

Jak vytváří statistika obrazy světa a života. II. díl

Část I. Teorie náhodných výběrů [odst. 1, 2,1-2,5,
3,1-3,11, 4,1-4,5]

In: Jaroslav Janko (author): Jak vytváří statistika obrazy světa a života. II. díl. (Czech). Praha: Jednota českých matematiků a fyziků, 1944. pp. 7–82.

Terms of use:

Persistent URL: <http://dml.cz/dmlcz/403076>

© Jednota českých matematiků a fyziků

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

ČÁST I.

TEORIE NÁHODNÝCH VÝBĚRŮ.

(Znak kvantitativní.)

(1) Úvod.

V každodenních záležitostech všedního života a v pracích svého zaměstnání jsme vedeni k tvoření úsudků, které zveřejňují poznatky získané na jistém omezeném počtu pozorovaných případů. Když hospodyně koupila deset housek v určitém obchodě a shledala, že pět z nich není tak čerstvých jak by si bylo přáti, rozhodne se, že je bude příště kupovati jinde. Když někdo čekal pětkrát v týdnu deset minut nebo déle na vůz elektrické dráhy na stanici, činí závěr, že dopravní možnost je poměrně chudá. Lékař vezme kapku krve pacientovy, rozředí ji a počítá pak pod mikroskopem krvinky v nepatrné částce zředěného roztoku, aby na podkladě tohoto materiálu poznal podstatné vlastnosti krve pacientovy pro svou diagnosu. Prodává-li se pšenice podle toho, jak je bohatá na lepek, posuzuje se podle malých množství jako vzorků vzatých z celé sklizně. Také obchod jiným zbožím se provádí odedávna pomocí vzorků. Jsou tedy v praktickém životě časté případy, kdy je nutno usuzovati na vlastnosti jistého souboru prvků, jež nemůžeme všechny pozorovat, na základě toho, co jsme zjistili na nějakém menším množství prvků z něho, nebo jak často říkáme na vzorku. Dospěli jsme zkušeností a intuicí k této víře, že nám může vzorek něco povědět o celém základním množství, z něhož byl vzat a že nám to může povědět tím lépe, čím je větší. Jestliže výrobce ložiskových kuliček kontroluje své výrobky na nejvyšší možnou zatížitelnost podle vzorků a usuzuje z toho na jakost celé výroby, činí totéž co experimentátor, který vychází z předpokladu, že je možno usuzovati z následků na příčiny a z pozorovaných zvláštních případů odvozovati věty obecněji platné. Říkáme

tomu v logice indukce čili závěr postupující od zvláštního k obecnému nebo úsudek z výběru na základní soubor.

Je to úsudek, kterým se rozšiřují výsledky odvozené z pozorovaného souboru na rozsáhlejší soubor obsahující prvky, jež nebyly v původně studovaném souboru. Soustředíme-li informaci o pozorovaném souboru do několika charakteristik, můžeme pak usuzovati podle těchto hodnot na hodnoty příslušných parametrů v rozsáhleším základním souboru. Metody, jimiž docilujeme zobecnění statistických výsledků odvozených z výběrů, nazýváme statistickou indukcí.

Bylo by možno uvést velkou řadu případů, kdy je možno nebo nutno dosáhnouti veškeré informace o základním souboru podle jednoho malého náhodného výběru či několika málo výběrů. Víme, že se mnohdy provede ve fyzice osm až dvanáct pokusů, aby se odvodil přírodní zákon určitého jevu; nebo, botanik zkoumá charakteristiky malého výběru několika rostlin určitého druhu a přisuzuje pak tomuto druhu vlastnosti, které získal z výběru. Podobně závisí předpovídání vývoje obchodu na informaci, která vyplyne z výběru. Možno říci, že skoro každá empirická formule je získána pomocí výběrových dat. Důvody užívání těchto metod mohou spočívat také v nutnosti úspory času a nákladu, jehož vyžadují taková šetření a zkoušky. Při zkoušení nebo kontrole některých vlastností předmětů musí býti dotyčné předměty zničeny nebo znehodnoceny jako je tomu na př. při zkoušení citlivosti fotografických desek, délky života žárovek, pevnosti trubek a pod.

Připouštíme zajisté, že tento postup od zvláštního k obecnému musí obsahovati jakési prvky nejistoty. To ovšem neznamená, že závěry statistické indukce nejsou zcela přesné, ježto máme prostředky, které jsou s to vyjádřiti přesně povahu a stupeň oné nejistoty. Doklady toho máme v aplikaci teorie pravděpodobnosti. Jednotlivý případ je sice nejistý, ale mohou býti odvozeny přesnou dedukcí pravděpodobnosti různých možných jevů nebo jejich kombinací.

Skutečnost, že v úsudcích, k nimž jsme dospěli indukcí, je jakási nejistota, nevylučuje možnost zcela přesných a jednoznačných závěrů. V interpretaci odvozených pravděpodobností mívá někdy původ nejistota, s níž se setkáváme. Logický základ a soustavu měr pro vyvozování úsudků statistickou indukcí nám poskytuje teorie náhodného výběru. Hlavní výsledky a aplikaci její v případě alternativního znaku jsme vyložili již v I. díle.

(2) Náhodné výběry ze známého základního souboru.

Abychom odvodili výsledky upotřebitelné co nejlépe v praxi, tedy za podmínek co nejméně omezujících, budeme nejprve zkoumat, jaké hodnoty charakteristik poskytují všechny možné výběry ze známého základního souboru a vyslovíme tedy konkrétněji první úkol teorie výběru. Známe základní soubor čili jeho rozdělení četností, které je vyjádřeno několika parametry. Budeme uvažovati první tři parametry (průměr, směrodatnou odchylku a šikmost). Tímto základním souborem je podána informace o nějakém studovaném jevu. Chceme nyní popsat tento jev charakteristikami všech výběrů určitého rozsahu, které mohou vzniknouti ze základního souboru. V našem případě to tedy budou výběrové průměry, směrodatné odchylky a šikmosti, jejichž rozdělení četností budeme zkoumat; abychom našli souvislost mezi nimi a momenty základního souboru.

Známe rozdělení četností uvažovaného základního souboru (první dva sloupce tab. 1).

Parametry tohoto rozdělení četností, odvozené z prvních tří momentů budou

$$\begin{aligned} \bar{x} = \mu'(x, 1) &= 7,0000, \quad \sigma(x) = \sqrt{\mu(x, 2)} = 2,002, \\ \alpha(x, 3) &= \frac{\mu(x, 3)}{\sigma^3(x)} = 0,0000. \end{aligned} \quad (1)$$

Tabulka 1.

x_i	$n(x_i)$	\bar{x}_i	$n(\bar{x}_i)$	$f(\bar{x}_i)$
(1)	(2)	(3)	(4)	(5)
0	1	6,625—6,674	0	0,000
1	2	6,675—6,724	2	0,005
2	9	6,725—6,774	12	0,030
3	28	6,775—6,824	20	0,050
4	66	6,825—6,874	22	0,055
5	121	6,875—6,924	46	0,115
6	175	6,925—6,974	57	0,142 ₅
7	197	6,975—7,024	66	0,165
8	175	7,025—7,074	55	0,137 ₅
9	121	7,075—7,124	50	0,125
10	66	7,125—7,174	40	0,100
11	28	7,175—7,224	15	0,037 ₅
12	9	7,225—7,274	11	0,027 ₅
13	2	7,275—7,324	3	0,007 ₅
14	1	7,325—7,374	1	0,002 ₅
Součet	1001	Součet	400	1,0000

Ze základního souboru tohoto rozdělení četností vezmeme 400 výběrů, z nichž každý bude rozsahu $r = 200$ prvků. Můžeme to provést na př. tak, že napíšeme pro každý prvek základního souboru lístek (celkem 1001), na němž je poznamenána hodnota znaku a vytáhneme z nich 200 lístků, které tvoří jeden výběr. Můžeme postupovati buď tak, že každý jednotlivý lístek hned zase vrátíme zpět, nebo vrátíme teprve vytažených 200 najednou. Tuto okolnost však zatím necháme stranou.

Vypočítáme pro každý výběr rozsahu 200 prvků průměr výběrový \bar{x}_i , takže dostaneme 400 hodnot \bar{x}_i ($i = 1, 2, \dots, 400$), které budeme považovati za prvky nového souboru rozsahu $n = 400$; je to tedy soubor výběrových průměrů. Každý z těchto výběrových průměrů můžeme považovati za odhad parametru \bar{x} základního souboru. Soubor výběrových průměrů má své určité rozdělení četností, které

si uvedeme v třídách intervalech délky $h = 0,05$, jejichž střed označíme ${}_0\bar{x}_i$, takže dostaneme sloupec 3 a 4, tab. 1.

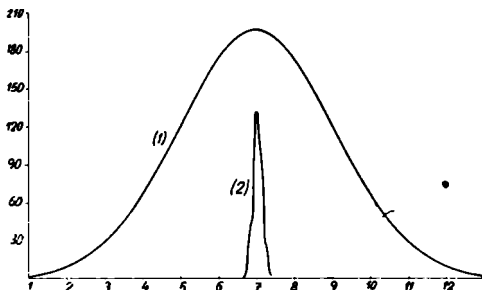
První tři charakteristiky tohoto rozdělení četností jsou

$${}_0\bar{x}_P = 7,005, \quad \sigma_P = 0,121, \quad \bar{\alpha}_3 = 0,0001. \quad (2)$$

Porovnáme-li nyní výsledky, jež jsme dostali pro parametry základního souboru a pro charakteristiky souboru výběrových průměrů, vidíme, že průměr základního souboru 7,0000 je velmi blízko průměru výběrových průměrů 7,005. Směrodatná odchylka v základním souboru je 2,002, kdežto v rozdělení výběrových průměrů je jen 0,121. Tato malá hodnota směrodatné odchylky rozdělení výběrových průměrů ukazuje, že všechny výběrové průměry leží velmi blízko svého průměru a také průměru základního souboru. Variační obor proměnné x je v základním souboru $x_{15} - x_1 = 14$, kdežto pro rozdělení výběrových průměrů je pouze ${}_0\bar{x}_{15} - {}_0\bar{x}_1 = 0,70$, což opět ukazuje, jak blízko leží hodnoty proměnné v rozdělení výběrových průměrů kolem průměru jejich a kolem průměru základního souboru. V tomto příkladě jsou všechny výběrové průměry ve vzdálenosti nejvýše 0,35 od průměru základního souboru. Porovnáme-li tuto hodnotu se směrodatnou odchylkou $\sigma_P = 0,121$, vidíme, že podle počtu pravděpodobností ([1], str. 38) je téměř nemožno dostati průměr výběru rozsahu $r = 200$, aby se lišil od průměru základního souboru o 0,5. Ukazuje tudíž tento příklad, že průměr z dosti rozsáhlého výběru je asi tak dobrý jako průměr základního souboru.

Není-li tedy základní soubor znám, je průměr z jednoho výběru nebo z několika málo výběrů asi tak dobrý, jak je vůbec možno získat. Nedáme se proto do velké práce s obstaráváním dalších výběrů nebo rozsáhlejšího výběru. Výsledky a právě provedené úvahy vidíme snadno v grafickém znázornění relativních četností obou rozdělení (obr. 1). Kdybychom utvořili průměr všech možných průměrů výběrových, rovnal by se přesně průměru základního souboru, jak níže dokážeme.

Zcela obdobně bychom nyní mohli na tomto příkladě zjistiti charakteristiky souboru 400 výběrových rozptylů nebo směrodatných odchylek a souboru výběrových šikmostí. Je to však postup zcela obdobný, proto se budeme



Obr. 1. Rozdělení četností základního souboru (1). Rozdělení četností výběrových průměrů (2).

věnovati hned řešení obecnému, kde jej provedeme pro všechny tyto soubory výběrových charakteristik. Na předcházejícím příkladu vidíme tedy, že při numerickém popisu daného rozdělení četností musíme rozeznávat dvě hlediska. Jednak můžeme považovat popis daného rozdělení za cíl pro sebe, jednak jej můžeme uvažovati jako výběr, reprezentující větší soubor, t. zv. základní. Obyčejně je důležité toto druhé hledisko, neboť v případech, kdy je buď nepraktické nebo nemožné pozorovati nebo měřiti studovaný znak na všech prvcích základního souboru, musíme z něho vzítí výběr a podle něho usuzovati na celý základní soubor.

Přistoupíme tedy k obecnému řešení otázky jak dobře výběr popisuje základní soubor za předpokladu, že základní soubor má rozsah N a jsou z něho brány výběry rozsahu r , při čemž tedy $r \leq N$, takže můžeme dostati celkem $\binom{N}{r}$ různých výběrů. Každý z těchto výběrů má nějaký první mo-

ment čili průměr, takže soubor těch $\binom{N}{r}$ výběrů dává rozdělení četností všech výběrových průměrů. Podobně má každý výběr druhý moment, takže soubor druhých výběrových momentů má také své rozdělení četností; stejně je tomu pro třetí a další momenty výběrové. Bude tudíž naším úkolem nyní srovnati momenty počítané z výběru s momenty základního souboru.

Je-li N konečné, říkáme, že tvoříme výběry z konečného základního souboru (na př. obyvatelé nějakého státu). Je-li N nekonečné, tvoříme výběry z nekonečného základního souboru (barometrický tlak v různých bodech atmosféry). V mnohých případech je N tak veliké, že můžeme prakticky považovati základní soubor za nekonečný a do výsledků tím nepronikne chyba, která by musela býti uvažována. Jsou také případy, kdy nevíme bezpečně je-li studovaný základní soubor konečný či nekonečný, jak je tomu na př. u souboru hvězd. Statistikové se dříve zabývali hlavně případem nekonečného základního souboru, ježto algebraické výpočty zvláště pro rozdělení vyšších momentů z konečného základního souboru jsou značně složité. Poněvadž se zde omezujeme na rozdělení četností několika prvních momentů, odvodíme výsledky pro N konečné a necháme-li pak v nich růsti N do nekonečna, vyplynou snadno zjednodušené výsledky pro výběry z nekonečného základního souboru.

(2, 1) Momenty rozdělení četností výběrových průměrů.

Průměr výběrových průměrů. Pro konečné N a r můžeme vybrati ze základního souboru, jehož prvky mají pozorované hodnoty proměnné x_1, x_2, \dots, x_N , celkem $\binom{N}{r} = \nu$ různých výběrů, z nichž každý má r prvků ($r \leq N$) a hodnota x_i se v něm vyskytuje jen jednou. Můžeme je nějakým způsobem seřadit, takže dostaneme ν výběrových průměrů

$\bar{x}_1, \bar{x}_2, \dots, \bar{x}_i, \dots, \bar{x}_\nu$. Bude tedy $\bar{x}_i = \frac{1}{r} \sum^{r,i} x_l$, kde součet $\sum^{r,i}$ značí součet všech r hodnot znaku i -tého výběru. Abychom stanovili jejich průměr $\mu'_1 = \frac{1}{\nu} (\bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_\nu)$ uvážíme, že počet prvků ve všech ν výběrech je dohromady $\nu \cdot r$, kdežto počet různých prvků je dán počtem prvků N základního souboru. Vystupuje tedy jeden určitý prvek v tomto celku všech výběrů $\frac{\nu r}{N}$ krát. To platí pro každý z N prvků základního souboru, takže součet všech hodnot znaku ve všech výběrech bude $(x_1 + x_2 + \dots + x_N) \frac{\nu r}{N}$, a poněvadž je těchto prvků celkem νr , je tudíž průměr

$$\mu'_1 = (x_1 + x_2 + \dots + x_N) \frac{1}{N} = \bar{x}. \quad (3)$$

Dostáváme tak první důležitou větu:

Průměr μ'_1 , čili první moment výběrových průměrů \bar{x}_i se rovná průměru základního souboru \bar{x} .

Rozptyl výběrových průměrů. Další charakteristikou rozdělení výběrových průměrů je jejich rozptyl μ_2 . Druhý moment výběrových průměrů, vzhledem k jejich průměru \bar{x} je dán podle definice výrazem

$$\mu_2 = \frac{1}{\nu} \sum_1^{\nu} (\bar{x}_i - \bar{x})^2 \quad (4)$$

kde $\bar{x}_i = \frac{1}{r} \sum^{r,i} x_j$, takže označíme-li odchylky hodnot znaku od průměru v základním souboru

$$x_j - \bar{x} = \xi_j \quad (5)$$

bude

$$\bar{x}_i - \bar{x} = \frac{1}{r} (\sum^{r,i} x_j - r\bar{x}) = \frac{1}{r} \sum^{r,i} \xi_j \quad (6)$$

a čtverec

$$(\bar{x}_i - \bar{x})^2 = \frac{1}{r^2} \left(\sum_{j=1}^{r,i} \xi_j \right)^2 = \frac{1}{r^2} \left[\sum \xi_j^2 + 2 \sum_{j,k} \xi_j \xi_k \right], \quad j \neq k.$$

Součet $\sum_{j,k}^{r,i}$ značí součet všech členů utvořených tak, že bereme součiny všech proměnných ξ v i -tém výběru po dvou.

Dosadíme-li do (4), dostaneme

$$\mu_2 = \frac{1}{\nu} \cdot \frac{1}{r^2} \left[\frac{r}{N} \nu \sum_{j=1}^N \xi_j^2 + 2 \frac{\binom{r}{2} \nu}{\binom{N}{2}} \sum_{j,k=1}^N \xi_j \xi_k \right], \quad j \neq k \quad (7)$$

neboť sčítáme čtverce odchylek všech hodnot znaku ve všech ν výběrech, takže jako v předchozím případě průměru také zde každý z N čtverců hodnot odchylek znaku od průměru se bude vyskytovat $\frac{r\nu}{N}$ krát. Podobně uvažujeme v případě sčítání druhého členu hranaté závorky. Součet $\sum_{j,k}^{r,i}$

obsahuje $\binom{r}{2}$ členů, takže v součtu pro všech ν výběrů je $\binom{r}{2} \nu$ sčítanců. Celkem je $\binom{N}{2}$ různých podvojných součinů N hodnot znaku, takže každý určitý podvojný součin $\xi_j \xi_k$ se bude vyskytovat $\binom{r}{2} \nu : \binom{N}{2}$ krát.

Pro součet odchylek od průměru a jejich čtverců můžeme psát rovnice

$$\sum_{j=1}^N \xi_j = 0, \quad \sum_{j=1}^N \xi_j^2 = N \mu(x, 2); \quad (8)$$

čtverec první z těchto rovnic bude

$$\left(\sum_{j=1}^N \xi_j \right)^2 = \sum_{j=1}^N \xi_j^2 + 2 \sum_{j,k=1}^N \xi_j \xi_k = 0 \quad (j \neq k),$$

takže

$$2 \sum_{j,k=1}^N \xi_j \xi_k = -N \mu(x, 2) \quad (9)$$

Budeme-li dále označovati

$$\vartheta_i = \frac{\binom{r}{i}}{\binom{N}{i}} = \frac{r(r-1)(r-2)\dots(r-i+1)}{N(N-1)(N-2)\dots(N-i+1)}, \quad (10)$$

dostaneme rovnici pro druhý moment vzhledem k průměru ve tvaru

$$\mu_2 = \frac{1}{r^2} N \cdot \mu(x, 2) [\vartheta_1 - \vartheta_2] \quad (11)$$

nebo

$$\mu_2 = \frac{1}{r} \mu(x, 2) \left(1 - \frac{r-1}{N-1} \right). \quad (12)$$

Pro směrodatnou odchylku výběrových průměrů σ_P pak najdeme odmocněním

$$\sigma_P = \sigma(x) \sqrt{\frac{N-r}{r(N-1)}} = \sigma(x) \sqrt{\frac{1}{r} \frac{N-1}{N-r}}. \quad (13)$$

Z toho výsledku je patrna druhá věta:

Rozptyl a tudíž také směrodatná odchylka výběrových průměrů je menší než rozptyl či směrodatná odchylka základního souboru.

