $F(x)$  normálneho rozdelenia, kde

$$n = 10 \text{ a } s = \sqrt{\frac{1}{7}}.$$

$z$	$n = 10$	$F(x)$
0,1	0,61462	0,60411
0,2	0,71846	0,70159
0,3	0,80423	0,78641
0,4	0,86970	0,85520
0,5	0,91609	0,90691
0,6	0,94732	0,94375
0,7	0,96747	0,96799
0,8	0,98007	0,98253
0,9	0,98780	0,99137
1,0	0,99252	0,99820
1,1	0,99539	0,99926
1,2	0,99713	0,99971
1,3	0,99819	0,99986
1,4	0,99885	0,99989
1,5	0,99926	0,99999

**Poznámka:** Celú tabuľku je možné nájsť v Studentovom článku [10] str. 21.

## 2.6 Použitie tabuľky v praxi

Student urobil experiment na porovnanie účinnosti dvoch liečiv. Hodnoty daných liečiv získal z tabuliek A. R. Cushnyho a A. R. Peeblesa z časopisu *Journal of Physiology* z roku 1904. V nich bol zaznamenaný spánok 10 pacientov pred použitím hypnotických látok a po podaní (1) Dextro – hyoscyaminu hydromrobidu, (2) Laevo – hyoscyaminu hydrobromidu. Tabuľka 2.6.1 nám udáva počet hodín spánku získaných navyše u pacientov, ktorí užívali liek (1) a (2):

2.6.1 Tabuľka udávajúca počet hodín spánku získaných navyše užívaním lieku hyoscyaminu hydromrobidu získaná z [10].

Pacient	1 (Dextro-)	2 (Laevo-)	Rozdiel (2-1)
1.	+0,7	+1,9	+1,2
2.	-1,6	+0,8	+2,4
3.	-0,2	+1,1	+1,3
4.	-1,2	+0,1	+1,3
5.	-1,0	-0,1	0
6.	+3,4	+4,4	+1,0
7.	+3,7	+5,5	+1,8
8.	+0,8	+1,6	+0,8
9.	0	+4,6	+4,6
10.	+2,0	+3,4	+1,4
aritmetický priemer $\bar{X}$	+0,75	+2,33	+1,58
smerodatná odchýlka $s$	+1,70	+1,90	+1,17

Najskôr nás zaujíma pravdepodobnosť, že liek (1) v priemere dáva nárast spánku.

$$\text{Keďže } z = \frac{\bar{X}}{s} = \frac{+0,75}{1,70} = 0,44, \text{ hľadáme v tabuľkách 2.5.1 hodnotu distribučnej}$$

funkcie  $F(z)$  pre  $z = 0,44$  a  $n = 10$ . Tá nám vyšla niekde medzi 0,8697 a 0,9161. Student presne spočítal, že je to 0,8873. Inak povedané, šance sú 0,887 k 0,113. (Čiže pravdepodobnosť, že liek (1) zaznamená nárast spánku je približne 0,8873). Teda, je veľmi pravdepodobné, že liek (1) zvýši dĺžku spánku.

Analogicky spočítame pravdepodobnosť, že liek (2) v priemere spôsobuje nárast spánku. Keďže  $z = \frac{\bar{X}}{s} = \frac{+2,33}{1,90} = 1,23$  a opäť  $n = 10$ , z tabuliek 2.5.1 dostávame pravdepodobnosť rovnú 0,9974. (t.j. šance sú približne 400 k 1)

Z toho sa zdá byť pravdivé, že liek (2) je účinnejší ako liek (1). Ideme túto hypotézu otestovať tak, že odčítame dané hodnoty lieku (1) od lieku (2). (viď tabuľka 2.6.1 stĺpec „rozdiel (2-1)“).

Spočítame  $z = \frac{\bar{X}}{s} = \frac{+1,58}{1,17} = 1,35$  a pre  $n = 10$  z tabuľky 2.5.1 dostávame hod-

notu 0,9985. Čiže pravdepodobnosť, že liek (2) je účinnejší ako liek (1) je okolo 0,9985. (t.j. šanca je približne 666 k 1) Pri takej vysokej pravdepodobnosti môžeme už predpokladať, že liek (2) je určite účinnejší ako liek (1).

Ďalšie zaujímavé experimenty, ktoré Student uskutočnil je možné nájsť v jeho článku [10].