Je-li N velké, takže můžeme psáti N místo $N-1$, přejde výraz (13) ve tvar

$$\sigma_P = \sigma(x) \sqrt{\frac{N-r}{rN}} = \sigma(x) \sqrt{\frac{1}{r} - \frac{1}{N}}.$$

Šikmost. Abychom našli šikmost rozdělení četnosti výběrových průměrů, odvodíme nyní jejich třetí moment

vzhledem k průměru podobně jako v předcházejícím odstavci. Podle definice

$$\mu_3 = \frac{1}{\nu} \sum_1^{\nu} (\bar{x}_i - \bar{x})^3,$$

což vzhledem ku (5) a (6) dává pomocí multinomického teorému

$$\mu_3 = \frac{1}{\nu} \sum_1^{\nu} \frac{1}{r^3} \left[\sum_1^{r,i} \xi_j^3 + 3 \sum_1^{r,i} \xi_j^2 \xi_k + 6 \sum_1^{r,i} \xi_j \xi_k \xi_l \right]. \quad (14)$$

Multinomický teorém je zobecněním binomického teorému, podle něhož je

$$(\xi_1 + \xi_2)^n = \sum \frac{n!}{a! b!} \xi_1^a \xi_2^b,$$

kde $n = a + b$ a součet se vztahuje na všechna možná rozložení čísla n ve dva sčítance, z nichž každý je číslo celé nezáporné. Rozšíří-li se počet členů v závorce, lze dokázat, že platí

$$(\xi_1 + \xi_2 + \dots + \xi_N)^n = \sum \frac{n!}{a! b! c! \dots} \xi_1^a \xi_2^b \xi_3^c \dots$$

Číslo n se nyní rozloží zase všemi možnými způsoby na N celých nezáporných sčítanců $n = a + b + c + \dots$ a součet se vztahuje na všechna tato rozložení čísla n .

Zcela obdobnými úvahami jako při odvozování rovnice (7) provedeme součty, takže dostaneme

$$\mu_3 = \frac{1}{r^3} \left[\partial_1 \sum_{j=1}^N \xi_j^3 + 3\partial_2 \sum_{j,k=1}^N \xi_j^2 \xi_k + 6\partial_3 \sum_{j,k,l=1}^N \xi_j \xi_k \xi_l \right] \quad (15)$$

a tyto součty můžeme vyjádřit pomocí momentů základního souboru na základě vztahů plynoucích ze sčítání třetího řádu

$$\sum_{j=1}^N \xi_j^3 = N\mu(x, 3), \quad \sum_{j=1}^N \xi_j^2 \sum_{k=1}^N \xi_k = 0 = \sum_{j=1}^N \xi_j^3 + \sum_{j,k} \xi_j^2 \xi_k,$$

$$\sum_{j,k=1}^N \xi_j \xi_k \sum_{l=1}^N \xi_l = 0 = \sum_{j,k=1}^N \xi_j^2 \xi_k + 3 \sum_{j,k,l=1}^N \xi_j \xi_k \xi_l,$$

čili

$$\sum_{j,k=1}^N \xi_j^2 \xi_k = -N\mu(x, 3), \quad 3 \sum_{j,k,l=1}^N \xi_j \xi_k \xi_l = N\mu(x, 3),$$

takže po dosazení do (15)

$$\mu_3 = \frac{1}{r^3} N\mu(x, 3) [\vartheta_1 - 3\vartheta_2 + 2\vartheta_3] \quad (16)$$

a rozvedeme-li výrazy v závorce podle (10), můžeme také po jednoduché úpravě psáti

$$\mu_3 = \frac{1}{r^2} \mu(x, 3) \frac{(N-r)(N-2r)}{(N-1)(N-2)}. \quad (17)$$

Pro šikmost

$$\bar{\varrho} = \frac{\mu_3}{\mu_2^{\frac{3}{2}}} = \frac{\mu_3}{\mu_2 \sigma_P},$$

tudíž plyne po dosazení z (17), (12), (13) výraz

$$\bar{\varrho} = \varrho(x) \frac{N-2r}{N-2} \sqrt{\frac{N-1}{r(N-r)}}. \quad (18)$$

Šikmost výběrových průměrů je menší než šikmost základního souboru.

Exces. Stejným postupem vypočítáme čtvrtý moment vzhledem k průměru \bar{x} , abychom mohli vyjádřit exces rozdělení četností výběrových průměrů. Podle multinomického teorému bude

$$\begin{aligned} \mu_4 = \frac{1}{r} \sum_{i=1}^r \frac{1}{r^4} & \left[\sum_{j=1}^{r,i} \xi_j^4 + 4 \sum_{j=1}^{r,i} \xi_j^3 \xi_k + 6 \sum_{j=1}^{r,i} \xi_j^2 \xi_k^2 + \right. \\ & \left. + 12 \sum_{j=1}^{r,i} \xi_j^2 \xi_k \xi_l + 24 \sum_{j=1}^{r,i} \xi_j \xi_k \xi_l \xi_m \right] \end{aligned}$$

čili

$$\mu_4 = \frac{1}{r^4} \left[\vartheta_1 \sum_{j=1}^N \xi_j^4 + 4\vartheta_2 \sum_{j,k=1}^N \xi_j^3 \xi_k + 6\vartheta_2 \sum_{j,k=1}^N \xi_j^2 \xi_k^2 + \right. \\ \left. + 12\vartheta_3 \sum_{j,k,l=1}^N \xi_j^2 \xi_k \xi_l + 24\vartheta_4 \sum_{j,k,l,m=1}^N \xi_j \xi_k \xi_l \xi_m \right]. \quad (19)$$

Použijeme nyní vztahů

$$\sum_{j=1}^N \xi_j^4 = N\mu(x, 4), \quad \sum_{j=1}^N \xi_j^3 \sum_{k=1}^N \xi_k = 0 = \sum_{j=1}^N \xi_j^4 + \sum_{j,k=1}^N \xi_j^3 \xi_k, \\ \left(\sum_{j=1}^N \xi_j^2 \right)^2 = \sum_{j=1}^N \xi_j^4 + 2 \sum_{j,k=1}^N \xi_j^2 \xi_k^2,$$

$$\sum_{j=1}^N \xi_j^2 \left(\sum_{k=1}^N \xi_k \right)^2 = 0 = \sum_{j=1}^N \xi_j^2 \left(\sum_{k=1}^N \xi_k^2 + 2 \sum_{k,l=1}^N \xi_k \xi_l \right) = \\ = \left(\sum_{j=1}^N \xi_j^2 \right)^2 + 2 \sum_{j,k=1}^N \xi_j^3 \xi_k + 2 \sum_{j,k,l=1}^N \xi_j^2 \xi_k \xi_l,$$

$$\left(\sum_{j=1}^N \xi_j \right)^4 = 0 = \sum_{j=1}^N \xi_j^4 + 4 \sum_{j,k=1}^N \xi_j^3 \xi_k + 6 \sum_{j,k=1}^N \xi_j^2 \xi_k^2 + \\ + 12 \sum_{j,k,l=1}^N \xi_j^2 \xi_k \xi_l + 24 \sum_{j,k,l,m=1}^N \xi_j \xi_k \xi_l \xi_m,$$

z nichž plyne

$$\sum_{j,k=1}^N \xi_j^3 \xi_k = -N\mu(x, 4), \quad 2 \sum_{j,k=1}^N \xi_j^2 \xi_k^2 = N^2\mu^2(x, 2) - N\mu(x, 4), \\ 2 \sum_{j,k,l=1}^N \xi_j^2 \xi_k \xi_l = 2N\mu(x, 4) - N^2\mu^2(x, 2), \\ 24 \sum_{j,k,l,m=1}^N \xi_j \xi_k \xi_l \xi_m = -6N\mu(x, 4) + 3N^2\mu^2(x, 2),$$

takže dosazením do (19) dostáváme

$$\mu_4 = \frac{1}{r^4} N\mu(x, 4) [\vartheta_1 - 7\vartheta_2 + 12\vartheta_3 - 6\vartheta_4] + \\ + 3N^2\mu^2(x, 2) [\vartheta_2 - 2\vartheta_3 + \vartheta_4]. \quad (20)$$

Pro pátý moment výběrových průměrů vzhledem k jejich průměru bychom dostali

$$\mu_5 = \frac{1}{r^5} N\mu(x, 5) [\vartheta_1 - 15\vartheta_2 + 50\vartheta_3 - 60\vartheta_4 + 24\vartheta_5] + \\ + 10N^2\mu(x, 2)\mu(x, 3) [\vartheta_2 - 4\vartheta_3 + 5\vartheta_4 - 2\vartheta_5]. \quad (21)$$

Vidíme, že výpočet vyšších momentů se stává stále složitějším, takže je pak nutno k zjednodušení výpočtů užívatí účelné symboliky, ale přes to i výsledky jsou pro praktické upotřebení velmi složité. Podstatné zjednodušení nastává pro případ, kdy základní soubor je nekonečného rozsahu.

(2,2) Momenty výběrových průměrů z nekonečného základního souboru. Necháme-li v rovnicích (3), (12), (13), (17), (20), (21), růsti $N \rightarrow \infty$ do nekonečna, a při tom zůstává r konečné, dostaneme pro výběrové momenty výrazy

$$\mu'_1 = \bar{x}, \quad \mu_2 = \frac{1}{r} \mu(x, 2), \quad (22)$$

$$\mu_3 = \frac{1}{r^2} \mu(x, 3), \quad \mu_4 = \frac{1}{r^3} [\mu(x, 4) + 3(r-1)\mu^2(x, 2)], \quad (23)$$

$$\mu_5 = \frac{1}{r^4} [\mu(x, 5) + 10(r-1)\mu(x, 3)\mu(x, 2)], \dots$$

odkud dostáváme pro směrodatnou odchylku výběrových průměrů

$$\sigma_P = \frac{\sigma(x)}{\sqrt{r}} \quad (24)$$

a rovnice (22) a (23) můžeme snadno uvést na tvar

$$\mu_2 = \frac{1}{r} \mu(x, 2),$$

$$\mu_3 = \frac{1}{r^2} \mu(x, 3), \quad (25)$$

$$\mu_4 - 3\mu_2^2 = \frac{1}{r^3} [\mu(x, 4) - 3\mu^2(x, 2)],$$

$$\mu_5 - 10\mu_3\mu_2 = \frac{1}{r^4} [\mu(x, 5) - 10\mu(x, 3)\mu(x, 2)],$$

odkud plynou pro charakteristiky směrodatné proměnné (I, str. 22 a 23) výrazy

$$\begin{aligned} \alpha_{P,2} &= 1, \\ \alpha_{P,3} &= \frac{1}{r^{\frac{1}{2}}} \alpha(x, 3), \end{aligned} \quad (26)$$

$$\alpha_{P,4} - 3 = \frac{1}{r} [\alpha(x, 4) - 3],$$

$$\alpha_{P,5} - 10\alpha_{P,3} = \frac{1}{r^{\frac{3}{2}}} [\alpha(x, 5) - 10\alpha(x, 3)],$$

.....

Budeme-li nyní psáti v (25)

$$\begin{aligned} \lambda_2 &= \mu_2, & \lambda(x, 2) &= \mu(x, 2), \\ \lambda_3 &= \mu_3, & \lambda(x, 3) &= \mu(x, 3), \\ \lambda_4 &= \mu_4 - 3\mu_2^2, & \lambda(x, 4) &= \mu(x, 4) - 3\mu^2(x, 2), \\ \lambda_5 &= \mu_5 - 10\mu_3\mu_2, & \lambda(x, 5) &= \mu(x, 5) - 10\mu(x, 3)\mu(x, 2), \end{aligned} \quad (27)$$

vidíme, že podle rovnic (25) je rozdělení výběrových průměrů z nekonečného základního souboru charakterisováno jednoduchým vztahem mezi λ -funkcemi

$$\lambda_n = \frac{1}{r^{n-1}} \lambda(x, n). \quad (28)$$

Jinou cestou objevil Thiele důležitost těchto λ -funkcí, které tolik přispěly k rozvoji teorie matematické statistiky a nazývají se Thieleho semiinvarianty.

Dále přicházíme k semiinvariantům směrodatné proměnné čili standardisovaným semiinvariantům Thieleho, píšeme-li v (26)

$$\begin{aligned} \gamma_3 &= \alpha_{P,3}, & \gamma(x, 3) &= \alpha(x, 3), \\ \gamma_4 &= \alpha_{P,4} - 3, & \gamma(x, 4) &= \alpha(x, 4) - 3, \\ \gamma_5 &= \alpha_{P,5} - 10\alpha_{P,3}, & \gamma(x, 5) &= \alpha(x, 5) - 10\alpha(x, 3), \end{aligned} \quad (29)$$

mezi nimiž platí vztah

$$\gamma_n = \frac{1}{r^{\frac{n}{2}-1}} \gamma(n, x). \quad (30)$$

Necháme-li zde růsti rozsah výběru $r \rightarrow \infty$ do nekonečna, bude $\lim_{r \rightarrow \infty} \gamma_n = 0$, z čehož podle (26) plyne

$$\alpha_{P,3} = 0, \quad \alpha_{P,4} = 3, \quad \alpha_{P,5} = 0, \dots$$

a dále bychom dostali obecně

$$\alpha_{P,2n} = \frac{(2n)!}{2^n(n!)}, \quad \alpha_{P,2n+1} = 0,$$

což jsou momenty normálního rozdělení četností [I, (47), (48)].

Z toho vidíme, že pro velká r jsou momenty tohoto rozdělení četností shodné s momenty normálního rozdělení. To znamená, že bereme-li velké výběry z nekonečného základního souboru, můžeme očekávat, že rozdělení výběrových průměrů se těsně přiblíží normálnímu.

(2,2,1) Příklad. 1. Stanovme momenty výběrových průměrů ze základního souboru nekonečného, jedná-li se o alternativní znak s četností p . Označíme pozorovaný znak jedničkou a jeho nepřítomnost nulou, takže hodnota znaku $x_1 = 1$ má relativní četnost p , $x_2 = 0$ má relativní četnost $q = 1 - p$. Průměr $\bar{x} = p$. Bude tudíž n -tý moment vzhledem k průměru

$$\mu(x, n) = \sum (x_i - \bar{x})^n f_i$$

vyjádřen

$$\begin{aligned} \mu(x, n) &= (1-p)^n p + (-p)^n (1-p) = \\ &= pq [q^{n-1} + (-1)^n p^{n-1}]. \end{aligned}$$

Výběry rozsahu r z tohoto základního souboru mají vzhledem k počtu prvků s pozorovaným znakem rozdělení binomické $\binom{r}{x} p^x q^{r-x}$.

Momenty výběrových průměrů pro proměnnou x pak stanovíme podle rovnic (22) a (23); dostáváme

$$\begin{aligned}\mu'_1 &= p, & \mu_2 &= \frac{1}{r} pq, \\ \mu_3 &= \frac{1}{r^2} pq (q^2 - p^2), \\ \mu_4 &= \frac{1}{r^3} pq (q^3 + p^3) + \frac{3}{r^3} p^2 q^2 (r - 1).\end{aligned}$$

Příklad 2. Vypočítejme momenty výběrových průměrů ze základního souboru:

a) S rozdělením normálním.

Vzhledem k rovnicím [I, (48), str. 80] můžeme podle rovnic (22) a (23) psáti

$$\begin{aligned}\mu'_1 &= \bar{x}, & \mu_2 &= \frac{1}{r} \sigma^2(x), & \mu_3 &= 0, \\ \mu_4 &= \frac{1}{r^3} [3\sigma^4(x) + 3(r-1)\sigma^4(x)] = \frac{3}{r^2} \sigma^4(x).\end{aligned}$$

Z rovnic (27) bychom se přesvědčili, že všechny semiinvarianty vyššího stupně než druhého jsou identicky rovny nule.

b) S rozdělením podle Pearsonovy křivky typu III.

Typ III ze systému křivek Pearsonových (Janko [1], str. 42), lze pomocí momentů směrodatné proměnné uvést na tvar

$$y = y_0 \left(1 + \frac{\alpha(x, 3)}{2} t \right)^{\frac{4}{\alpha^2(x, 3)} - 1} e^{-\frac{2}{\alpha(x, 3)} t}.$$

Počátek souřadnic je v průměru a první momenty směrodatné proměnné vzhledem k němu jsou $\alpha(x, 0) = 1$, $\alpha(x, 1) = 0$, $\alpha(x, 2) = 1$.

Pro další momenty pak je odvozena rekurentní formule

$$\alpha(x, n+1) = n \left[\alpha(x, n-1) + \frac{\alpha(x, 3) \alpha(x, n)}{2} \right]. \quad (31)$$

Z toho tedy plyne vyjádření dalších momentů pomocí $\alpha(x, 3)$, takže

$$\begin{aligned} \alpha(x, 4) &= 3 \left(1 + \frac{\alpha^2(x, 3)}{2} \right), \\ \alpha(x, 5) &= 2\alpha(x, 3) \left[5 + \frac{3\alpha^2(x, 3)}{2} \right], \dots \end{aligned}$$

a semiinvarianty směrodatné proměnné jsou podle (29)

$$\begin{aligned} \gamma(x, 4) &= \frac{3\alpha^2(x, 3)}{2} = \left(\frac{\alpha(x, 3)}{2} \right)^2 3!, \\ \gamma(x, 5) &= \frac{2 \cdot 3\alpha^3(x, 3)}{2} = \left(\frac{\alpha(x, 3)}{2} \right)^3 4!. \end{aligned} \quad (32)$$

Semiinvarianty γ_4 a γ_5 rozdělení průměrů jsou pak dány podle rovnice (30).

Vypočítáme však momenty směrodatné proměnné podle rovnic (26) za použití rekurentního vztahu (31). Tak dostaneme

$$\begin{aligned} \alpha_{P,2} &= 1, & \alpha_{P,3} &= \frac{1}{r^{\frac{1}{2}}} \alpha(x, 3), & \alpha_{P,4} &= \frac{3}{r} \left[r + \frac{\alpha^2(x, 3)}{2} \right], \\ \alpha_{P,5} &= \frac{2}{r^{\frac{3}{2}}} \alpha(x, 3) \left[5r + \frac{3}{2} \alpha^2(x, 3) \right], \end{aligned}$$

jako momenty rozdělení průměrů, jestliže se berou výběry rozsahu r z nekonečného základního souboru Pearsonova typu III. Pro čtvrtý a pátý moment se můžeme přesvědčiti, že vyhovují podmínce

$$\alpha_{P,n+1} = n(\alpha_{P,n-1} + \frac{\alpha_{P,3}}{2} \alpha_{P,n});$$

poněvadž i další momenty této podmínce vyhovují, vidíme, že rozdělení průměrů je tu dáno také Pearsonovým typem III.

(2,3) Momenty rozdělení četností výběrových rozptylů. Rozdělení četností průměrů bylo známo již v minulém století. Je však jednou z nejvýznamnějších složek pokroku tohoto století ve statistice, že se začalo studovati rozdělení četností rozptylů.

Průměr. Než přistoupíme k výpočtu těchto charakteristik, je dobře zvláště si uvědomiti, že se jedná o druhé momenty kolem průměru každého jednotlivého výběru. Tak rozptyl výběru čili druhý moment kolem průměru tohoto výběru je podle definice

$$\sigma_{x,i}^2 = \frac{1}{r} \sum^{r,i} (x_j - \bar{x}_i)^2, \quad (33)$$

kde se součet vztahuje na všech r prvků i -tého výběru.

Tento výraz upravíme, abychom do něho zavedli jen hodnoty znaku

$$\sum^{r,i} (x_j - \bar{x}_i)^2 = \sum^{r,i} x_j^2 - 2\bar{x}_i \sum^{r,i} x_j + r\bar{x}_i^2 = \sum^{r,i} x_j^2 - r\bar{x}_i^2,$$

kde jsme psali $r\bar{x}_i = \sum^{r,i} x_j$, takže bude dále

$$\sum^{r,i} x_j^2 - \frac{1}{r} \left(\sum^{r,i} x_j \right)^2 = \frac{1}{r} \left[r \sum^{r,i} x_j^2 - \sum^{r,i} x_j^2 - 2 \sum_{j,k}^{r,i} x_j x_k \right]$$

a tedy

$$\sigma_{x,i}^2 = \frac{1}{r^2} \left[(r-1) \sum^{r,i} x_j^2 - 2 \sum_{j,k}^{r,i} x_j x_k \right].$$

Sečteme-li všech ν výběrových rozptylů a dělíme jejich počtem, dostaneme průměr $\bar{\sigma}^2$ podobnými úvahami jako při

odvození rovnice (7)

$$\bar{\sigma}^2 = \frac{1}{\nu} \sum_{i=1}^{\nu} \sigma_{x,i}^2 = \frac{1}{\nu} \cdot \frac{1}{r^2} \left[\frac{r\nu}{N} (r-1) \sum_{j=1}^N x_j^2 - 2 \frac{\binom{r}{2}^{\nu}}{\binom{N}{2}} \sum_{j,k} x_j x_k \right],$$

takže vzhledem k rovnicím (7) a (9) bude

$$\bar{\sigma}^2 = \frac{1}{r^2} N \mu(x, 2) [(r-1) \vartheta_1 + \vartheta_2] = \frac{1}{r^2} N^2 \vartheta_2 \mu(x, 2), \quad (34)$$

což lze psát

$$\bar{\sigma}^2 = \frac{r-1}{r} \frac{N}{N-1} \mu(x, 2). \quad (35)$$

Průměrný rozptyl výběrový je menší než rozptyl základního souboru.

Z rovnic (12) a (35) je zřejmo, že platí vztah

$$\mu_2 + \bar{\sigma}^2 = \mu(x, 2). \quad (36)$$

Rozptyl rozdělení četností výběrových rozptylů. Pro zjednodušení výpočtu druhého momentu $\mu(\sigma_{x,i}^2, 2)$ výběrových rozptylů kolem jejich průměru (35) by bylo třeba užití vhodného symbolického počtu, jak jsme se již dříve zmínili. Pro náš účel se však zde spokojíme sdělením výsledku

$$\begin{aligned} \mu(\sigma_{x,i}^2, 2) &= \frac{1}{r^4} N \mu(x, 4) [(r-1)^2 \vartheta_1 - (r^2 - 6r + 7) \vartheta_2 - \\ &- 4(r-3) \vartheta_3 - 6\vartheta_4] + \frac{1}{r^4} N^2 \mu^2(x, 2) [(r^2 - 2r + 3) \vartheta_2 + \\ &+ 2(r-3) \vartheta_3 + 3\vartheta_4 - N^2 \vartheta_2^2] \end{aligned}$$

čili

$$\begin{aligned} \mu(\sigma_{x,i}^2, 2) &= \frac{N(r-1)(N-r)}{r^3(N-1)^2(N-2)(N-3)} \\ &\{ (N-1)(rN - N - r - 1) \mu(x, 4) + \\ &+ [(3-r)N^2 - 6N + 3r + 3] \mu^2(x, 2) \}. \quad (37) \end{aligned}$$

Druhou odmocninou tohoto výrazu je pak dána směrodatná odchylka výběrových rozptylů $\sigma(\sigma_{x,i^2})$.

Výrazy pro další momenty jsou ještě mnohem rozsáhlejší a nebudeme je uvádět, ježto také praktická cena jejich je tím omezena.

(2,4) Momenty výběrových rozptylů z nekonečného základního souboru. Průměr výběrových rozptylů je dán rovnicí (35), v níž necháme nyní růsti rozsah $N \rightarrow \infty$ do nekonečna, takže potom

$$\bar{\sigma}^2 = \frac{r-1}{r} \mu(x, 2). \quad (38)$$

Rozptyl jejich dostaneme obdobně z rovnice (37), z níž pro $N \rightarrow \infty$ vyplývá

$$\begin{aligned} \mu(\sigma_{x,i^2}, 2) &= \frac{r-1}{r^3} [(r-1)\mu(x, 4) - (r-3)\mu^2(x, 2)] = \\ &= \frac{r-1}{r^3} \sigma^4(x) [(r-1)\alpha(x, 4) - (r-3)] \end{aligned} \quad (39)$$

a směrodatnou odchylku můžeme tedy psát ve tvaru

$$\sigma(\sigma_{x,i^2}) = \frac{\sigma^2(x)}{r} \sqrt{\frac{r-1}{r} [(r-1)\alpha(x, 4) - r + 3]}. \quad (40)$$

Uvedeme ještě třetí moment výběrových rozptylů $\mu(\sigma_{x,i^2}, 3)$ z nekonečného základního souboru.

$$\begin{aligned} \mu(\sigma_{x,i^2}, 3) &= \frac{r-1}{r^5} \sigma^6(x) [(r-1)^2 \alpha(x, 6) - \\ &- 3(r-1)(r-5)\alpha(x, 4) - 2(3r^2 - 6r + 5)\alpha^2(x, 3) + \\ &+ 2(r^2 - 12r + 15)]. \end{aligned}$$

Velmi značné zjednodušení nastává pro případ základního souboru s normálním rozdělením četností. Výpočtem semi-invariantů lze prokázat, že rozdělení výběrových rozptylů kolem průměru je Pearsonova křivka typu III.

Charakteristiky rozdělení četností třetích a čtvrtých výběrových momentů.

Uvedeme jen výsledky pro průměr a rozptyl třetích výběrových momentů kolem průměru z nekonečného základního souboru $N \rightarrow \infty$.

Průměr třetích výběrových momentů

$$\begin{aligned}\mu'(\mu_{3,i}, 1) &= \frac{(r-1)(r-2)}{r^2} \mu(x, 3) = \\ &= \frac{r-1}{r^2} (r-2) \sigma^3(x) \alpha(x, 3)\end{aligned}$$

rozptyl

$$\begin{aligned}\sigma^2(\mu_{3,i}) &= \frac{(r-1)(r-2)}{r^5} \sigma^6(x) [(r-1)(r-2) \alpha(x, 6) - \\ &- 3(r-2)(2r-5) \alpha(x, 4) - (r-2)(r-10) \alpha^2(x, 3) + \\ &+ 3(3r^2 - 12r + 20)].\end{aligned}$$

Průměr rozdělení čtvrtých momentů

$$\begin{aligned}\mu'(\mu_{4,i}, 1) &= \frac{r-1}{r^3} [(r^2 - 3r + 3) \mu(x, 4) + \\ &+ 3(2r - 3) \mu^2(x, 2)].\end{aligned}$$

Odvozením charakteristik výběrových pomocí parametrů základního souboru jsme získali možnost řešiti v případech, kde známe rozdělení četností dotyčné charakteristiky otázku, jaká je pravděpodobnost, že na př. průměr náhodného výběru bude v daných mezích. Tak víme, že bude v intervalu $\pm 2\sigma_P$ s pravděpodobností 0,955, nebo v intervalu $\pm 3\sigma_P$ s pravděpodobností 0,997, když se jedná o velký výběr z nekonečného základního souboru, t. j. výběr, který byl vzat tak, že se každý prvek po zjištění příslušné hodnoty znaku vrátil zpět do základního souboru čili jeho složení zůstávalo nezměněno. V tom případě totiž můžeme považovat rozdělení četností průměrů za normální, jak jsme si odvodili. Obdobně můžeme postupovati, známe-li rozdělení

četností výběrových rozptylů, nebo můžeme-li o něm předpokládati, že je vyjádřeno na př. Pearsonovou křivkou typu III. Integrály této křivky jsou dány buď Pearsonovými tabulkami neúplné Γ -funkce nebo výhodněji pro směrodatnou proměnnou tabulkami Salvosovými.

Budeme potřebovati hlavně k praktickému použití směrodatné odchylky některých charakteristik, které zde ještě uvedeme, ale již bez odvození.

Směrodatná odchylka rozdělení směrodatných odchylek výběrových z normálního základního souboru je

$$\sigma_{\sigma} = \frac{\sigma(x)}{\sqrt{2r}} \quad (41)$$

a nesmí se jí ovšem užívat bez tohoto zřetele k formě rozdělení četností základního souboru. V obecném případě by totiž bylo nutno pracovati s výrazem

$$\sigma_{\sigma} = \frac{\sigma(x)}{\sqrt{2r}} \left(1 + \frac{\frac{\mu(x, 4)}{\mu^2(x, 2)} - 3}{2} \right)^{\frac{1}{2}}. \quad (42)$$

Je-li exces $\left[\frac{\mu(x, 4)}{\mu^2(x, 2)} - 3 \right]$ malý, pak odmocnina výrazu v kulaté závorce je přibližně rovna $1 + \frac{1}{4} \left[\frac{\mu(x, 4)}{\mu^2(x, 2)} - 3 \right]$ a toto číslo se liší od jednotky o více než 5 procent, je-li $\frac{\mu(x, 4)}{\mu^2(x, 2)}$ menší než 2,8 nebo větší než 3,2.

Dále uvedeme ještě směrodatnou odchylku rozdělení četností šikmosti ve výběrech z normálního základního souboru,

která je $\sqrt{\frac{6}{r}}$ a excesu, která je $2 \sqrt{\frac{24}{r}}$.

(2,4,1) Příklad 1. Základní soubor má rozsah 1000 prvků; průměr měřeného znaku je 140 jednotek. Jaký průměr mů-

žeme očekávat, vezmeme-li z něho náhodný výběr rozsahu $r = 500$, nebo $r = 100$. Vyložte, zda by k odpovědi pomohla znalost směrodatné odchylky rozdělení četnosti výběrových průměrů. Odpověď: můžeme očekávat průměr blízký hodnotě $\bar{x} = 140$ jednotek a to v prvním případě bližší než v druhém. Kdybychom znali směrodatnou odchylku σ_P , která by byla v každém z obou případů jiná, mohli bychom říci, že pravděpodobná odchylka od hodnoty $\bar{x} = 140$ bude $\pm 0,6745 \sigma_P$.

Příklad 2. Základní soubor je rozsahu $N = 1000$ prvků, průměr hodnot pozorovaného znaku je 67,6 cm a směrodatná odchylka 2,5 cm. Jaká je směrodatná odchylka všech možných výběrových průměrů, je-li rozsah každého výběru 200, resp. 500 nebo 800 prvků? Použijeme-li výrazu (13) dostaneme $\sigma_P = 0,158; 0,08; 0,04$. Znázorněte graficky výraz (13) pro svrchu uvedené N a $\sigma(x)$ a proveďte diskusi. Jak velký musí být rozsah výběrů r , aby σ_P bylo menší než $\frac{1}{2}\sigma(x)$ resp. $\frac{1}{10}\sigma(x)$?

$$\sigma_P = \frac{\sigma(x)}{\sqrt{r \frac{N-1}{N-r}}} < \frac{\sigma(x)}{2} \text{ znamená, že } r \frac{N-1}{N-r} > 4 \text{ čili}$$

$$r > \frac{4N}{N+3} \doteq 3,99. \text{ V druhém případě } r \frac{N-1}{N-3} > 100 \text{ čili}$$

$$r > \frac{100N}{N+99} \doteq 90,99.$$

Výraz pro σ_P dává pro některá r tyto hodnoty σ_P :

r	1	10	100	200	500	800	N
σ_P	$\sigma(x) = 2,5$	0,79	0,24	0,14	0,08	0,04	0

Z tohoto přehledu a jeho grafického znázornění je patrné, že σ_P je klesající funkcí r a v našem případě klesá od hodnoty směrodatné odchylky základního souboru $\sigma(x) = 2,5$ tak, že pro $r = 4$ má hodnotu menší než její polovina, pro $r = 91$

menší než desítina a dospěje k nule, je-li rozsah výběru totožný s rozsahem základního souboru.

Příklad 3. Známe z tab. 1 (str. 10) rozdělení četností základního souboru rozsahu $N = 1001$ podle pozorovaného znaku. Průměr je $\bar{x} = 7,000$, směrodatná odchylka $\sigma(x) = 2,002$. Šikmost nemusíme brát v úvahu, neboť je rovna nule. Jaká je pravděpodobnost, že ve výběru rozsahu $r = 200$ bude průměr větší než 7,124, nebo větší než 7,174, nebo 7,500?

$$\sigma_P = 2,002 \sqrt{\frac{801}{200 \cdot 1000}} = \frac{2,002}{15,8015} = 0,127,$$

$$t_1 = \frac{7,124 - 7,000}{0,127} = 0,976, \quad t_2 = \frac{0,174}{0,127} = 1,370,$$

$$t_3 = \frac{0,500}{0,127} = 3,937.$$

Vzhledem k povaze základního souboru můžeme použít tabulky integrálu Laplace-Gaussova ([1], str. 38), takže

$$p = 0,5 - \frac{\alpha(t)}{2}$$

a dostáváme

$$p_1 = 0,165$$

$$p_2 = 0,085$$

$$p_3 = 0,00004.$$

Ve skutečnosti bylo podle sloupce (7) tabulky 1. větších průměrů výběrových než 7,124 celkem 70 ze 400, tedy $\sum f = 0,175$; větších než 7,174 bylo 30, takže $\sum f = 0,075$ a větších než 7,500 se nevyskytl.

Pozorovaná směrodatná odchylka byla $\sigma'_P = 0,121$.

Příklad 4. Základní soubor rozsahu $N = 10\,000$ prvků má průměr $\bar{x} = 69,7$ jednotek a $\sigma(x) = 7,4$. Jaký rozsah musí mít výběr, aby směrodatná odchylka průměrů byla menší než 2, 1, 0,5 jednotky?

Obecně má býti splněna nerovnost

$$\sigma_P = \frac{\sigma(x)}{\sqrt{r \frac{N-1}{N-r}}} < k,$$

odkud $\sigma^2(x) (N-r) < k^2 r (N-1)$,

$$r > \frac{\sigma^2(x) N}{k^2(N-1) + \sigma^2(x)}.$$

Pro $k = 2$ dostáváme $r > \frac{7,4^2 \cdot 10\,000}{4,9999 + 7,4^2} = 13,67$ a podobně pro $k = 1$ je $r > 54,47$ a pro $k = 0,5$ je $r > 214,37$. Poněvadž rozsah výběru je číslo celé, vidíme, že musí býti v prvním případě $r \geq 14$, v druhém případě $r > 54$ a ve třetím $r > 214$.

Příklad 5. Je-li v základním souboru nekonečného rozsahu šikmost 1,1, jak musí býti velký náhodný výběr, aby šikmost rozdělení všech možných výběrových průměrů byla menší než 0,1? Nakreslete křivku $\frac{1,1}{\sqrt{r}}$ a proveďte rozbor.

Podle rovnic (26) $\alpha_{P,3} = \frac{1}{r^{\frac{1}{2}}} \alpha(x, 3)$ a podle úlohy má býti

splněna podmínka $0,1 < \frac{\alpha(x, 3)}{\sqrt{r}}$, tedy $0,1 < \frac{1,1}{\sqrt{r}}$, takže

$r < \frac{1,21}{0,01} = 121$. Křivka s počátku rychle klesá od hodnoty

1,1 pro $r = 1$ a blíží se asymptoticky nule.

(2,5) Podstatná informace v parametrech a charakteristikách. Než přikročíme k výkladu inverzního úkolu odhadování parametrů základního souboru podle zjištěných charakteristik výběrových, osvětlíme si otázku, jak mnoho informace podávají jednotlivé charakteristiky o souboru, z něhož byly stanoveny. Budeme předpokládati k tomuto účelu, že

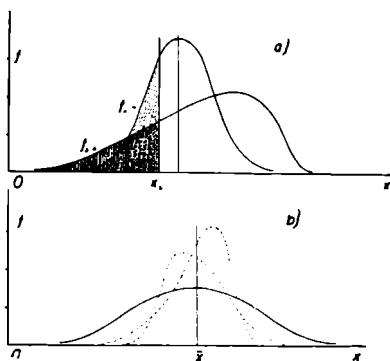
celá informace o určitém souboru vzhledem k pozorovanému znaku je obsažena v původní množině čísel x_1, x_2, \dots, x_N seřazených vzestupně podle velikosti (uspořádaně), tedy neseskupených do tříd (viz na př. I. díl, str. 19). Mnoho z této celé informace je obsaženo v několika málo charakteristikách. Máme tedy za úkol udati, v jakých charakteristikách musíme podati zhuštěně informaci, aby z ní vyplývalo těsné přiblížení k uspořádané posloupnosti čísel x_i , což znamená jinými slovy, aby bylo co možná s nejlepším přiblížením vyjádřeno procento prvků z celého počtu N , které spadají do určitého intervalu (x_i, x_{i+k}) hodnot znaku.

Zhuštěním informace o pozorovaném materiálu pomocí skupinového rozdělení četností jsme se zabývali v I. dílu (str. 29). Představme si, že uvedeme jen jednu hodnotu relativní četnosti p_k , která udává zlomek z celého počtu N pozorovaných hodnot x_i , které jsou menší než x_k .

Je zřejmo, že tím dáváme jen velmi malou část celé informace, neboť tuto hodnotu p_k mohou vykazovati rozdělení četností, která se od sebe zcela liší (viz v obr. 2 plochy f_1 a f_2).

Také velmi malou část celé informace podáváme, uvedeme-li jen průměr \bar{x} a rozsah pozorovaného souboru N . Touž hodnotu \bar{x} mohou totiž míti úplně odlišná rozdělení četností (obr. 2b).

Vidíme tudíž, že jedna charakteristika nemůže sama podati mnoho z celé informace obsažené v původní posloup-



Obr. 2. Množství informace obsažené jen a) v relativní četnosti, b) v průměru.

nosti. Teprve uvedením dvou nebo tří charakteristik můžeme dosáhnouti značně úplného vystižení rozdělení četností.

Dva parametry \bar{x} , $\sigma(x)$ nám udávají spolu s rozsahem N , že víc než $1 - \frac{1}{\tau^2}$ z celkového počtu r prvků souboru je v mezích $\bar{x} \pm \tau \sigma(x)$ (kde $\tau \geq 1$) (viz I. díl, str. 72) a to bez jakékoliv výhrady a bez ohledů na tvar rozdělení četností. Můžeme-li podle některých okolností, na př. podle původu materiálu, souditi na tvar rozdělení četností, dosáhneme výstižnějšího odhadu počtu prvků v uvedeném intervalu. Na př. pro normální rozdělení četností jsou dotyčné zlomky uvedeny v procentech v I. díle, str. 83.

Tak se vyskytuje normální rozdělení často při kontrole hromadné výroby nějakého předmětu, kde prvky pozorovaného souboru byly vyrobeny za týchž podstatných podmínek a rovněž pozorování znaku vyplynulo z měření za týchž podstatných podmínek. Pak se obyčejně říká, že data byla získána za kontrolovaných podmínek. V takových případech bývá možno pro většinu praktických účelů předpokládati, že rozdělení četností je normální nebo mírně nesymetrické.

Uvádíme-li kromě průměru a směrodatné odchylky ještě šikmost $\alpha(x, 3)$, přispíváme málo k vystižení rozdělení četností v symetrických intervalech kolem průměru a tedy k podání celé informace obsažené v pozorovaných datech. V takovém případě je třeba pomáhat si plochami křivky nesymetrické a v intervalech, jejichž hranice nejsou souměrně rozloženy nad a pod průměrem. Praktikové přikládají tomuto postupu význam tehdy, je-li rozsah pozorovaného souboru větší než $N = 250$.