### 3 Objavené nepresnosti v Studentovom článku

Prvou Studentovou nepresnosťou pri uvažovaní bolo, že na ukázanie gama rozdelenia náhodnej veličiny  $s^2$ , spočítal prvé štyri momenty veličiny  $s^2$  a tie mu vyšli rovnako ako pri gama rozdelení, teda usúdil, že môžeme predpokladať gama rozdelenie náhodnej veličiny  $s^2$ . Hoci to tvrdenie všeobecne pre všetky rozdelenia neplatí (viď kap. 2.2 str. 14).

Ďalšieho omylu sa dopustil, keď vo vete o transformácii namiesto nezávislosti  $\bar{X}$  a  $s$ , ukázal nekorelovanosť  $\bar{X}$  a  $s$  a následne aj nekorelovanosť  $(\bar{X})^2$  a  $s^2$ . Čo v podstate neimplikuje nezávislosť náhodných veličín  $\bar{X}$  a  $s$ .

Student počítal s testovacou štatistikou

$$z = \frac{\bar{X} - \mu}{s},$$

kde

$$s^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

V dnešnej dobe je už zrejماً testovacia štatistika

$$t = \sqrt{n} \frac{\bar{X} - \mu}{S},$$

kde

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Na ňu však prišiel v roku 1925 štatistik Ronald A. Fisher a to tak, že použil transformáciu  $t = z\sqrt{n-1}$ , ktorú si neskôr osvojil Student.

Ďalším významným článkom, ktorým sa Student zaoberal a ktorému sa budeme v krátkosti venovať je *Pravdepodobná chyba korelačného koeficientu* [9].

## 4 Pravdepodobná chyba korelačného koeficientu

Student sa pokúsil týmto článkom zaoberať problémom významu korelačných koeficientov odvodených z malých výberov a taktiež chcel motivovať matematikov, ktorí majú čas aj schopnosti sa daným problémom tiež zaoberať.

Predpokladá sa náhodný výber  $(X_1, Y_1)^T, \dots, (X_n, Y_n)^T$  z populácie (kde  $(X_i, Y_i)^T$  sú stĺpcové vektory) s dvojrozmerným normálnym rozdelením  $N_2\left((\mu_X, \mu_Y)^T, \begin{pmatrix} \sigma_X^2 & \text{cov}(X, Y) \\ \text{cov}(X, Y) & \sigma_Y^2 \end{pmatrix}\right)$ .

Definujme výberový korelačný koeficient  $R$  vzorcom

$$R = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}.$$

(Vid' [3] str. 93).

Ďalej definujme korelačný koeficient  $r$  ako

$$r = \frac{\text{cov}(X_i, Y_i)}{\sqrt{(\text{var } X_i)(\text{var } Y_i)}}.$$

My budeme zisťovať pravdepodobnosť, že  $r$  leží medzi danými dvoma medzami. (Tzv. *interval spoľahlivosti*  $P(D \leq r \leq H)$ , kde  $D$  je dolná hranica a  $H$  je horná hranica).

Aby sme vyriešili tento problém, musíme vedieť dve veci: (1) rozdelenie hodnôt  $R$  získaných z výberu populácie (ktorá nám dáva  $r$ ) a (2) *apriórnu pravdepodobnosť* (vid' kap. 1.3), že  $r$  leží medzi danými dvoma medzami.

Student uvažoval, že rozdelenie  $R$  by mohlo mať len jednu z dvoch hustôt. A to buď  $R$  ma rovnomerné rozdelenie  $R[-1, 1]$ , potom hustota  $R$  je

$$f(r) = \begin{cases} \frac{1}{2}, & \text{ak } r \in (-1, 1), \\ 0, & \text{inak.} \end{cases}$$

Alebo hustota náhodnej veličiny  $R$  môže byť rovná

$$f(r) = C \cdot (1 - r^2), \quad \text{ak } r \in (-1, 1), \\ = 0, \quad \text{inak.}$$

V súlade s jeho pokusmi a skúsenosťami stanovil apriórne rozdelenie ako

$$f(r) = \frac{3}{4} \cdot (1 - r^2), \quad \text{ak } r \in (-1, 1), \\ = 0, \quad \text{inak.}$$

Pre veľké výbery štatistickí Pearson a Filon ukázali, že rozdelenie  $R$  odpovedá asymptoticky normálnemu rozdeleniu. No malými výbermi sa nik nezaoberal.

Student sa snažil tento problém vyriešiť a spravil nasledujúci experiment. Vzal 3000 hodnôt výšky postavy a dĺžky prostredného prsta na ľavej ruke kriminálnikov a urobil náhodný výber po štyroch z tejto populácie. Čím dostal 750 hodnôt  $R$  z populácie, ktorej skutočný korelačný koeficient  $r$  bol 0,66.