Zkušenosti s rozděleními četností pro fyzikální vlastnosti materiálu a výrobků v továrnách vyráběných hromadně potvrzují, že se vystačí se způsoby právě uvedenými pro hrubé odhady pozorovaných procent nějakého rozdělení četností aspoň pro symetrické intervaly kolem průměru, ježto jsou to zpravidla rozdělení s jedním maximem. Není ovšem možno očekávati, že bychom s nimi vystačili v případě bi-

modálního rozdělení, které je výsledkem dvou různých množství podmínek, leč že bychom mohli rozštěpiti soubor ve dvě skupiny dat, z nichž každá by měla rozdělení jednovrcholové.

Záleží na účelu, k němuž hodláme použití pozorovaných dat, abychom mohli rozhodnouti, které charakteristiky mají býti uvedeny v konkrétním případě. Uvádíme je proto, abychom z nich mohli odvoditi hledané závěry nebo, aby jich mohl užítí statistický spotřebitel. Chceme tedy, aby obsahovaly podstatnou informaci. Co tvoří podstatnou informaci v určitém případě závisí na povaze otázek, které máme zodpověděti a na povaze hypotés, které chceme učiniti na základě informace z dat k tomu cíli dosažitelných. Říkáme, že nějaká skupina charakteristik obsahuje podstatnou informaci danou pozorovanými daty, když můžeme pomocí těchto charakteristik odpověděti na položené otázky tak, že by další rozbor dat naše odpovědi prakticky nepozměnil.

Praxe pak dospívá k názoru, že při studiu jednoho znaku, vzniklého za týchž podstatných podmínek na prvcích souboru, obsahuje průměr, směrodatná odchylka a rozsah souboru podstatnou informaci ve většině případů. Bývá tomu tak, když se zajímáme o průměrnou jakost materiálu a variabilitu průměrů postupných výběrů, nebo když srovnáváme tyto vlastnosti určitého materiálu s jiným. Kdybychom potřebovali k zodpovědění položených otázek znáti procenta těch prvků z celkového počtu pozorovaných prvků, kde hodnoty studovaného znaku jsou větší nebo menší než určitá daná hodnota, pak může býti podstatná informace obsažena v tabulce třídniho rozdělení četností.

Totéž, co zde bylo řečeno o parametrech, tedy pro základní soubor, platí o charakteristikách pro náhodný výběr.

(3) Náhodné výběry z neznámého základního souboru.

Zabývali jsme se úkolem: najíti pravděpodobnost, že charakteristika náhodného výběru bude v daných mezích, když známe parametry základního souboru.

Prakticky důležitější je ve statistice úkol obrácený: najítí pravděpodobnost, že parametry základního souboru se neliší od výběrových charakteristik určitého pozorovaného výběru o více než o daný zlomek jejich. Otázka může být také jinak položena. Je známa z pozorování charakteristika výběrová; máme udati meze, v nichž musí být hledaný parametr základního souboru s určitou napřed danou pravděpodobností.

Nebo konkrétně, ale méně přesně můžeme říci, že tato otázka znamená

a) do jaké míry se shoduje průměr výběru s průměrem základního souboru,

b) jak dobře se shoduje rozptyl a vyšší momenty výběru s příslušnými v základním souboru,

c) jak těsně blízko je křivka rozdělení četností výběru u křivky rozdělení četností základního souboru.

Vypočítáme-li z velkého výběru nějakou charakteristiku, třeba průměr \bar{x}_i , přisuzujeme jí obvykle velkou přesnost. Můžeme-li vzítí ze základního souboru nekonečného rozsahu více výběrů a pro každý vypočítati průměr, budou se zpravidla rozdíly mezi nimi zmenšovat, když rozsah výběrů poroste. Podle věty Bienaymé-Čebyševovy ([1], str. 70) usuzujeme, že rozdíl $|\bar{x}_i - \bar{x}|$ průměru pozorovaných výběrů a očekávaného průměru v základním souboru bude menší než libovolné kladné číslo $\eta = \tau\sigma_P$ s pravděpodobností

$$P = 1 - P_\tau > 1 - \frac{\sigma_P^2}{\eta^2}.$$

Poněvadž $\sigma_P^2 = \frac{1}{r} \sigma^2(x)$ — podle (22) — a rozptyl v základním souboru předpokládáme jako veličinu kladnou, blíží se $P > 1 - \frac{\sigma^2(x)}{r\eta^2}$ s rostoucím rozsahem výběru r k jednotce čili pravděpodobnost, že se průměr výběrový \bar{x}_i liší libovolně málo od průměru v základním souboru \bar{x} může se libovolně přiblížiti jednotce. Za těchto okolností říkáme, že

průměr pozorovaných dat (platí to také pro jiné charakteristiky výběrové) konverguje stochasticky k své očekávané hodnotě, t. j. k hodnotě v základním souboru (k příslušnému parametru).

(3,1) Charakteristiky konsistentní, efficientní, sufficientní. Můžeme to vyjádřit také tak, že nejvhodnějším odhadem parametru je charakteristika vypočítaná z náhodného výběru, jejíž očekávaná (průměrná) hodnota dává hledaný výsledek, t. j. hodnotu parametru. Takové charakteristiky se nazývají konsistentními čili souhlasnými odhady parametru.

Obecně existují různé charakteristiky, které mohou mít touž očekávanou hodnotu: na př. $\mathfrak{E}(x_i) = \bar{x}$ a také

$$\mathfrak{E}\left(\frac{1}{r} \sum_{i=1}^r x_i\right) = \bar{x}.$$

Zavádí se proto další kritérium, které dovoluje mezi několika konsistentními charakteristikami téhož parametru vybrati poměrně nejvhodnější. Za takovou se považuje ta charakteristika z konsistentních, která má normální rozdělení četností pro výběry velkého rozsahu r a poměrně nejmenší rozptyl: nazývá se efficientní čili vydatná nebo výstižná. Pro objasnění uvažujeme prvky $x_1, x_2, \dots, x_j, \dots, x_r$ výběru rozsahu r . Touž očekávanou hodnotu \bar{x} budou mít výrazy

$$x_1, \frac{x_1 + x_2}{2}, \dots, \frac{1}{r-1} \sum_{j=1}^{r-1} x_j, \frac{1}{r} \sum_{r=1}^r x_j$$

a jsou tedy konsistentní. Jejich rozptyly vzhledem k očekávané hodnotě však jsou

$$\mathfrak{E}(x_1 - \bar{x})^2 = \mu(x, 2),$$

$$\mathfrak{E}\left(\frac{x_1 + x_2}{2} - \bar{x}\right)^2 = \left(1 - \frac{1}{N-1}\right) \frac{\mu(x, 2)}{2},$$

.....

$$\mathbb{E} \left(\frac{1}{r-1} \sum_{j=1}^{r-1} x_j - \bar{x} \right)^2 = \left(1 - \frac{r-2}{N-1} \right) \frac{\mu(x, 2)}{r-1},$$

$$\mathbb{E} \left(\frac{1}{r} \sum_{i=1}^r x_i - \bar{x} \right)^2 = \left(1 - \frac{r-1}{N-1} \right) \frac{\mu(x, 2)}{r},$$

takže výběrový průměr $\frac{1}{r} \sum_{j=1}^r x_j$ považovaný za odhad parametru \bar{x} má poměrně nejmenší rozptyl. Vzhledem k tomu pak, že rozdělení četností spěje při rostoucím r k normálnímu, je efficientní čili vydatný. Vydatnost ostatních srovnávaných výrazů se měří obráceným poměrem jejich rozptylů k rozptylu nejvydatnějšího. Tak bude na př. vydatnost x_1 u srovnání s průměrem $\frac{1}{r} \sum_{j=1}^r x_j$ vyjádřena zlomkem

$$\left(1 - \frac{r-1}{N-1} \right) \frac{\mu(x, 2)}{r} : \mu(x, 2) = \frac{N-r}{(N-1)r}.$$

Lze také říci, že v tomto smyslu obsahuje výběrový průměr myslitelně nejúplnější informaci o parametru základního souboru u srovnání s ostatními výrazy, které mohou být z dat výběru rozsahu r počítány.

Nemůže tudíž výpočet těchto výrazů jako $\frac{1}{r-1} \sum_{j=1}^{r-1} x_j, \dots$ přispěti ničím novým k informaci podávané průměrem. Charakteristiky tohoto druhu mají také své pojmenování jako sufficientní čili vyčerpávající.

Je zřejmo, že můžeme vydatnost charakteristik vyjadřovat také v procentech té nejvydatnější. Tak můžeme porovnáním rozptylů průměru a mediánu zjistiti ([1], str. 184), že vydatnost mediánu klesá při rostoucím rozsahu výběru z normálního základního souboru na 80% při $r = 4$ a pak dále asymptoticky na 63%. Z toho na př. vyplývá, že medián obsahuje jen asi 63% té informace, kterou podává průměr a to z výběru rozsahu již asi $r = 25$. Objeví nám tedy průměr

z výběru o rozsahu $r=63$ změnu v poloze základního souboru stejně dobře jako medián z výběru rozsahu teprve $r = 100$. Podobně můžeme srovnati dvě charakteristiky rozptylu a to se směrodatnou odchylkou σ průměrnou odchylku ϑ . Vydatnost směrodatné odchylky je v případě normálního rozdělení četností o 12% vyšší než průměrné odchylky. Je tudíž výhodnější vynaložiti trochu více práce výpočtu výrazu $\frac{N-1}{N} \cdot \frac{r}{r-1} \bar{\sigma}^2$ — vzhledem k rovnici (35) — než zvět-

šiti počet pozorování asi o 14%, což by bylo nutné, kdybychom chtěli dosáhnouti pomocí průměrné odchylky té přesnosti jako při směrodatné odchylce, které se proto užívá téměř výhradně pro přesnější výpočty statistické.

(3,2) Odhad průměru. V praxi statistické musíme zpravidla udati určitou hodnotu charakteristiky vypočítanou z náhodného výběru, kterou lze považovati za nejvhodnější odhad velikosti neznámého parametru. Jedná-li se konkrétně o průměr, uvědomíme si, že průměr všech možných výběrových průměrů byl týž jako průměr základního souboru; to jsme viděli v případě, kdy základní soubor byl znám (str. 14). Také jsme si vyložili (str. 11), že průměr výběru velkého rozsahu se mnoho neliší od průměru základního souboru. Ježto v praxi nebereme nikdy všechny možné výběry, nýbrž jeden nebo dva výběry dostačujícího rozsahu, užijeme jednoho zjištěného průměru nebo průměru dvou výběrů jako odhadu průměru základního souboru. Tento výběrový průměr bude representovati tedy průměr neznámého základního souboru a bude považován za jemu blízký, ježto variační obor rozdělení četností všech možných výběrových průměrů je velmi malý [viz formuli (13) nebo příklad (2,4,1,3), str. 31]. Proto se přijímá výběrový průměr za dobrý odhad průměru základního souboru.

(3,3) Odhad směrodatné odchylky. Potřebujeme nyní směrodatnou odchylku všech možných výběrových průměrů čili směrodatnou odchylku jejich rozdělení četností, když zá-

kladní soubor neznáme. Tato směrodatná odchylka nám bude naznačovat jak je přesný průměr, který jsme dostali z jednoho nebo několika málo výběrů. V případě, kdy základní soubor byl znám, jsme našli směrodatnou odchylku výběrových průměrů (13) a vyjádřili ji pomocí známé směrodatné odchylky základního souboru.

Tohoto výrazu však nemůžeme užití, když základní soubor neznáme a tedy směrodatná odchylka jeho také není známa. Stojíme tudíž před problémem odhadu tohoto parametru. Vezmeme tedy ku pomoci průměr všech možných rozptylů výběrových (35). Musíme si připomenout, že každý rozptyl výběrový je počítán vzhledem k svému výběrovému průměru; z těchto rozptylů se pak vzal průměr. V tom je podstatný rozdíl proti výpočtu rozptylu výběrových průměrů vzhledem k průměru základního souboru (12) resp. (13).

Průměr všech výběrových rozptylů tedy je

$$\bar{\sigma}^2 = \frac{r-1}{r} \frac{N}{N-1} \sigma^2(x). \quad (43)$$

Kdybychom mohli vzít aspoň několik výběrů a našli průměr jejich rozptylů, byl by lepším odhadem veličiny $\bar{\sigma}^2$, než vezmeme-li jeden náhodný výběr o r prvcích, který může dát hodnotu rozptylu, která není daleko od $\bar{\sigma}^2$, ale také nemusí. Rozhodně však budeme považovati veličinu $\bar{\sigma}^2$ za nejlepší odhad rozptylu kolem průměru kteréhokoliv náhodného výběru nebo průměru rozptylů několika výběrů. Poněvadž $\bar{\sigma}^2$ neznáme, budeme považovati obráceně za jeho odhad rozptyl výběrový $\sigma_{x,v}^2$ nějakého výběru dosti velkého rozsahu a položíme tedy

$$\sigma_{x,v}^2 = \frac{(r-1)N}{r(N-1)} \sigma^2(x). \quad (44)$$

Jestliže veškerá informace, kterou máme, je z výběru, musíme podle ní odhadnouti s dobrým přiblížením parametry základního souboru, zde tedy $\sigma(x)$. Vypočítáme tudíž z poslední rovnice

$$\sigma^2(x) = \frac{r(N-1)}{N(r-1)} \sigma_{x,v}^2$$

a tento odhad dosadíme nyní do rovnice pro směrodatnou odchylku σ_P výběrových průměrů (13), takže dostaneme

$$\sigma_P = \sigma_{x,v} \sqrt{\frac{N-r}{N(r-1)}}, \quad (45)$$

kde $\sigma_{x,v}$ je směrodatná odchylka pozorovaného náhodného výběru. Vyskytuje se ovšem v této formuli N čili rozsah základního souboru. Výraz (45) dává nejvhodnější odhad směrodatné odchylky rozdělení četností výběrových průměrů, když základní soubor neznáme.

Je-li základní soubor nekonečného rozsahu, redukuje se pro $N \rightarrow \infty$ poslední výraz (45) na

$$\sigma_P = \sigma_{x,v} \cdot \sqrt{\frac{1}{r-1}} = \frac{\sigma_{x,v}}{\sqrt{r-1}}. \quad (46)$$

Označíme-li ζ_j odchylku hodnoty pozorovaného znaku od průměru výběrového, je $\sigma_{x,v}^2 = \frac{1}{r} \sum_{j=1}^r \zeta_j^2$, takže můžeme rovnici (46) psát

$$\sigma_P = \sqrt{\frac{\sum_{j=1}^r \zeta_j^2}{r-1}} : \sqrt{r}.$$

Čitatel tohoto zlomku

$$\sigma(x, v) = \sqrt{\frac{\sum_{j=1}^r \zeta_j^2}{r-1}}$$

se považuje za nejlepší odhad směrodatné odchylky základního souboru, z něhož náhodný výběr o r prvcích byl vzat.

Můžeme tudíž také směrodatnou odchylku výběrových průměrů vyjádřit vztahem

$$\sigma_P = \frac{\sigma(x, v)}{\sqrt{r}}. \quad (47)$$

Výraz (45) — a výrazy z něho odvozené — je přibližným odhadem, neboť nevíme jak blízko je směrodatná odchyłka jednoho pozorovaného výběru $\sigma_{x,v}$, třeba velkého, odmocnině z průměru všech možných rozptylů výběrových (ovšem z téhož základního souboru) $\sqrt{\overline{\sigma^2}}$, za niž jsme ji položili v rovnici (44).

Činíme jen předpoklad, že je jí blízko, což je přijatelnou hypotézou pro výběry velkého rozsahu.

Rovnice (46) a (47) podávají odhad směrodatné odchyłky průměru z náhodného výběru rozsahu r , má-li základní soubor nekonečný nebo prakticky velmi velký rozsah. Užívá se jich tedy, když je základní soubor neznámý.

(3,3,1) Příklad. Jeden z náhodných výběrů, uvažovaných v tab. 1 s rozsahem $r = 200$ vykázal průměr $\bar{x}_i = 7,03$ a součet čtverců odchylek od průměru $\sum_{j=1}^r \zeta_j^2 = 796,0$. Jest najít odhad směrodatné odchyłky základního souboru, z něhož byl výběr vzat a odhad směrodatné odchyłky výběrových průměrů.

Směrodatná odchyłka výběru je

$$\sigma_{x,v} = \sqrt{\frac{796,0}{200}} = 1,99.$$

Nejvhodnější odhad směrodatné odchyłky základního souboru

$$\sigma(x, v) = \sqrt{\frac{796,0}{199}} = 2,00.$$

Odhad směrodatné odchyłky výběrových průměrů

$$\sigma_P = \frac{2,00}{\sqrt{200}} = 0,14.$$

Jsou tedy meze jedné směrodatné odchylky pro průměr základního souboru

$$\bar{x}_i \pm \sigma_P = 7,03 \pm 0,14,$$

v nichž je s pravděpodobností 0,683.

(3,4) Metoda největší věrohodnosti. V předcházejícím jsme se pokusili o odhad průměru a směrodatné odchylky základního souboru. Nyní ukážeme jednu z obecných metod, jimiž lze odhadnouti parametry rozdělení četností v základním souboru podle pozorovaného výběru a to takovou, že odhady provedené její pomocí jsou efficientní (za předpokladu, že v daném případě efficientní charakteristika parametru existuje). Ze základního souboru nekonečného rozsahu máme výběr rozsahu r , kde hodnoty znaku pozorovaného na jednotlivých prvcích jsou $x_1, x_2, \dots, x_j, \dots, x_r$. Relativní četnost hodnot náhodné proměnné x_i budiž v základním souboru p_i . Jsou-li jednotlivé hodnoty náhodné proměnné vzájemně nezávislé, jak tomu je při rozsahu $N = \infty$, bude pravděpodobnost, že se ve výběru současně vyskytnou hodnoty $x_1, x_2, \dots, x_j, \dots, x_r$ právě v tomto pořadí dána podle věty o násobení pravděpodobností součinem $p_1 \cdot p_2 \cdot \dots \cdot p_r$.

Nezáleží-li pak na pořadí, v němž se vyskytnou jednotlivé hodnoty x_i , nemusíme rozeznávat jednotlivé permutace a pravděpodobnost, že se vyskytne právě takový výběr, jaký máme, bude $P = c \cdot p_1 p_2 \dots p_r$, kde konstanta c je z kombinatoriky známa.

Předpokládejme, že lze každou relativní četnost v základním souboru p_i vyjádřiti pomocí příslušné hodnoty znaku x_i a parametrů základního souboru čili, že známe tvar rozdělení četností. Pro zjednodušení výkladu tedy předpokládejme, že rozdělení četností v základním souboru je normální, takže se v něm vyskytují jen dva parametry $\bar{x}, \sigma(x)$; potom bude

$$p_i = \frac{1}{\sigma(x)\sqrt{2\pi}} e^{-\frac{(x_i - \bar{x})^2}{2\sigma^2(x)}} \quad (48)$$

a hledaná pravděpodobnost je

$$P = c \frac{1}{(\sigma(x)\sqrt{2\pi})^r} e^{-\frac{1}{2\sigma^2(x)} \{(x_1 - \bar{x})^2 + \dots + (x_r - \bar{x})^2\}}. \quad (49)$$

Když hodnoty parametrů neznáme, pak tento výraz pro nějak předpokládané hodnoty parametrů se nazývá věrohodností (vraisemblance, likelihood) předpokládaných hodnot.

Metoda největší věrohodnosti spočívá pak v tom, že se mají zvolit pro parametry takové hodnoty, pro něž bude věrohodnost maximem, čili hodnoty nejvěrohodnější. Místo maxima výrazu (49) můžeme určit maximum jeho logaritmu, neboť větší číslo má větší logaritmus. Označme předpokládané hodnoty \bar{x}_v , $\sigma(x, v)$, potom bude přirozený logaritmus věrohodnosti

$$\lg P = \lg c - \frac{r}{2} \lg 2\pi - r \lg \sigma(x, v) - \frac{1}{2\sigma(x, v)^2} \{(x_1 - \bar{x}_v)^2 + \dots + (\bar{x}_r - \bar{x}_v)^2\}. \quad (50)$$

Pro určení maxima položíme první derivace podle předpokládaných parametrů rovny nule, tedy

$$\frac{\partial \lg P}{\partial \bar{x}_v} = 0, \quad \frac{\partial \lg P}{\partial \sigma(x, v)} = 0$$

a dostaneme snadno

$$\bar{x}_v = \frac{1}{r} (x_1 + \dots + x_r)$$

čili výběrový průměr je odhadem maximální věrohodnosti pro \bar{x} . Podobně dostaneme z druhé rovnice, že směrodatná odchylka $\sigma(x, v)$ je nejvěrohodnějším odhadem $\sigma(x)$.

Můžeme nyní formulovati metodu maxima věrohodnosti zcela obecně tím, že místo zvláštní funkce udávající tvar rozdělení četností (48) označíme obecně $p_i = \varphi(x_i, \Theta_1, \Theta_2, \dots)$, kde Θ_i jsou parametry. Potom bude

$$P = c \prod_{i=1}^r \varphi(x_i, \Theta_1, \Theta_2, \dots),$$

$$\lg P = \lg c + \sum_{i=1}^r \{\lg \varphi(x_i, \Theta_1, \Theta_2, \dots)\}$$

a poněvadž c je konstanta, jedná se o určení maxima funkce

$$L = \sum_{i=1}^r \{\lg \varphi(x_i, \Theta_1, \Theta_2, \dots)\}.$$

Když na př. máme tabulku skupinového rozdělení četností, kde jednotlivé třídní četnosti jsou n_1, n_2, \dots, n_l a odhad četnosti v j -té třídě pro předpokládaný základní soubor označíme v_j , bude o konstantní člen zkrácený logaritmus věrohodnosti

$$L = \sum_{j=1}^l n_j \lg v_j, \quad (51)$$

neboť v tomto případě je

$$P = c \prod_{j=1}^l v_j^{n_j}$$

a dále

$$\lg P = \lg c + \sum_{j=1}^l n_j \lg v_j.$$

(3,4,1) Příklad. V pozorovaném výběru rozsahu r jsou třídní četnosti n_0, n_1, \dots, n_l . O rozdělení četností v základním souboru učiníme hypotézu, že je vyjádřeno Poissonovou exponentiélou, takže absolutní četnosti jsou $\frac{e^{-\lambda} \lambda^x}{x!} \cdot r$ a jejich

logaritmy $x \cdot \lg \lambda - \lambda - \lg \frac{x!}{r}$, kde $x = 0, 1, 2, \dots, l$.

Hledáme nejvěrohodnější odhad pro parametr λ a označíme jej λ_v . Podle (51) dostáváme

$$L = \sum_{x=1}^l n_x \left\{ x \lg \lambda_v - \lambda_v - \lg \frac{x!}{r} \right\}$$

a položíme-li první derivaci podle λ_v rovnu nule, je

$$\frac{\sum x n_x}{\lambda_v} - \sum n_x = 0 \quad \text{čili} \quad \lambda_v = \frac{\sum x n_x}{\sum n_x}.$$

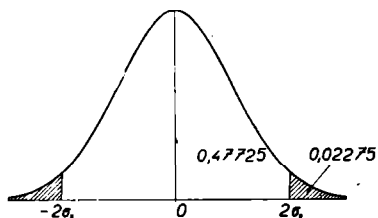
Je tudíž průměr pozorovaného rozdělení četností efficientním odhadem λ , což je zajímavé potud, že také druhý moment tohoto rozdělení je možným odhadem λ [I. díl, rovnice (57)], takže jsme neměli důvodu považovati jej za horší odhad.

(3,5) Testy významnosti. Viděli jsme již, že užíváme ve statistice metody pracovních hypotés jako v tak mnohých případech vědeckého bádání. Použili jsme na př. předpokladu, že rozdělení četností v základním souboru je normální.

Učinili jsme to za tím účelem, abychom převedli zkušenosti získané pozorováním určitého výběru na rozsáhlejší soubor základní, který nemůžeme celý vyšetřiti. Hypotetický soubor s určitými parametry je tu postulátem. Hypotéza je pak přijata, jestliže si vhodným způsobem ověříme, že je možno důvodně soudit, že takový výběr, jaký jsme pozorovali, mohl býti vzat ze základního souboru s těmi vlastnostmi, jež odpovídají hypotéze. Není však určité ostré hranice mezi výběry, které by mohly a které by nemohly vyjít ze základního souboru představeného hypotésou, ježto každá výběrová charakteristika má své rozdělení četností, takže vykazuje výběrové odchylky od své očekávané hodnoty, která je parametrem základního souboru. Je možno udati pouze pravděpodobnost, že by z něho mohl vyjít takový výběr, jako je právě pozorovaný. Je-li tato pravděpodobnost malá, hypotéza se zamítne; je-li velká, hypotéza se přijme a odchylka mezi výběrem a hypotetickým základním souborem se připisuje t. zv. náhodnému kolísání výběrovému, které nám na př. pro výběrový průměr znázorňuje tab. 1 ve sloupci 3 a 4.

Můžeme-li důvodně předpokládat, že základní soubor má normální rozdělení četností pozorovaného znaku s určitým průměrem a směrodatnou odchylkou, je rozdělení četnosti výběrových průměrů pro určité r normální a můžeme z integrálu Laplace-Gaussova odvoditi pravděpodobnost, že ně-

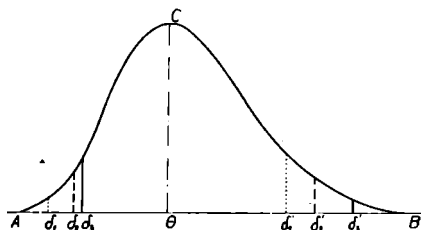
jaký výběrový průměr se odchýlí od parametru, t. j. od průměru v základním souboru více než o danou veličinu [viz příklad (2,4,1,3)]. Tato pravděpodobnost je na př. 0,0027, že se bude odchýlovat nejméně o $\pm 3\sigma_P$, a je tedy malá. Odchyluje-li se tudíž pozorovaný výběrový průměr o trojnásobnou směrodatnou odchylku nebo více, soudíme, že nelze důvodně zařaditi takovou odchylku mezi náhodné a hypotézu zamítáme. Pravděpodobnost, že nějaká náhodná odchylka přesáhne určitou danou hodnotu se nazývá stupeň významnosti a vyjadřuje se často v procentech. Tak na př. odchylka $\pm 3\sigma_P$ je na 0,27 procentním stupni významnosti, nebo odchylka $\pm \sigma_P$ je na 32procentním stupni významnosti. Je otázka, na kterém stupni máme považovati odchylku za statisticky významnou, t. j. tak velkou, že leží na dosti nízkém stupni pravděpodobnosti, aby vedla k zamítnutí hypotézy o základním souboru.



Obr. 3. Testování pomocí normální křivky.

Stalo se zvykem vzítí zřetel k decimálnímu systému a voliti pro tuto pravděpodobnost 0,05 (někdy i 0,01), tedy 5procentní stupeň významnosti. Pro normální křivku odpovídá tomuto stupni přibližně odchylka velikosti dvojnásobné směrodatné odchylky, neboť plocha oddělená pořadnicí v tom bodě je 0,0227 celé plochy křivky na straně kladných odchylek a rovněž tolik na záporné straně (obr. 3). Dvojnásobku směrodatné odchylky $2\sigma_x$ odpovídá také přibližně trojnásobek pravděpodobné chyby $0,6745\sigma_x$. Uvažujeme stupeň významnosti 0,05 pro kladnou i zápornou odchylku, takže je složen ze dvou hodnot 0,025 na každé z obou stran od parametru, ježto není obyčejně důvodu, abychom soudili, že odchylka nebo rozdíl má míti jen jedno nebo jen druhé znaménko.

Je-li rozdělení četnosti testované charakteristiky nesouměrné, je možno několikerým způsobem stanoviti 5% celé plochy křivky. Tak můžeme oddělit 5% celé plochy stejně velkými odchylkami na obě strany od průměru jak je na obr. 4 $\overline{\Theta\delta_1} = \overline{\Theta\delta'_1}$. Nebo můžeme uvažovati část plochy od pořadnice vztyčené v hodnotě parametru Θ (od něhož odchylky testujeme) na-



Obr. 4. Testování pomocí křivky nesymetrické.

pravo zvlášť a nalevo rovněž samostatně a od každé z nich oddělíme 5% plochy, takže pořadnice v δ_2 odděluje 5% plochy $AC\Theta$ a pořadnice v δ'_2 také 5% plochy ΘCB ; znamená to tedy také 5% celé plochy. Konečně můžeme najíti na každé z obou

stran od hodnoty Θ odchylku δ_3 resp. δ'_3 , v níž vztyčená pořadnice odděluje 2,5% celé plochy, takže dohromady je to opět 5% celé plochy. Rozdíly ve výsledcích posledních dvou způsobů nejsou prakticky důležité, ale třetí způsob se zdá výhodnějším.

(3,5,1) Významnost rozdílu mezi dvěma výběrovými průměry. Máme-li rozhodnouti, zda rozdíl mezi průměry dvou pozorovaných náhodných výběrů můžeme pokládati za náhodný či za statisticky významný, musíme především znáti výběrové rozdělení rozdílů mezi průměry náhodných výběrů z téhož základního souboru.

Vydjeme pak od hypotézy, že pozorované dva výběry rozsahů r_1 resp. r_2 jsou ze základních souborů s týmž průměrem a počítáme pravděpodobnost pozorovaného rozdílu $d = \bar{x}_1 - \bar{x}_2$. Bude-li tato pravděpodobnost menší než 0,05, zamítneme hypotézu.

Představme si, že vezmeme z jednoho základního souboru veliký počet výběrů rozsahu r_1 a také z druhého základního

souboru veliký počet výběrů rozsahu r_2 . Pro každý výběr stanovíme průměr a potom pro každý pár výběrů stanovíme rozdíl mezi jejich dvěma průměry, takže dostaneme tolik rozdílů, kolik je párů výběrů a z nich se vytvoří určité rozdělení četností průměrových rozdílů. Bylo odvozeno obecně a je normální s průměrem rovným nule.

Směrodatná odchylka σ_d tohoto rozdělení četností rozdílů d je větší než směrodatná odchylka výběrových průměrů $\sigma_{P,1}$ resp. $\sigma_{P,2}$, neboť víme [I, rovnice (67')], že $\sigma_d^2 = \sigma_{P,1}^2 + \sigma_{P,2}^2$, jsou-li oba výběry na sobě nezávislé. Pomocí rozptylů základních souborů $\sigma(x_1)$, $\sigma(x_2)$, jakožto jejich parametrů, dostáváme vzhledem k rovnici (24)

$$\sigma_d^2 = \frac{\sigma^2(x_1)}{r_1} + \frac{\sigma^2(x_2)}{r_2}. \quad (52)$$

Pojmeme-li do naší hypotézy, že se oba základní soubory shodují čili, že také $\sigma^2(x_1) = \sigma^2(x_2)$, pak při stejném rozsahu obou výběrů $r_1 = r_2$ bude také $\sigma_{P,1}^2 = \sigma_{P,2}^2 = \sigma_P^2$ čili $\sigma_d = \sqrt{2}\sigma_P$. Kriterium dvojnásobné směrodatné odchylky $2\sigma_d$ pro stupeň významnosti 0,05 vede v tomto případě testování výběrových rozdílů k pracovnímu pravidlu, že statisticky významné jsou rozdíly větší než trojnásobek směrodatné odchylky jednotlivého průměru $3\sigma_P$, neboť $2\sigma_d = 2\sqrt{2}\sigma_P$ a $2\sqrt{2} \doteq 3$. Předpokládali jsme v této úvaze, že známe směrodatnou odchylku $\sigma(x)$ základního souboru. Není-li známa, nahradí se odhadem z výběru, čímž je ovšem dáno omezení pro aplikaci teorie velkých výběrů. Když tedy vycházíme od hypotézy, že pozorované dva výběry jsou z téhož základního souboru, užijeme za odhad vhodné kombinace obou výběrových rozptylů; za takovou můžeme užítí výrazu

$$\sigma_{x,v}^2 = \frac{r_1\sigma_{x,1}^2 + r_2\sigma_{x,2}^2}{r_1 + r_2},$$

který dosadíme do (52) za $\sigma^2(x_1)$ a také za $\sigma^2(x_2)$, které jsou stejné, takže dostaneme

$$\begin{aligned} & \frac{r_1\sigma_{x,1}^2 + r_2\sigma_{x,2}^2}{r_1(r_1 + r_2)} + \frac{r_1\sigma_{x,1}^2 + r_2\sigma_{x,2}^2}{r_2(r_1 + r_2)} = \\ & = \frac{(r_1 + r_2)r_1\sigma_{x,1}^2 + (r_1 + r_2)r_2\sigma_{x,2}^2}{r_1r_2(r_1 + r_2)} \end{aligned}$$

a odhad σ_d^2 tedy bude

$$\sigma_{d,v}^2 = \frac{\sigma_{x,1}^2}{r_2} + \frac{\sigma_{x,2}^2}{r_1}.$$

(3,5,2) Příklad. Dva výběry rozsahu $r_1 = r_2 = 200$ vykazují průměry $\bar{x}_1 = 6,68$, $\bar{x}_2 = 7,37$ a směrodatné odchylky $\sigma_{x,1} = 2,05$, $\sigma_{x,2} = 2,02$.

Máme považovati rozdíl mezi těmito průměry $d = 0,69$ za statisticky významný?

Odhad směrodatné odchylky bude

$$\sigma_{d,v} = \sqrt{\frac{4,20 + 4,08}{200}} = 0,203.$$

Vidíme tedy, že $d > 2\sigma_{d,v}$ a že tedy rozdíl nepovažujeme za náhodný, neboť jen v 5 případech ze sta dostaneme rozdíly větší než 0,406 (srovnej s tabulkou 1 sloupec 3 a 4).

(3,6) Ověřování hypotés. Objasníme ještě odpověď na otázku jak voliti stupeň statistické významnosti pro ověřování hypotés. Podle předešlého výkladu vidíme, že chceme postupovat podle pravidla: hypotéza je přijatelná, když pozorovaná odchylka (nebo rozdíl) je pod daným stupněm významnosti; je-li nad ním, je zamítnuta. Při tomto postupu můžeme rozhodnout nesprávně

1. tím, že zamítneme správnou hypotézu,
2. tím, že přijmeme nesprávnou hypotézu,

nebo můžeme rozhodnout správně tím,

3. že přijmeme správnou hypotézu, nebo
4. že zamítneme nesprávnou hypotézu.

Víme, že žádná hypotéza nemůže být dokázána s nějakou konečnou platností, takže naše rozhodování má ráz pokusný a naší snahou je především, abychom v dlouhé řadě případů statistické praxe jen v málo případech zamítli správnou hypotézu, abychom tedy poměr počtu případů nesprávného rozhodnutí 1. k počtu všech rozhodnutí 1. a 3. udělali libovolně malým. Toho můžeme dosáhnouti stanovením vysoké hranice odchylky a tedy vysokého stupně významnosti, již odpovídá nízká pravděpodobnost. Navržené pravidlo 5procentního stupně významnosti vede tedy k zamítnutí správné hypotézy průměrně jednou v každých 20 pokusech, pro něž jsme zvolili správnou hypotézu. Kdybychom přijali 1procentní stupeň významnosti, takže pravděpodobnost by byla 0,01, přihodilo by se nám takové nesprávné rozhodnutí průměrně jen jednou v každých 100 pokusech. — Tím bychom sice snížili četnost chyby tohoto druhu, ale zvýšili bychom možnost chyby druhého druhu, neboť jsme tím usnadnili možnost přijmouti hypotézu a tedy také přijmout nesprávnou hypotézu čili učiniti chybu (2). Tak tedy bude volba stupně statistické významnosti a tím stanovení kritéria pro přijetí nebo zamítnutí hypotézy jistým kompromisem mezi nebezpečím dvou druhů chyb. Nelze však stanovit poměr chyb (2) k celkovému počtu případů, v nichž byla zvolena nějaká nesprávná hypotéza; ten bude záležitosti značně na tom, jak je hypotéza blízko pravdě a jak je test přísný. Proto je volba kritické hranice významnosti do veliké míry ovládána velikým množstvím vědomostí z oboru, v němž se šetření provádí a jistou vědeckou tradicí. V důsledku této tradice se na př. doporučují jednoduché hypotézy, t. j. takové, které zahrnují málo konstant před složitými.

Bylo by možno zmenšiti vliv chybných úsudků druhého druhu aniž se dotkneme prvního druhu, kdybychom redukovali směrodatnou odchylku charakteristiky, kterou testujeme. Pro určitý stupeň významnosti totiž vede menší směrodatná odchylka σ k příslušné menší odchylce $\tau\sigma$, která leží právě na naší zvolené hranici významnosti a tedy vede

k menší odchylce, kterou nesprávně odmítneme jako nevýznamnou. Směrodatnou odchylku pak lze redukovat [viz rovnice (24), (39), (40), (41)] zvětšením rozsahu výběru.

Musíme mít také stále na paměti, že označíme-li podle zvoleného pravidla nějakou odchylku za nevýznamnou, znamená to spíše, že její významnost není prokázána. Statistická významnost nedává posudek o velikosti nebo praktické důležitosti nějaké odchylky či rozdílu, neboť ty mohou být posouzeny jen na základě vědomostí z oboru, do něhož předmět šetření spadá. Je-li nějaký průměr výběrový významně odlišný od hodnoty průměru v hypotetickém souboru, nebo je-li rozdíl mezi dvěma výběrovými průměry významný, nemůže statistická teorie osvědčiti příčinu této odchylky. Příčina může být v tom, že základní soubor je skutečně různý od hypotetického, nebo výběr není vskutku reprezentativním a náhodným. Může tam být skutečný rozdíl mezi dvěma základními soubory nebo rozdíl ve výběrové technice.

Testování nesouměrnosti pozorovaného rozdělení četnosti vzhledem k normálnímu lze provést pomocí míry $\sqrt{\beta_1} = \alpha_{x,3}$, která má v prvním přiblížení směrodatnou odchylku $\sqrt{\frac{6}{r}}$, jak jsme již uvedli, a rozdělení četnosti blízké normálnímu, takže dvojnásobná směrodatná odchylka je prakticky na 5procentní hranici významnosti.

(3,7) Náhodný výběr malého rozsahu. Vyložili jsme vhodné metody k testování statistické významnosti, jichž lze užívatí pro výběry velkých rozsahů. Nesmíme jich však užívatí, testujeme-li na př. významnost mezi průměry malých výběrů. Proto byla odvozena teorie, která dává přesné testy bez ohledu na rozsah výběru, tedy vhodná pro malé výběry. Spočívá v užití přesných výběrových rozdělení četností místo přibližných, jichž se užilo pro velké výběry. Tím jsme se sice zbavili vlivu výběrového rozsahu, ale ve většině dalších výsledků zůstává omezení, že se aplikují jen na veličiny s normálním rozdělením četností.