Keď vzal výšky postáv z nejakého výberu a dĺžky prostredných prstov ľudí iného výberu, dostal 750 hodnôt  $R$  z populácie, ktorej skutočný korelačný koeficient  $r$  bol rovný 0 (pretože výška jedného jedinca a dĺžka prostredného prsta iného jedinca, sú nezávislé náhodné veličiny, čiže korelácia je nulová).

Na základe získaných údajov z experimentov, s ktorými presnejšie pracuje v jeho článku [9] a matematickou úvahou prišiel na to, že ak je korelácia medzi dvoma normálne rozdelenými náhodnými veličinami nulová, hustota

$$f(r) = r_0 (1 - r^2)^{\frac{n-4}{2}}, \quad \text{kde } n \geq 3,$$

nám dáva presné rozdelenie výberového korelačného koeficientu  $R$  získaného z náhodného výberu  $n$  jednotlivcov.  $r_0$  spočítame z rovnice

$$r_0 \int_{-1}^1 (1 - r^2)^{\frac{n-4}{2}} dr = N,$$

kde  $N$  je ako vo vzťahu (2.12) v kapitole 2. 2. V súčasnosti je známe, že  $N$  je rovné 1, pretože  $r_0$  sa volí tak, aby to bola normalizujúca konštanta, čiže aby platilo, že hustota

$$f(r) = r_0 \int_{-1}^1 (1 - r^2)^{\frac{n-4}{2}} dr = 1. \quad (4.1)$$

Spočítame teraz integrál zo vzťahu (4.1). Pred tým si ale potrebujeme zadefinovať pojem *beta funkcia*.

Pre  $a > 0$ ,  $b > 0$  definujeme beta funkciu  $B(a, b)$  vzorcom

$$B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx.$$

Medzi gama funkciou a beta funkciou platí vzťah

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}. \quad (4.2)$$

Vrátime sa späť k integrálu  $\int_{-1}^1 (1-r^2)^{\frac{n-4}{2}} dr$ :

$$\int_{-1}^1 (1-r^2)^{\frac{n-4}{2}} dr = 2 \int_0^1 (1-r^2)^{\frac{n-4}{2}} dr.$$

Použijeme substitúciu

$$\begin{aligned} y &= r^2, \\ \sqrt{y} &= r, \\ \frac{1}{2} \frac{1}{\sqrt{y}} dy &= dr, \end{aligned}$$

a dostaneme

$$I = \int_{-1}^1 (1-r^2)^{\frac{n-4}{2}} dr = \int_0^1 y^{-1/2} (1-y)^{\frac{n-4}{2}} dy = \int_0^1 y^{(1/2)-1} (1-y)^{\frac{n-2}{2}-1} dy = B\left(1/2, \frac{n-2}{2}\right),$$

kde  $n \geq 3$ .

Čiže nakoniec

$$I = B\left(1/2, \frac{n-2}{2}\right) = \frac{\Gamma(1/2)\Gamma\left(\frac{n-2}{2}\right)}{\Gamma\left(\frac{n-2+1}{2}\right)} = \sqrt{\pi} \frac{\Gamma\left(\frac{n-2}{2}\right)}{\Gamma\left(\frac{n-1}{2}\right)}, \quad n \geq 3.$$

Teda dostávame

$$r_0 = \frac{N}{\int_{-1}^1 (1-r^2)^{\frac{n-4}{2}} dr} = \frac{1}{\frac{\sqrt{\pi} \Gamma\left(\frac{n-2}{2}\right)}{\Gamma\left(\frac{n-1}{2}\right)}} = \frac{\Gamma\left(\frac{n-1}{2}\right)}{\sqrt{\pi} \Gamma\left(\frac{n-2}{2}\right)}, \quad \text{kde } n \geq 3.$$



Takže výsledná hustota výberového korelačného koeficientu  $R$  pri výbere o rozsahu  $n \geq 3$  z dvojrozmerného normálneho rozdelenia s nulovým korelačným koeficientom  $r$  je rovná

$$f(r) = \frac{\Gamma\left(\frac{n-1}{2}\right)}{\sqrt{\pi} \Gamma\left(\frac{n-2}{2}\right)} (1-r^2)^{\frac{n-4}{2}}.$$

V ďalšej kapitole sa budeme venovať štatistikom, ktorí na Studentovu tvorbu neskôr naväzovali.