Charakteristiky odvozené z velkých výběrů dávají spolehlivé odhady parametrů v základním souboru, kdežto malé výběry poskytují chudé odhady. Každá charakteristika má své rozdělení četností. Soudilo se dlouho klamně, že tvar tohoto rozdělení závisí na rozsahu výběru, ale ve skutečnosti se nejedná o počet všech prvků r , nýbrž o počet nezávislých prvků. Máme-li na př. krabičku zápalek, z nichž každá je nějak označena třeba počtem čárek, které jsme na ni tužkou udělali a máme rozdělit zápalky do pěti skupin, můžeme měnit čtyři skupiny, ale pátá je vždy již určena celkovým počtem zápalek. Máme tedy zde jen čtyři volné cesty, jimiž můžeme prováděti libovolné skupiny. Počítáme-li z nějakého výběru s rozsahem 25 prvků průměr a máme-li vzítí z téhož základního souboru druhý náhodný výběr 25 prvků, který by měl též výběrový průměr jako první, můžeme vzítí libovolně jen 24 prvků tedy $(r - 1)$, neboť poslední je již danou podmínkou určen. Pro tyto případy byl zaveden pojem „stupně volnosti“, takže v tomto případě máme na př. 24 stupňů volnosti pro odhad charakteristiky za uvedené podmínky.

Shledali jsme již v odstavci (3,3), že nejlepší odhad rozptylu $\sigma^2(x)$ dostaneme, dělíme-li součet čtverců odchylek od průměru počtem stupňů volnosti $r - 1$ a nikoliv počtem pozorování r .

Počet stupňů volnosti se zde rovná počtu odchylek zmenšenému o počet konstant určených z výběru, jichž bylo použito k pevnému stanovení bodu, od něhož jsou odchylky měřeny, tedy o jednu, ježto se našel jen průměr z výběru.

(3,7,1) Příklad. K osvětlení vlivu, který má užití počtu stupňů volnosti při odhadu rozptylu, bylo vzato náhodně $r = 200$ čísel dvojciferných mezi 10 a 50. Jsou to náhodná čísla z tabulek Tippettových (viz str. 73). Z těchto čísel bylo utvořeno 10 skupin, které budeme nazývatí varietami, po $s = 20$ číslech (tab. 2). Můžeme si představit celkovou variaci všech těchto čísel ze dvou složek, jednak ze složky variace uvnitř skupin (variet), jednak ze složky mezi skupinami (varietami).

Tabulka 2.

Běžné číslo	I	II	III	IV	V	VI	VII	VIII	IX	X
1	29	45	14	25	11	47	28	35	18	25
2	45	39	49	32	32	36	24	27	16	39
3	16	29	29	18	46	42	10	18	34	24
4	37	11	31	28	44	36	27	44	18	30
5	50	12	19	20	28	38	11	25	30	24
6	10	37	20	44	40	21	42	33	29	36
7	26	44	49	41	27	41	22	49	35	31
8	19	15	15	10	28	26	30	11	35	10
9	24	50	11	43	27	17	17	17	42	13
10	10	14	22	19	11	50	33	39	50	43
11	10	22	23	10	48	30	44	26	21	27
12	48	20	41	13	21	39	32	29	11	20
13	22	46	40	31	44	21	23	16	45	39
14	13	14	12	45	16	46	25	47	18	30
15	28	21	39	39	36	22	27	10	31	18
16	32	15	43	23	42	34	16	20	26	11
17	10	37	31	11	12	50	20	12	34	46
18	11	26	34	22	48	13	47	42	22	43
19	30	29	49	35	30	46	38	50	24	44
20	37	22	37	49	30	47	12	34	42	24
	507	548	608	558	621	702	528	584	581	577

Oba prameny variace budou vyrovnány, byla-li tato čísla vzata náhodně; nebyly by vyrovnány, kdyby některé skupiny (variety) měly na př. všechna malá čísla a jiné zase všechna velká čísla. Náhodný výběr čísel zajišťuje, že tento případ nenastane. Naše úvaha o variaci znamená, že rozptyly uvnitř skupin, mezi skupinami a rozptyl celkový budou stejné až na odchylky v mezích náhodného výběru. Je-li rozptyl mezi skupinami velmi blízko celkovému rozptylu, takže se mu skoro rovná, musí být také rozptyl uvnitř skupin skoro přesně roven celkovému rozptylu.

Uřídíme tedy rozptyl uvnitř každé variety a průměr jejich by měl být velmi blízko rozptylu pro celý výběr, je-li naše metoda správnou. Sestavíme výpočet rozptylů do následu-

jící tabulky 3. Poněvadž v každé skupině je 20 čísel, bude tam 19 stupňů volnosti pro odhad rozptylu. Dělili jsme tedy v sloupci (7) součet čtverců počtem stupňů volnosti a v sloupci (8) rozsahem souboru.

K sestavení tabulky 3 se doporučuje napsati ještě pomocnou tabulku čtverců čísel tabulky 2, která bude míti jako součty v posledním řádku čísla, jež potřebujeme v 3. sloupci tab. 3. Výpočet sloupce 6 se provádí podle vztahu

$$\sum(x - \bar{x}_i)^2 = \sum x^2 - \frac{1}{s} (\sum x)^2,$$

takže je třeba poznamenati si také pomocný sloupec čtverců čísel druhého sloupce tab. 3, abychom z něho mohli psáti pro $s = 20$ sloupec 5.

Tabulka 3.

Varieta	Σx	Σx^2	\bar{x}_i	$\frac{1}{s} (\Sigma x)^2$	$\Sigma(x - \bar{x}_i)^2$	$\frac{\Sigma(x - \bar{x}_i)^2 : 19}{}$	$\frac{\Sigma(x - \bar{x}_i)^2 : 20}{}$
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
I	507	16 159	25,35	12852,45	3306,55	174,0289	165,3275
II	548	18 110	27,40	15015,20	3094,80	162,8842	154,7400
III	608	21 602	30,40	18483,20	3118,80	164,1474	155,9400
IV	558	18 620	27,90	15568,20	3051,80	160,6210	152,5900
V	621	22 189	31,05	19282,05	2906,95	152,9974	145,3475
VI	702	27 208	35,10	24640,20	2567,80	135,1474	128,3900
VII	528	16 132	26,40	13939,20	2192,80	115,4105	109,6400
VIII	584	20 306	29,20	17052,80	3253,20	171,2210	162,6600
IX	581	19 043	29,05	16878,05	2164,95	113,9447	108,2475
X	577	19 045	28,85	16646,45	2398,55	126,2395	119,9275
I—X	5814	198 414	Průměr sloupce (7) resp. (8)			147,6642	140,2810

Pro samostatný výpočet celkového rozptylu použijeme jednak 199 stupňů volnosti, jednak celý rozsah výběru $r = 200$.

Výpočet $\sum(x - \bar{x})^2$ pak provedeme opět podle vztahu

$$\sum (x - \bar{x})^2 = \sum x^2 - \frac{1}{r} (\sum x)^2, \text{ kde } \sum x^2 = 198\,414, \frac{1}{r} (\sum x)^2 =$$

$$= 33\,802\,596 : 200 = 169\,012,98, \text{ takže } \sum (x - \bar{x})^2 = 29401,02$$

a dále $\frac{1}{99} \sum (x - \bar{x})^2 = 147,74$, kdežto $\frac{1}{200} \sum (x - \bar{x})^2 =$
 $= 147,00$.

Dostáváme tedy čtyři výsledky:

Pomocí stupňů volnosti:

Průměrný rozptyl uvnitř skupin	147,66
Celkový rozptyl výběru	147,74

Pomocí rozsahu výběru:

Průměrný rozptyl uvnitř skupin	140,28
Celkový rozptyl výběru	147,00

a vidíme, že při užití počtu stupňů volnosti dává průměrný rozptyl uvnitř skupin hodnotu rovnou 99,94% celkového rozptylu, kdežto pomocí rozsahu výběru je jen 95,43% celkového rozptylu. V tomto případě je tudíž skutečná hodnota podhodnocena o 4,57%. Osvětlili jsme tak správnost odhadu rozptylu pomocí počtu stupňů volnosti.

(3,8) Významnost průměrů. t-test. Při testování statistické významnosti odchylky výběrového průměru od předpokládané hodnoty základního souboru (hypotézy) jsme použili té skutečnosti, že poměr odchylky k její směrodatné odchylce (čili odchylka vyjádřená ve směrodatné odchylce jako jednotce) má normální rozdělení četnosti se směrodatnou odchylkou rovnou jednotce.

Označíme-li odchylku $d = \bar{x}_v - \bar{x}$, je tento poměr $d : \sigma_P$, kde $\sigma_P = \sigma(x) : \sqrt{r}$. Předpokládá se tu tedy, že směrodatná odchylka základního souboru $\sigma(x)$ je známa. Chceme-li však v praxi ověřiti významnost nějakého průměru, je odhad $\sigma(x, v)$ získaný z výběru vše, co známe. Jedná-li se o výběr velkého rozsahu, je tento odhad $\sigma(x, v)$ dostatečně blízky svému parametru $\sigma(x)$ a je možno dřívějšího postupu užití. Je-li výběr malý, je třeba jisté úpravy vzhledem k nepřes-

nosti, kterou zahrnujeme tím, že užitíme $\sigma(x, v)$ místo $\sigma(x)$. Poměr, jímž testujeme významnost odchylky je nyní $t = d : \frac{\sigma(x, v)}{\sqrt{r}}$. Rozdělení četnosti hodnot t nám umožňuje

provést test přesně bez omezení na velký rozsah výběru. Ale toto rozdělení se rozhodně liší od normálního rozdělení četností hodnot $d : \frac{\sigma(x)}{\sqrt{r}}$, když rozsah výběru a tedy počet

stupňů volnosti je malý a je prakticky s ním shodné pro velký rozsah výběru. To poprvé vyzvedl „Student“ (1908) a příslušné rozdělení t pro výběry z normálního základního souboru a počet n stupňů volnosti t je dáno funkcí

$$y = y_0 \left(1 + \frac{t^2}{n} \right)^{-\frac{n+1}{2}}$$

Je patrné, že tedy výběrové rozdělení hodnot t je symetrické vzhledem ku $t = 0$ a závisí jen na n , počtu stupňů volnosti, jimiž byla odhadnuta směrodatná odchylka, takže v tomto případě $n = r - 1$, neboť

$$\sigma^2(x, v) = \frac{\sum (x - \bar{x})^2}{r - 1}.$$

t -křivky odpovídající funkci y mají modus spadající do průměru v $t = 0$, neboť výraz $\left(1 + \frac{t^2}{n} \right)^{-1}$ klesá, když t roste, a

na obou stranách jdou větve do nekonečna, jsou jen špičatější ($\beta_3 > 3$) než normální křivka. Také čtenář vidí, že pro

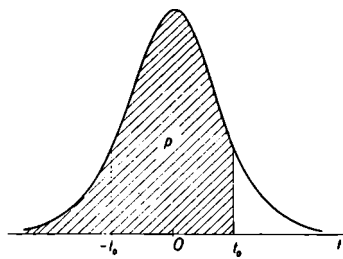
$n \rightarrow \infty$ spěje výraz $\left(1 + \frac{t^2}{n} \right)^{-\frac{n+1}{2}}$ ku $e^{-\frac{t^2}{2}}$, takže je zřej-

mo, že pro velká n je t -rozdělení normální. Pravděpodobnost, že při náhodných výběrech dostaneme nějakou hodnotu t , která není větší než t_0 , je dána obsahem plochy omezené křivkou, osou t a pořadnicí vztyčenou v bodě t_0 ; stručně

říkáme plochou křivky až k pořadnici vztyčené v bodě t_0 , čili

$$p(t_0) = \int_{-\infty}^{t_0} y_0 \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}} dt$$

obdobně rovnici [I, (50)]. Označíme-li p obecně pravděpodobnost, že pozorovaná hodnota nepřekročí určitou mez t_0 ,



Obr. 5. Testování pomocí t -rozdělení.

kdežto P bude pravděpodobnost, že pozorovaná hodnota překročí t_0 a to bez ohledu na znaménko, potom bude p plocha křivky nalevo od pořadnice v t_0 (obr. 5), P bude plocha napravo od t_0 plus plocha nalevo od minus t_0 čili (obr. 5) dvojnásobek plochy napravo od t_0 , ježto se jedná o křivky symetrické. Bude tedy $P = 2(1 - p)$.

Sestavíme si několik hodnot pravděpodobností p a P pro srovnání s křivkou normální do tabulky 4.

Tabulka 4.

t	p			$P = 2(1 - p)$		
	$n = 10$	$n = 15$	norm.	$n = 10$	$n = 15$	norm.
0	0,500	0,500	0,500	1,000	1,000	1,000
0,6745	0,742	0,745	0,750	0,516	0,510	0,500
1,0	0,830	0,833	0,841	0,340	0,334	0,318
2,0	0,963	0,968	0,977	0,074	0,064	0,046
2,6	0,987	0,990	0,995	0,026	0,020	0,010
3,0	0,993	0,995	0,999	0,014	0,010	0,002

Testujeme-li na 5% stupni významnosti, znamená to, že $P \geq 0,05$, čili $p \geq 0,975$. Jeví se pro praktickou potřebu

testování významnosti účelným sestavením pro různé stupně volnosti tabulku hodnot t (ležících na některých stupních významnosti na př. 0,1; 0,05; 0,01. Vliv větší špičatosti t -křivek proti křivce normální je vážný pro výběry rozsahu menšího než $r = 20$.

Tabulka 5.

Hodnoty t .

$P \backslash n$	1	2	3	5	10	15	20	25	∞
0,10	6,31	2,92	2,35	2,02	1,81	1,75	1,73	1,71	1,64
0,05	12,71	4,30	3,18	2,57	2,23	2,13	2,09	2,06	1,96
0,01	63,66	9,93	5,84	4,03	3,17	2,95	2,85	2,79	2,58

Pro velké výběry leží hodnota $t = 2,0$ na 5% stupni významnosti: jiné hodnoty t , které jsou na téměř stupni, vidíme v tabulce 5 pro $P = 0,05$.

Poznámka: Hodnota t je v podstatě poměr nějaké charakteristiky s normálním rozdělením četností a průměrem v nule k odhadu směrodatné odchylky této charakteristiky provedenému s n stupni volnosti. Každý takový poměr, ať vzniká jakkoliv, má totéž výběrové rozdělení četností jako t .

(3,8,1) Příklad. Pevnost nějakého materiálu byla měřena na deseti kusech náhodně vybraných a byly zjištěny hodnoty znaku 63, 63, 66, 67, 68, 69, 70, 70, 71, 71 jednotek. Pomocí dat tohoto náhodného výběru jest ověřiti hypotézu, že průměrná pevnost celého materiálu je 66 jednotek.

Předpokládáme, že rozdělení četností v základním souboru je normální. Výběrový průměr je $\bar{x}_v = 67,8$ a odhad směrodatné odchylky $\sigma(x, v) = 3,011$. Podle rovnice $t = d : \frac{\sigma(x, v)}{\sqrt{r}}$

bude $t = \frac{67,8 - 66}{3,011} \sqrt{10} = 1,89$. Testujeme-li na 5% stupni významnosti, vidíme z tab. 5, že pro $n = 9$ je hraniční hodnota t větší než 2,23, takže naši pozorovanou hodnotu nepo-

važujeme za statisticky významnou a nemusíme naši hypotézu zamítnouti.

(3,9) Významnost rozdílu mezi průměry. Pomocí t -testu můžeme také ověřovati statistickou významnost rozdílu mezi dvěma výběrovými průměry. Omezíme se na případ, v němž činíme hypotézu, že výběry jsou ze základních souborů, majících společnou směrodatnou odchylku $\sigma(x)$ a též průměr \bar{x} . Poněvadž předpokládáme také normální rozdělení četností základního souboru, znamená tato hypotéza, že výběry jsou z téhož základního souboru. Podle zkušenosti nejsou testy příliš citlivé na mírné odchylky od normálního rozdělení. Při úvaze o těchto testech vycházíme z představy nekonečného základního souboru diferencí, jejichž průměr se rovná nule. Z pozorování jsme dostali dva výběry, jejichž průměry jsou různé, takže tedy vykazují určitou diferenci. Tážeme se, v jakém procentu případů takových dvojic výběrů dostaneme průměrně diferenci tak velkou jako je pozorovaná nebo větší. Testy tohoto druhu musíme rozdělit do dvou skupin.

a) Nejprve máme případ dvou výběrů různého rozsahu r_1 resp. r_2 , jejichž prvky jsou na sobě úplně nezávislé, takže hodnoty proměnné netvoří dvojice k sobě nějak vázané. Jsou-li jejich průměry \bar{x}_1 resp. \bar{x}_2 , pak diference $d = \bar{x}_1 - \bar{x}_2$ má normální rozdělení četností kolem nuly a odhad její směrodatné odchylky σ_d provedeme podle rovnice [I, (67')], která praví, že rozptyl rozdílu dvou proměnných nezávislých se rovná součtu rozptylů každé z nich. Pro odhad rozptylu v základním souboru použijeme kombinace součtu čtverců odchylek od jejich průměrů, kterou dělíme počtem stupňů volnosti $r_1 + r_2 - 2$, neboť dva průměry byly stanoveny, takže bude

$$\sigma^2(x, v) = \frac{\sum_{r_1} (x_1 - \bar{x}_1)^2 + \sum_{r_2} (x_2 - \bar{x}_2)^2}{r_1 + r_2 - 2},$$

směrodatné odchytky průměrů budou (24)

$$\sigma_{P_1} = \frac{\sigma(x, v)}{\sqrt{r_1}}, \quad \sigma_{P_2} = \frac{\sigma(x, v)}{\sqrt{r_2}}$$

a směrodatná odchytka rozdílů průměrů výběrových

$$\begin{aligned} \sigma_d &= \sqrt{\sigma_{P_1}^2 + \sigma_{P_2}^2} = \sqrt{\frac{\sigma^2(x, v)}{r_1} + \frac{\sigma^2(x, v)}{r_2}} = \\ &= \sigma(x, v) \sqrt{\frac{r_1 + r_2}{r_1 r_2}}. \end{aligned} \quad (53)$$

Hodnota t tudíž bude

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sigma(x, v)} \sqrt{\frac{r_1 r_2}{r_1 + r_2}} \quad (54)$$

a má rozdělení podle křivky t v malých výběrech, ve velkých pak normální rozdělení s jednotkovou směrodatnou odchylkou. Jisté zjednodušení nastává ještě, jsou-li oba výběry stejného rozsahu $r_1 = r_2 = r$, potom

$$\left. \begin{aligned} \sigma^2(x, v) &= \frac{\sum_r (x - \bar{x}_1)^2 + \sum_r (x - \bar{x}_2)^2}{2(r-1)}, \\ t &= \frac{\bar{x}_1 - \bar{x}_2}{\sigma(x, v)} \sqrt{\frac{r}{2}}. \end{aligned} \right\} \quad (55)$$

Do tabulky 5, hodnot t vstupujeme s počtem stupňů volnosti $n = r_1 + r_2 - 2$.

b) Druhý případ musíme rozeznávat, nejsou-li proměnné nezávislé čili každá hodnota proměnné x_1 je sdružena nějakou logickou cestou s příslušnou hodnotou proměnné x_2 a tvoří tedy dvojice. V takových případech budou mít oba výběry stejný rozsah, takže bude r dvojic hodnot proměnných. Rozptyl nemůžeme stanovit jako v předešlém případě, nýbrž

$$\sigma^2(d, v) = \frac{\sum [(x_1 - \bar{x}_1) - (x_2 - \bar{x}_2)]^2}{2(r-1)}$$

a tento výraz můžeme upravit na tvar

$$\frac{1}{2(r-1)} \sum [(x_1 - x_2) - (\bar{x}_1 - \bar{x}_2)]^2.$$

Vzhledem k tomu, že součet čtverců lze rozvésti

$$\begin{aligned} & \sum_r (x_1 - x_2)^2 - 2(\bar{x}_1 - \bar{x}_2) \sum_r (x_1 - x_2) + r(\bar{x}_1 - \bar{x}_2)^2 = \\ & = \sum_r (x_1 - x_2)^2 - 2(\bar{x}_1 - \bar{x}_2) r \left(\frac{\sum x_1}{r} - \frac{\sum x_2}{r} \right) + r(\bar{x}_1 - \bar{x}_2)^2 = \\ & = \sum_r (x_1 - x_2)^2 - r(\bar{x}_1 - \bar{x}_2)^2, \end{aligned}$$

dostáváme

$$\sigma^2(d, v) = \frac{1}{2(r-1)} \left[\sum (x_1 - x_2)^2 - \frac{(\sum x_1 - \sum x_2)^2}{r} \right]. \quad (56)$$

Hodnota t tedy bude

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sigma(d, v)} \sqrt{\frac{r}{2}}. \quad (57)$$

(3,9,1) Příklad 1. Ve dvou náhodných výběrech rozsahu $r_1 = 9$, $r_2 = 7$ byly zjištěny průměry $\bar{x}_1 = 196,42$, $\bar{x}_2 = 198,82$, takže diference $d = 2,40$. Byly vzaty nezávisle, takže není důvodu k předpokladu, že jsou v nějaké závislosti. Součty čtverců odchylek od průměrů jsme stanovili $\sum (x_1 - \bar{x}_1)^2 = 26,94$, $\sum (x_2 - \bar{x}_2)^2 = 18,73$, což dává dohromady 45,67. Potom je $\sigma(x, v) = \sqrt{\frac{45,67}{14}} = 1,81$ a tudíž

$$t = \frac{2,40}{1,81} \cdot \sqrt{\frac{9}{14}} = 2,62.$$

Testujeme-li na 5% stupni významnosti, najdeme v tabulce 5, že pro $n = 14$ je přibližně $t = 2,15$, takže rozdíl mezi průměry považujeme za významný.