## 5 Štatici, ktorých ovplyvnila Studentova tvorba

V období pred Fisherom nikto nevenoval veľkú pozornosť Studentovým článkom. Prvým významným štatistikom, ktorého ovplyvnili Studentove práce, bol už vyššie zmienený Ronald A. Fisher. Ten nielen oživil Studentove články, ale tiež v nich opravil chyby a nedostatky. Keď sa hovorí o „Studentovom t-rozdelení“, máme v podstate namysli skôr Fisherove myšlienky ako Studentove. Napríklad t-rozdelenie o  $n - 1$  stupňov voľnosti, ako rozdelenie založené na normálnom rozdelení, no tiež jeho aplikácia v regresnej analýze, boli výplodom Fisherovho myslenia. Aplikáciu „Studentovho t-rozdelenia“ Fisher zahrnul vo svojom článku [4]. Student sa vo svojich prácach venoval len jednovýberovými problémami, no Fisher v r. 1922 ukázal aplikáciu t-rozdelenia aj na dvojjvýberové problémy, kde v oboch výberoch predpokladal rovnaký rozptyl náhodných veličín.

t-rozdelenie bolo neskôr v roku 1931 Haroldom Hotellingom (matematickým štatistikom a významným ekonomickým teoretikom) zovšeobecnené aj na testovanie vektoru stredných hodnôt s mnohorozmerným normálnym rozdelením.

Studentove práce sa dostali do povedomia aj štatistikom z ruskej ríše. Významným štatistikom, ktorého zaujala Studentova tvorba, bol Evgenii Evgenevič Sluckij. Jeho meno sa stalo významné v oblasti ekonometrie a matematickej štatistiky, hlavne kvôli štúdiu stochasticky závislých náhodných procesov. Sluckij venoval dlhé roky štúdiu publikácií Karla Pearsona, čím sa dostal aj k prácam Studenta. Neskôr napísal knihu [8], ktorej názov nám už sám o sebe napovedá, že sa v nej zaoberal sčasti (a dokonca z väčšej časti) teóriou korelácie, špeciálne pravdepodobnej chyby. To ho počas písania priviedlo k Studentovmu článku [9], z ktorého potom vo svojej knihe citoval niektoré výsledky a závery.

Ďalšími osobnosťami z Ruska, ktorí venovali pozornosť hlavne druhému Studentovmu článku [9], boli Alexander Alexandrovič Čuprov (štatistik) a Andrej Andrejevič Markov (matematik). Tí sa zaoberali prevažne šiestimi typmi Pearsonových kriviek (tie je možné nájsť v Pearsonovom článku [7]), čo v podstate do istej miery súviselo aj s t-rozdelením Studentovej štatistiky.

## 6 Návrat k súčasnosti

### 6.1 Studentovo t-rozdelenie

V súčasnej štatistike sa Studentovo t-rozdelenie zavádza nasledovne:

**Veta (súčasnosť)** *Nech  $X_1, \dots, X_n$  je náhodný výber z  $N(\mu, \sigma^2)$ , kde  $n \geq 2$  a  $\sigma^2 > 0$ , nech*

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2} \quad a \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i,$$

*potom*

$$T = \sqrt{n} \frac{\bar{X} - \mu}{S},$$

*kde  $\bar{X}$  a  $S$  sú nezávislé náhodné veličiny, má  $t$  – rozdelenie o  $n - 1$  stupňov voľnosti s hustotou*

$$f(t) = \frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{n-1}{2}\right) \sqrt{\pi(n-1)}} \left(1 + \frac{t^2}{n-1}\right)^{-\frac{n}{2}}. \quad (6.1)$$

*Dôkaz.* Vid' Andělova kniha [3] str. 74.

Pre porovnanie uvedieme, ako by asi sformuloval vetu Student pred 100 rokmi, na základe zmienených faktov v jeho článku.

**Veta (Student pred 100 rokmi)** *Nech  $X_1, \dots, X_n$  je náhodný výber z  $N(0, \sigma^2)$ , kde  $n \geq 2$ , ( $X_i$  je definované ako v kapitole 2.2), nech*

$$S = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2} \quad a \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i,$$

*potom*

$$T = \frac{\bar{X}}{S},$$

kde  $\bar{X}$  a  $S$  sú nekorelované náhodné veličiny, má hustotu

$$f(t) = \frac{N^2 n^{\frac{n-1}{2}}}{\sqrt{2\pi} \cdot 2^{\frac{n-1}{2}} \Gamma\left(\frac{n-1}{2}\right)} \frac{2^{\frac{n-2}{2}} \Gamma\left(\frac{n}{2}\right)}{n^{\frac{n-1}{2}}} (1+t^2)^{-\frac{n}{2}} =$$

$$= \frac{\Gamma\left(\frac{n}{2}\right)}{\sqrt{\pi} \Gamma\left(\frac{n-1}{2}\right)} (1+t^2)^{-\frac{n}{2}}.$$
(6.2)

Všimnime si, že súčasný vzorec sa veľmi nelíši od toho spred 100 rokov. V podstate sa odlišuje od toho Studentovho len o premennú  $n-1$ . Už vyššie sme spomenuli, že túto nepresnosť opravil Fisher v roku 1925 a dospel tým k hustote, s ktorou sa počíta v štatistike dodnes.