Příklad 2. Týž druh pšenice vyrostlé ve dvou různých oblastech se zkoumá na obsah proteinu. Z první oblasti bylo pět vzorků s výsledky 12,6; 13,4; 11,9; 12,8; 13,0, z druhé oblasti sedm vzorků s výsledky 13,1; 13,4; 12,8; 13,5; 13,3; 12,7; 12,4. Je tedy v první oblasti průměr $\bar{x}_1 = 12,740$ a v druhé $\bar{x}_2 = 13,029$. Není-li možno opatřiti další vzorky, jest podle těchto ověřiti, je-li rozdíl mezi průměry významný.

$$\begin{array}{r} \Sigma(x_1 - \bar{x}_1)^2 = 1,2320 \\ \Sigma(x_2 - \bar{x}_2)^2 = 0,9943 \\ \hline \text{Celkem} \quad \quad \quad 2,2263 \end{array}$$

$$\sigma(x, v) = \sqrt{\frac{2,2263}{10}} = 0,472$$

$$t = \frac{0,289}{0,472} \sqrt{\frac{3}{1} \frac{5}{2}} = 1,047.$$

Z tab. 5 je patrné pro $n = 10$, že tato hodnota $t = 1,05$ není významnou na hranici 1% ani 5%, tedy rozdíl mezi průměry $\bar{x}_1 - \bar{x}_2 = -0,289$ můžeme pokládati za nevýznamný.

Příklad 3. Byl zkoumán vliv dvou příbuzných krmiv A, B na vývoj párů zvířat a výsledek je obsažen v těchto číslech

Pár	1	2	3	4	5	6	7	8
A	49,2	53,3	50,6	52,0	46,8	50,5	52,1	53,0
B	51,5	54,9	52,2	53,3	51,6	54,1	54,2	53,3

Průměry jsou $\bar{x}_A = 50,94$, $\bar{x}_B = 53,14$.

Zkoumejme významnost difference $d = 2,20$ mezi průměry

a) za předpokladu, že hodnoty pozorované nejsou vázány na dvojice,

b) za předpokladu, že tvoří dvojice.

ad a) Za předpokladu, že výběry jsou na sobě nezávislé, bude

$$\Sigma(x - \bar{x}_A)^2 = 32,7587$$

$$\Sigma(x - \bar{x}_B)^2 = 11,1387$$

$$\text{Dohromady } 43,8974$$

$$\sigma^2(x, v) = 43,8974 : 14 = 3,1355$$

$$\sigma(x, v) = 1,77$$

$$t = \frac{2,20}{1,77} \sqrt{\frac{64}{16}} = 2,486.$$

Z tab. 5 vidíme, interpolujeme-li lineárně pro $n = 14$, že hodnota $t = 2,49$ je při testování na hranici 5% významnou, kdežto na hranici 1% není významnou.

ad b) Patří-li stejně očíslované dvojice obou výběrů k sobě, pak dostáváme

$$\Sigma(x_1 - x_2)^2 = 52,60, \quad \Sigma x_1 = 407,5, \quad \Sigma x_2 = 425,1,$$

$$\Sigma x_1 - \Sigma x_2 = -17,6, \quad \frac{(\Sigma x_1 - \Sigma x_2)^2}{r} = 38,72,$$

$$\sigma^2(d, v) = \frac{52,60 - 38,72}{14} = 0,9914, \quad \sigma(d, v) = 0,9957,$$

$$t = \frac{2,20}{0,9957} \cdot 2 = 4,42.$$

V tomto případě je hodnota $t = 4,42$ významnou při testování jak na 5% tak na 1% hranici. Považujeme tedy rozdíl $d = 2,02$ ve výběrech, v nichž pozorované hodnoty tvoří dvojice za významný.

(3,10) Rozšíření t -testu na tři výběry. Kdybychom měli srovnávat tři výběry nezávislých pozorování, bylo by to jednoduché rozšíření prvního testu. Odhad rozptylu by byl

$$\sigma^2(x, v) = \frac{\Sigma(x_1 - \bar{x}_1)^2 + \Sigma(x_2 - \bar{x}_2)^2 + \Sigma(x_3 - \bar{x}_3)^2}{r_1 + r_2 + r_3 - 3},$$

takže směrodatné odchylky průměrů jsou

$$\frac{\sigma(x, v)}{\sqrt{r_1}}, \quad \frac{\sigma(x, v)}{\sqrt{r_2}}, \quad \frac{\sigma(x, v)}{\sqrt{r_3}}$$

a hodnoty t pro difference mezi průměry budou

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sigma(x, v)} \sqrt{\frac{r_1 r_2}{r_1 + r_2}}, \quad t = \frac{\bar{x}_1 - \bar{x}_3}{\sigma(x, v)} \sqrt{\frac{r_1 r_3}{r_1 + r_3}},$$

$$t = \frac{\bar{x}_2 - \bar{x}_3}{\sigma(x, v)} \sqrt{\frac{r_2 r_3}{r_2 + r_3}}.$$

Úloha: Testujte významnost rozdílů mezi průměry variet uvedenými ve 4. sloupci tabulky 3 na 5% hranici významnosti

- a) mezi kterýmikoliv dvojicemi,
- b) mezi některými trojicemi.

Poznámka: Testujeme-li významnost nějaké průměrové difference, činíme tak na základě určité pravděpodobnosti, že dostaneme tak velké nebo větší difference než je pozorovaná a to buď znaménka kladného nebo záporného. Tak víme, že hodnota t na pětiprocentním stupni významnosti je ta, která odetíná svou pořadnicí 2,5% celkové plochy křivky na pravé straně od parametru a tolikéž na levé straně (obr. 5). Představme si nyní případ, že máme testovati výsledek nějaké setby, při níž bylo semeno nějakým postupem mořeno. Víme již, že tento postup semení prospívá, takže úrodu zvyšuje; chceme srovnati jeho vliv s výsledkem kontrolní setby, kde tohoto postupu nebylo užito. V takovém případě můžeme uvažovati jen kladné odchylky a pracovati s hranicí významnosti v bodě, kde pořadnice t odetíná 5% plochy křivky jen na kladné straně. Zde se zdá logicky oprávněnějším založiti test na pravděpodobnosti, že dostaneme kladnou diferencii tak velkou nebo větší než je pozorovaná. Podle tabulky hodnot t pak budeme testovati na 5% stupni významnosti, budeme-li bráti ekvivalentní hodnotu t na řádce desetiprocentní, t. j. 0,10.

(3,11) Významnost rozdílů mezi rozptyly. Testujeme-li rozdíly mezi rozptyly výběrů velkých rozsahů, můžeme předpokládati, že rozdíly $\sigma(x_1v) - \sigma(x_2v)$ výběrových směrodatných odchylek mají normální rozdělení se směrodatnou odchylkou

$$\sqrt{\frac{\sigma^2(x)}{2r_1} + \frac{\sigma^2(x)}{2r_2}},$$

neboť podle (41) je rozptyl směrodatných odchylek výběrových z normálního základního souboru $\frac{\sigma^2(x)}{2r}$. Neznáme-li rozptyl základního souboru $\sigma^2(x)$, z něhož jsme výběry vzali, nahradíme jej pozorovanými $\sigma^2(x_1v)$ resp. $\sigma^2(x_2v)$. Ale v případě výběrů malých rozsahů jsou zase chyby vznikající touto aproximací závažné. Proto se uchylujeme k jinému postupu. Zavádíme nový index

$$z = \frac{1}{2} (\lg \sigma^2(x_1, v) - \lg \sigma^2(x_2, v)) = \lg \frac{\sigma(x_1, v)}{\sigma(x_2, v)}, \quad (59)$$

kde rozptyly $\sigma^2(x_1, v)$ a $\sigma^2(x_2, v)$ jsou počítány pomocí n_1 resp. n_2 stupňů volnosti. Rozdělení četnosti tohoto indexu z bylo odvozeno a má tvar

$$y = y_0 e^{nz} (n_1 e^{2z} + n_2)^{-\frac{1}{2}(n_1 + n_2)};$$

obsahuje jako proměnné jen z , n_1 , n_2 , a je tedy nezávislé na směrodatné odchylce základního souboru $\sigma(x)$. Poněvadž nepotřebuje předpokladů o přibližném vyjádření jejím, lze ho vhodně použít na případy malých výběrů. Abychom si postup osvětlili, všimněme si, že index z se může pohybovat mezi $+\infty$ a $-\infty$, má hodnoty záporné, když $\frac{\sigma(x_1, v)}{\sigma(x_2, v)} < 1$

a kladné pro $\frac{\sigma(x_1, v)}{\sigma(x_2, v)} > 1$. Tvar rozdělení četností je nesy-metrický kromě případu $n_1 = n_2$. Kladná část křivky pro $z = \frac{\sigma(x_1, v)}{\sigma(x_2, v)}$ je zřejmě táž jako záporná část křivky pro $z = \frac{\sigma(x_2, v)}{\sigma(x_1, v)}$. Postačí tedy pro jakoukoliv kombinaci stupňů volnosti pravděpodobnostní integrály jen pro kladné odchylky a ostatní dostaneme záměnou n_1 a n_2 . Zjednodušíme po-

stup, když se zápornými hodnotami z nepočítáme, ale vezmeme rozdíl logaritmů vždy tak, že je kladný a za n_1 zvolíme ten počet stupňů volnosti, s nímž jsme počítali větší rozptyl. Uvedeme několik hodnot z pro testování na hranici 5%; podrobnější tabulky najde čtenář v [1], str. 150.

Tabulka 6.

$n_1 \backslash n_2$	8	10	15	20	30	60	∞
8	0,618	0,561	0,486	0,447	0,409	0,370	0,331
12	0,595	0,535	0,453	0,412	0,369	0,326	0,280
24	0,568	0,504	0,414	0,367	0,318	0,265	0,209
∞	0,537	0,466	0,363	0,306	0,242	0,164	0,000

a hodnoty z na jednocentní hranici

Tabulka 7.

$n_1 \backslash n_2$	8	10	15	20	30	60	∞
8	0,898	0,810	0,694	0,636	0,577	0,519	0,460
12	0,867	0,774	0,650	0,586	0,522	0,457	0,391
24	0,832	0,732	0,596	0,525	0,452	0,375	0,291
∞	0,790	0,682	0,527	0,442	0,348	0,235	0,000

Pro vylíčený zjednodušený postup si musíme uvědomiti, že tyto hodnoty nejsou příslušnými hranicemi významnosti v našem smyslu. Podle našich úvah leží na 5% stupni významnosti hodnoty z, v nichž vztyčené pořadnice odetínají plochy, z nichž každá je 0,025 celé plochy. Proto leží pětiprocentní hranice významnosti někde mezi pětiprocentními a jednocentními hodnotami uvedených tabulek (6 a 7). Je-li počet stupňů volnosti n_1 a n_2 velký, nebo přibližně sobě rovný i když ne velký, blíží se rozdělení z normálnímu se směrodatnou odchylkou $\sqrt{\frac{1}{2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$ a potom, jak víme,

leží hodnoty z větší než dvojnásobek této směrodatné odchylky nad pětiprocentní hranicí významnosti. Je-li na př. $n_1 = n_2 = 20$, bude směrodatná odchylka 0,224 a z rovnající se dvojnásobku 0,448 je na pětiprocentní hranici čili mezi 0,382 a 0,545, jež dostáváme lineární interpolací z hořejších tabulek.

(3,11,1) Příklad. K objasnění uvedených metod testování při malých výběrech poslouží výsledky studia vlivu insulinu na králíky. Vliv byl měřen procentem svalového glykogenu (hodnota znaku) u 11 zvířat, která byla pod vlivem insulinu a u 10, která nebyla pod tímto vlivem takže slouží za kontrolu. Máme tak měření hodnoty znaku pro dva výběry: a) v kontrole, b) po insulinu.

Hodnota znaku
ve výběru

a)	b)
0,19	0,15
0,18	0,13
0,21	0,00
0,30	0,07
0,66	0,27
0,42	0,24
0,08	0,19
0,12	0,04
0,30	0,08
0,27	0,20
—	0,12

Průměr ve výběru

a) je $\bar{x}_1 = 0,273$,

ve výběru

b) je $\bar{x}_2 = 0,135$.

Vzhledem k tomu, že variace znaku od jednoho prvku výběru ke druhému jsou veliké, není rozdíl mezi průměry $d = 0,138$ velký. Přezkoumáme proto napřed jeho významnost.

Součet čtverců odchylek

$$\Sigma(x_1 - \bar{x}_1)^2 = 0,2530$$

$$\Sigma(x_2 - \bar{x}_2)^2 = 0,0715$$

Součet 0,3245,

takže

$$\sigma^2(x, v) = \frac{0,3245}{19} = 0,01708, \quad \sigma(x, v) = 0,1307.$$

Vzhledem k tomu, že $r_1 = 10$, $r_2 = 11$, bude

$$t = \frac{0,138}{0,1307} \sqrt{\frac{110}{21}} = 2,49.$$

Z tabulky 5 vidíme, že při $n = r_1 + r_2 - 2 = 19$ bude na pětiprocentní hranici významnosti $t = 2,09$, takže považujeme naši hodnotu rozdílu za významnou a vliv insulinu na svalový glykogen za skutečný.

Testujeme nyní ještě významnost rozdílu ve variabilitě. Rozptyly jsou $\sigma^2(x_1, v) = 0,02811$, $\sigma^2(x_2, v) = 0,00715$, takže $z = \frac{1}{2} \lg 3,93 = 0,684$.

Musíme nyní najít hodnoty z v tabulkách 6 a 7 pro počet stupňů volnosti $n_1 = 9$, $n_2 = 10$. Poněvadž v nich nejsou hledané hodnoty pro $n_1 = 9$ uvedeny, nýbrž jen pro $n_1 = 8$ a pak až pro $n_1 = 12$, musíme je určit interpolací. Vyjdeme při tom z té skutečnosti, že při témž n_2 jsou změny hodnoty z přibližně úměrné $\frac{1}{n_1}$, t. j. převrácené hodnotě počtu stupňů volnosti n_1 .

Tak najdeme napřed z tab. 6 pro $n_2 = 10$ hodnoty z ; pro $n_1 = 8$ je v první řádce $z_1 = 0,561$ a pro $n_1 = 12$ na druhé řádce $z'_1 = 0,535$, takže rozdíl $\Delta z = 0,561 - 0,535 = 0,026$.

Abychom provedli lineární interpolaci vzhledem ku $\frac{1}{n_1}$ na-

jdeme rozdíl $\frac{1}{n_1} - \frac{1}{n'_1} = \Delta = 0,1250 - 0,0833 = 0,0417$ a

dostaneme úměru $\left(\frac{1}{n_1} - \frac{1}{n'_1} \right) : \Delta = x : \Delta z$ čili

$$x = \frac{0,0139}{0,0417} \cdot 0,026 = 0,009,$$

z čehož plyne, že hledaná hodnota z na pětiprocentní hranici je $z_5 = 0,561 - 0,009 = 0,552$.

Abychom našli z_1 , vezmeme z tab. 7 pro $n_2 = 10$ hodnoty z $0,810 - 0,774 = 0,036$ a opět $x = \frac{0,0139}{0,0417} \cdot 0,036 = 0,012$ a $z_1 = 0,810 - 0,012 = 0,798$.

Naše hodnota z leží sice mezi pětiprocentní hranicí z_5 a jednaprocentní z_1 , ale není jisto, zda je na 5% stupni významnosti, takže můžeme jen říci, že tato data nás vedou k domněnce, že vlivem insulinu mají procenta glykogenu pravidelnější průběh, ale k rozhodnutí by bylo třeba ještě dalších pozorování.

(4) Reprezentativní metoda.

Teorie náhodného výběru tvoří základ t. zv. reprezentativního statistického šetření, jehož hlavním cílem je podati s nejmenším nákladem co nejvíce informací o základním souboru. Reprezentativní šetření je takové, které zkoumá část celého uvažovaného souboru, aby z ní odvodilo úsudky o celku. Částečných šetření se užívá ve velmi různých oborech, kde není možno nebo účelno provést vyšetření všech prvků souboru, z toho důvodu, že se uspoří peněz a práce nebo se zjednoduší a urychlí šetření i zpracování. Uplatňuje se tak reprezentativní metoda nejen ve vědách přírodních a v technické kontrole výroby průmyslové i zemědělské, nýbrž i v četných oborech hospodářské a sociální stránky života. Tak zjistíme mzdy textilních dělníků nebo kovodělníků reprezentativním šetřením podle určitého výběru a nevyšetřujeme mzdy všech textilních resp. kovodělníků. A jako v otázkách mzdových, tak postupujeme obdobně v otázkách cenových, složení rodin a rozvržení jejich vydání, sledování výroby v různých odvětvích co do množství a jakosti.

V mnohých případech se přeceňoval význam úplného čili vyčerpávajícího statistického šetření, neboť k účelům, pro něž bylo šetření provedeno, stačí hrubší, okrouhlá čísla, uvá-

žíme-li, že přesnost získaných čísel úplným šetřením bývá fikcí.

I kdyby tato čísla byla zcela přesná, jsou takovými jen v okamžiku, k němuž bylo provedeno šetření; čím jsme dále od tohoto okamžiku, tím jsou odchylky od skutečnosti větší. Poněvadž pak ve většině případů v praxi není třeba zcela přesných čísel a na statistických veličinách lpí větší či menší chyby, jak jsme viděli, doporučuje se vždy uvážiti, nemůžeme-li dostati čísla stejné přesnosti nesrovnatelně rychleji a levněji cestou reprezentativní metody čili v obchodní praxi obvyklou a osvědčenou metodou vzorku a má-li v daném případě smysl snažiti se o úplnou přesnost a úplnost materiálu.

Výklad α reprezentativní metodě je vhodno provésti se dvou hledisek. Jednak lze prováděti výběr metodou náhodného výběru, jednak podle principu uváženého čili záměrného výběru.