## 6.2 Výberový korelačný koeficient

Student vo svojom článku [9] prišiel na to, že hustota výberového korelačného koeficientu  $R$  pri náhodnom výbere o rozsahu  $n \geq 3$  z dvojrozmerného normálneho rozdelenia s nulovým korelačným koeficientom  $r$  by mala byť daná vzorcom

$$f(r) = r_0 (1-r^2)^{\frac{n-4}{2}},$$
(6.3)

kde

$$r_0 = \frac{\Gamma\left(\frac{n-1}{2}\right)}{\sqrt{\pi} \Gamma\left(\frac{n-2}{2}\right)}.$$

V súčasnej štatistike (vid' Andělovu knihu [3]) je jasné, že hustota  $R$  je daná rovnicou

$$f(r) = \frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n-2}{2}\right) \sqrt{\pi}} (1-r^2)^{\frac{n-4}{2}},$$
(6.4)

kde

$$-1 < r < 1 \text{ a } n \geq 3.$$

Čiže sa Student vo výpočte nemýlil.

## 7 Záver

Student vo svojich prácach urobil pár vážnejších chýb, s ktorými jeho počítanie nebolo celkom korektné, no aj tak mu patrí obrovský obdiv. Svojimi výpočtami a závermi ovplyvnil množstvo štatistikov a matematikov, ktorí nakoniec Studentove postupy vylepšili a opravili natoľko, že našli využitie v súčasnej štatistike. Taktiež je potrebné spomenúť meno Arthur Guinness, ktorý Studentovi umožnil stráviť dva semestre akademického roku v Pearsonovom biometrickom laboratóriu na univerzite v Londýne, kde sa Student veľa vecí naučil. Ďalšou dôležitou osobou bol Ronald A. Fisher, bez ktorého by Studentove práce ostali bez povšimnutia.

No aj napriek vážnejším chybám sa v podstate Studentove výsledky až tak nelíšia od tých súčasných. Stačí porovnať hustotu (6.1) so (6.2) a hustotu (6.3) so (6.4), kde jasne vidieť len veľmi malé rozdiely. Ak však vezmeme do úvahy fakt, že pred 100 rokmi neexistovali žiadne počítače a ani štatistika nebola natoľko vyvinutá ako dnes a že Student bol v podstate zamestnanec pivovaru, určite prijmem názor, že Studentovi patrí obrovská pocta a Studentovo t-rozdelenie nesie právom jeho meno.

# Literatúra

- [1] Airy, G. B.: *On the Algebraical and Numerical Theory of Errors of Observations and the Combination of Observations*, Macmillan and Co., Cambridge and London 1861.
- [2] Anděl, J.: *Matematická statistika*, SNTL, Praha 1985.
- [3] Anděl, J.: *Základy matematické statistiky*, Matfyzpress, Praha 2005.
- [4] Fisher, R. A.: *Statistical Methods for Research Workers*, Oliver & Boyd., Edinburgh and London 1925.
- [5] Hotteling, H.: *British Statistics and Statisticians Today*, Journal of the American Statistical Association **25** (1930), 186 – 190.
- [6] Pearson, E. S., Plackett R. L., Barnard G. A.: *‘Student’, A Statistical Biography of William Sealy Gosset*, Oxford University Press, Oxford 1990.
- [7] Pearson, K.: *Contributions to the mathematical theory of evolution, II: Skew variation in homogeneous material*, Philosophical Transactions of the Royal Society of London, A **186** (1895), 343 – 414.
- [8] Slutsky, E. E.: *The Theory of Correlation and Elements of the Study of Curves of Distribution*, Izvestiia Kievskago Kommercheskago Instituta, Kiev 1912.
- [9] Student: *Probable Error of a Correlation Coefficient*, Biometrika **6** (1908), 302 – 310.
- [10] Student: *The Probable Error of a Mean*, Biometrika **6** (1908), 1 – 25.