(4,1) Náhodný výběr s hlediska techniky výběrové. Metoda náhodného výběru může míti dvojí formu: a) neomezeného náhodného výběru, b) oblastního (stratifikovaného) náhodného výběru.

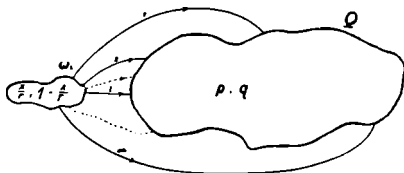
První forma je t. zv. klasická metoda; vybere se z určitého základního souboru jistý počet prvků tak, že pro každý prvek je stejná možnost, aby byl pojat do výběru. Vybere se tedy náhodně jako se konají na př. tahy z osudí. Při tom se postupuje buď tak, že se po tahu vrací prvek zpět, čímž vzniká způsob výběru z nekonečného základního souboru, nebo se vytažený prvek již nevrátí zpět; v tomto případě lze použití analogie s nekonečným základním souborem jen tehdy, má-li základní soubor veliký rozsah. Jsou-li podmínky náhodnosti a stejné možnosti splněny, není nebezpečí systematické chyby ve výběru.

V příslušné teorii náhodného výběru jsme rozeznávali dva hlavní případy:

a) odhadovali jsme četnost resp. pravděpodobnost p znaku, který se může u prvku vyskytnouti nebo nikoliv (obr. 6),

b) odhadovali jsme parametry rozdělení četností hodnot znaku. Určovali jsme pak pomocí směrodatné odchylky meze odpovídající jistým pravděpodobnostem.

První případ, který se týká alternativního znaku, byl řešen v podstatě rovnicemi (74), (75), (76) I. dílu. Pro druhý případ znaku kvantitativního je dáno řešení v předcházející kapi-



Obr. 6. Parametr a příslušná charakteristika.

tole. V obou případech jsme užili k provedení statistické indukce, metody nazývané nepříliš vhodně empirickou, která nám udává, jak docílíme z hodnoty charakteristiky na př.

$f_i = \frac{x}{r}$ nebo \bar{x}_i, σ_i pozorované v náhodném výběru ω_i (obr. 6)

nejlepší odhad příslušného parametru $p, \bar{x}, \sigma(x)$ v základním souboru Ω . Kromě toho jsme podali výklad druhé metody odhadu, t. zv. metody maximální věrohodnosti. Zbývá nám ještě všimnouti si blíže techniky náhodného výběru.

(4,2) Technika náhodného výběru. Výběr reprezentuje základní soubor stručně svými charakteristikami. Jejich hodnoty však podléhají náhodným odchylkám a vedle nich je tu nebezpečí odchylek systematických, které mají své zvláštní příčiny. Provádíme-li na př. statistiku mezd v kovodělném průmyslu pomocí několika závodů místo šetření ve všech závodech, může se státi, že určitá kategorie zaměstnanců má právě v těchto závodech z nějakého důvodu vyšší mzdu než v ostatních. Při zkoumání nákladů na spotřebu domácnosti dělnické či úřednické se užívá domácnostních účtů, jež vede určité množství rodin po dlouhé období; je jasno, že píše a

pečlivost, které je k této práci třeba, vyznačuje již jistou úroveň těchto rodin, takže k takové okolnosti jest přihlížeti při zevšeobecňování úsudků, které vyplynuly z poměrně malého reprezentativního souboru. V náhodném výběru mají býti tyto systematické odchylky vyloučeny. K technickému provedení náhodného výběru se tudíž obyčejně doporučuje nějaký mechanický postup, který by nezávisel na vlivu osoby, jež má výběr sestrojiti. Je-li na př. celý materiál dotazníkový nějakým způsobem seřazen, časově nebo místně, a má se z něho k určitému cíli poříditi pro rychlé odvození potřebného výsledku reprezentativní výběr náhodně, vezme se podle rozsahu na př. každý dvacátý dotazník, takže rozsah výběru tvoří potom přibližně dvacetinu celého souboru. Jsou-li prvky celého souboru zastupované dotazníky opatřeny pořadovými čísly, je možno opatřiti si miniaturu celého souboru tím, že všechna čísla napíšeme na jednotlivé lístky, jež dáme do osudí a pak z něho vytáhneme náhodně určitý počet lístků; do výběru pak zahrneme prvky souboru, které mají pořadová čísla, jež byla vytažena. Tohoto postupu nelze použítí je-li soubor, z něhož má býti proveden výběr, příliš rozsáhlý, neboť pak přesahuje tento postup často statisticky praktické možnosti; při rozdělení četností o dvaceti třídách s četnostmi mezi 1 až 60 je potřeba 1000 až 1200 lístků v osudí. Pro takové obsáhlejší případy byla sestavena t. zv. náhodná čísla Tippettova, jichž lze použítí k tvoření náhodných výběrů tak, že se vhodně přiřadí k prvkům základního souboru. Jsou to čísla vzata ze zpráv o jednom sčítání lidu a číslice jsou kombinovány po čtyřech. Byly pak provedeny pomoci statistických testů zkoušky, které potvrdily, že jsou to čísla vskutku prakticky náhodná. Budtež uvedena ukázkou tato čísla:

2693	1300	5356	7203	1396	1545	9524	4167
5624	3170	5911	7979	9792	3992	6641	2952
7691	6913	1089	3563	2762	3408	7483	2370
8776	4233	8126	6008	6107	1112	5246	0560
6446	8816	6111	7002	9025	1405	9143	2754

Abychom na př. vzali náhodný výběr rozsahu $r = 15$ ze základního souboru tab. 1 rozsahu 1001, očíslováme jeho prvky číslicemi od 1 do 1001 a najdeme nyní 15 čísel náhodně v mezích od 1 do 1001. Vezmeme tedy některé stránky čísel Tippettových a vybereme na nich prvních 15, která nejsou větší než 1001. Mezi našimi čísly nahoře by to bylo 560 a dále bychom na př. našli 423, 730, 918, 91, 17, 116, 708, 840, 638, 396, 29, 224, 717, 221. Číslování prvků v základním souboru jsme provedli na př. od nejnižší hodnoty znaku nahoru; pak můžeme náhodným prvkům přiřaditi hodnoty znaku, které budou 7, 7, 8, 10, 4, 3, 5, 8, 9, 8, 6, 3, 5, 8, 5, jak vidíme z tab. 1, sloupec (2). Postup je možno také jinak upravit, takže bychom se mohli přiblížit celkovému počtu Tippettových čísel tím, že přiřadíme každému jednotlivému prvku základního souboru 10 čísel. Potom by prvnímú intervalu odpovídala čísla 0000 až 0009, druhému 0010 až 0029 atd. Bližší výklad je uveden v [4]. Také bylo použito k sestrojování náhodných výběrů elektrických strojů třídících.

Máme-li vzíti jako vzorek na př. náhodný výběr r -kusů ze zásilky mnoha beden žárovek, zvolíme především náhodně některé bedny (můžeme pro každou hoditi mincí a bráti jen tu, pro niž padne rub) a v nich zase náhodně různé řady, z nichž žárovku k přezkoušení vezmeme. Kdybychom potřebovali náhodný výběr obyvatelů jedné z hlavních ulic města, můžeme vybrati některá čísla domů, jejichž obyvatelé budou tvořit žádaný výběr. Vyjdeme od některého libovolného domu a vezmeme třeba každý desátý. Nejsou-li tu nějaké zvláštní poměry v pravidelném seskupení zjišťovaných znaků, jako na př. příjem nebo počet členů rodiny, bude zvolená metoda nezávislá na vlastnostech souboru a výběr bude náhodný. Kdyby však v této ulici byl každý desátý dům rohový s velkým obchodem, nebyly by vyšetřované znaky nezávislé na metodě výběru, neboť se obchodní domy vyskytují s touž periodou, jaká byla zvolena pro výběr. Bývá ovšem často obtížno posouditi, zda při provedení výběru nebyl nějaký pramen systematických odchylek, který nemohl býti postřehnout.

Někdy vyžaduje účel šetření, aby byl proveden výběr oblastní. Rozdělí se tedy nejprve celý vyšetřovaný soubor podle určitého znaku na oblasti, a z nich se pak několik prvků vybere náhodně. Tak na př. při shora zmíněném šetření nákladů na spotřebu domácností dělnických na základě domácích účtů je účelno rozdělit všechny dělnické domácnosti podle zaměstnání přednosty domácnosti a sestaviti výběr tak, aby hlavní zaměstnání v něm byla zastoupena způsobem odpovídajícím přibližně struktuře obyvatelstva. V oblasti hlavních druhů zaměstnání vyberou se pak pokud možno náhodně zkoumané domácnosti. Pro vytváření oblastí (strata) může býti směrodatno i více znaků (město, venkov, počet dětí a pod.). Je ještě řada jiných způsobů konstrukce náhodného výběru. Pro vhodnou volbu techniky výběrové v určitém případě musí míti statistik dostatečné znalosti věcné v oboru, do něhož zkoumaný soubor patří a také dosti šťastné intuice.

(4,3) Splnění podmínek náhodného výběru. Praktický význam teorie náhodného výběru je v tom, že umožňuje měřit objektivně chyby odhadu a významnost hodnot zjištěných z náhodného výběru. Můžeme-li pak vzítí několik výběrů z jednoho základního souboru, lze zkoumati, zda rozdělení charakteristik je takové, jak je udává teorie. Odchyluje-li se významně, máme důvod k tomu, abychom zkoumali zvolenou výběrovou techniku a hledali, proč zavádí systematické odchylky. Tento postup předpokládá, že známe rozdělení četností základního souboru, neboť jinak je musíme jen odhadovat podle výběru a pak ovšem nemůžeme toho odhadu užítí ke kritice metody výběrové bez dalšího bližšího vyšetřování. Systematická odchylka od podmínek nutných k sestrojení náhodného výběru musí býti vyloučena dříve než je možno aplikovati výsledky, které plynou z teprve náhodných odchylek výběrových. Aby byly při praktickém provádění sestrojeny podmínky náhodného výběru z pramenných dat, je třeba odvozovati v určitých šetřeních zvláštní schemata k získání náhodného výběru, což závisí na povaze

oboru, v němž se šetření koná. Každý obor vědní nebo druh výroby či obchodu má své problémy při náhodném výběru.

Někde jsou vydána úřední ustanovení pro braní vzorků. Taková jsou na př. vyhlášena výsadní obilní společností o braní vzorků a zkoumání pšenice bohaté na lepek. Vzorky bere osoba úředně k tomu stanovená, která se nejprve přesvědčí, že určité vlastnosti pšenice jsou stejnoměrně vyrovnány v celém množství, z něhož se mají bráti vzorky. Potom se vezme při partiích do deseti pytlů z každého pytle vzorek po 150 g, při partiích do 20 pytlů nejméně 10 vzorků z různých nahodile vzatých pytlů, při partiích do 50 pytlů nejméně 15 vzorků z 15 různých pytlů po 100 g atd. Takto vzaté vzorky se smíchají a stejnoměrně rozprostřou. Celkové množství se potom rozdělí na pět stejných dílů, z nichž se 3 znovu spolu promíchají a tato směs tvoří konečný vzorek, který se rozdělí na tři části nejméně po 250 g, z nichž se každá zapečetí a jedna odevzdá příslušnému zkušebnímu ústavu, druhou dostane prodávající a třetí kupující. Pytle se po odebrání vzorků zaplombují předepsaným způsobem.

(4,4,1) Určení rozsahu výběru při znaku alternativním. Středem praktického provedení reprezentativní metody je určení velikosti výběru čili počtu prvků, které mají býti vzaty z celkového souboru a podrobeny vlastnímu zkoumání. Na této veličině závisí reprezentativní síla výběru a rozhoduje o tom, jak je výsledek dobrý a pravdivý. Viděli jsme, že neusuzujeme z hodnoty relativní četnosti znaku v jednom výběru nebo z průměru jednoho výběru na příslušnou hodnotu v základním souboru, nýbrž vždy jen na jistý obor, v němž může hledaný parametr ležet. Tento obor závisí na stupni pravděpodobnosti, s níž chceme počítat a na počtu prvků zahrnutých do výběru.

Když jsme si předepsali stupeň pravděpodobnosti a rozhodli se pro určité hranice odchylek, pak můžeme rozhodnouti otázku, kolik prvků musí býti pojato do výběru, abychom mohli s určitým stupněm pravděpodobnosti předpokládati, že na př. výběrem zjištěná relativní četnost f bude se

odchylovat od p nahoru či dolů nejvýše o určitou napřed stanovenou veličinu.

Je-li tato napřed stanovená odchylka $z_0 = |f - p|$ a stupeň pravděpodobnosti dán číslem 0,966, vidíme, že odpovídá 1,5násobnému modulu (tab. 8).

Tabulka 8.

t	$\alpha(t)$	$t = j\sqrt{2}$	$\alpha(t)$
1	0,683	$\sqrt{2}$	0,843
2	0,955	$1,5\sqrt{2}$	0,966
3	0,997	$2\sqrt{2}$	0,995

Vyjádříme-li odchylku z_0 v jednotce modul, máme známý vztah

$$\gamma = z_0 \sqrt{\frac{r \frac{N-1}{N-r}}{2p(1-p)}}, \text{ takže } z_0 = \gamma \sqrt{\frac{2p(1-p)}{r \frac{N-1}{N-r}}}$$

v případě, že se jedná o výběr ze základního souboru konečného a vyňatý prvek se nevrací zpět. Můžeme-li prakticky považovati rozsah základního souboru za nekonečný, je výraz modulu

$$\sqrt{\frac{2p(1-p)}{r}}$$

Je vhodno stanovit odchylku z_0 jako zlomek hodnoty p , tedy

$$\delta = \frac{z_0}{p}; \text{ je tedy } z_0 = \delta p = \gamma \sqrt{\frac{2p(1-p)}{r \frac{N-1}{N-r}}}$$

a pro známé hodnoty δ, γ můžeme odtud zjistit jak velký musí být r , čili jaký zlomek z N prvků musí být vzat do

výběru. Jednoduchou úpravou dostáváme t. zv. výběrovou rovnici

$$\frac{r}{N} = \frac{1}{1 + \frac{(N-1)p\delta^2}{2\gamma^2(1-p)}} \quad (60)$$

V našem případě dosadíme za $\gamma = 1,5$ a dostaneme jaký zlomek celého souboru nutno vzít, abychom mohli říci s pravděpodobností 0,966, že relativní četnost bude v mezích $p \pm z_0$ čili v mezích $p \pm \delta p$.

Z výběrové rovnice je zřejmo, že při stálém podílu $\frac{r}{N}$, který se také nazývá výběrovým koeficientem, se δ zmenšuje, jestliže N a p roste. Z toho tedy plyne, že pro znak s větší relativní četností budeme očekávat menší zlomek δ než pro znak s menší relativní četností při stejném poměru rozsahu výběrového r ku N .

Obráceně, má-li být δ stejné, musíme vzít pro znak s menší relativní četností větší rozsah výběrový.

Abychom si učinili představu, jakého rozsahu výběrového ze základního souboru rozsahu $N = 100\,000$ je třeba při pravděpodobnosti 0,966 pro různé hodnoty p , aby relativní četnosti f byly v mezích $p \pm \delta p$, sestavíme si několik hodnot do tabulky 9, z níž vidíme na př., že musíme vzít výběr rozsahu 448 prvků ze základního souboru rozsahu 100 000, abychom mohli očekávat s předpokládanou pravděpodobností 0,966, že relativní četnost nebude kolísati víc než 10% kolem $p = 0,5$, že tedy bude mezi 0,45 a 0,55. Když pak dělíme počet všech prvků základního souboru číslem v tabulce, dostaneme kolikátý prvek musí být vzat do výběru; v našem případě je to tedy $\frac{N}{r} = 100\,000 : 448 = 223$. Hlavní

potíž při praktickém užívání výběrové rovnice bývá v tom, že p neznáme. Výběrový koeficient $\frac{r}{N}$ závisí na čtyřech veličinách N , p , δ , γ . Skutečně danou veličinou je N , kdežto pro ostatní tři musíme zvolit určitý předpoklad, abychom

mohli stanovit rozsah výběru. Mezi nimi pak je vnitřní vázanost, neboť p je rozhodující pro velikost modulu a teprve za předpokladu γ -násobku tohoto modulu jako absolutní odchylky může být určeno δ . Je tudíž zřejmo, že rozsah výběru závisí na předpokladu, který učiníme o velikosti p . Vliv p na velikost výběrového koeficientu vynikne, uvážíme-li jeho extrémní hodnoty. Je-li $p = 1$, čili soubor je jednotný, takže všechny prvky mají uvažovaný znak, je z výběrové rovnice zřejmo, že je třeba výběru krajně nepatrného rozsahu vlastně žádného $\frac{r}{N} \rightarrow 0$, ježto každý prvek repre-

sentuje v tomto smyslu celý soubor. Pro $p = 0$ je $\frac{r}{N} = 1$,

čili za předpokladu, že se v celém základním souboru nevyskytuje zkoumaný znak, je třeba velmi rozsáhlého výběru, ba celého základního souboru. Všechny prvky nebo největší část základního souboru musí přijít do výběru, poněvadž jen tak dojde několik málo prvků se zkoumaným znakem k reprezentaci. Z toho vidíme, že veličina p a r se mění v obráceném poměru čili s poklesem p stoupá výběrovou rovnicí požadovaný rozsah výběru a obráceně. Význam této vlastnosti výběrové rovnice spočívá v tom, že statistik může určit potřebnou velikost výběru svým předpokladem o veličině p . Z předcházejících vývodů pak vyplývá, že musíme uvažovati tři případy.

Tabulka 9.

$p \backslash 100\delta$	50	25	10	5
0,01	1.754	6.667	33.333	—
0,10	162	645	4.000	14.286
0,25	54	216	1.333	5.263
0,5	18	72	448	1.785
0,6	12	48	299	1.190
0,7	8	31	193	769
0,8	5	18	112	448
0,9	2	8	50	200

a) Výběr je s hlediska jednoho alternativního znaku reprezentativní svým rozsahem; jestliže použitý předpoklad o p vzatý za základ výpočtu výběrového koeficientu se rovná parametru v základním souboru nebo je menší.

b) Slouží-li výběr k studiu několika alternativních znaků, jejichž parametry jsou p_i , je svým rozsahem reprezentativní, byl-li k určení rozsahu výběru zvolen předpoklad o nejmenším p_i .

c) Při studiu kombinací několika znaků je třeba přihlížeti k minimálním relativním četnostem, které kombinacemi vznikají a jež vedou k velikým rozsahům výběru. V tomto bodě se často hřeší v praxi statistických šetření.

(4,4,2) Určení rozsahu výběru při znaku kvantitativním. Pro znak kvantitativní, nabývající hodnot x_1, x_2, \dots, x_l obvykle považujeme výběr za reprezentativní, můžeme-li se spolehnouti s pravděpodobností předem určenou, že odchylky výběrových průměrů od průměru \bar{x} v základním souboru budou v určitých mezích, daných nějakým násobkem směrodatné odchylky nebo modulu. K výběrové rovnici dospějeme,

zvolíme-li opět $\delta = \frac{z_0}{\bar{x}}$; poněvadž směrodatná odchylka vý-

běrových průměrů je podle (13) $\sigma_P = \frac{\sigma(x)}{\sqrt{r \frac{N-1}{N-r}}}$, pak má-li

se odchylka z_0 rovnat nejvýše γ -násobnému modulu, budeme mít opět vztah

$$\delta \bar{x} = \gamma \sigma(x) \sqrt{2 \frac{N-r}{r(N-1)}},$$

z něhož dostáváme pro výběrový koeficient

$$\frac{r}{N} = \frac{1}{1 + \frac{\bar{x}^2}{\gamma^2 \sigma^2(x)} \frac{N-1}{2} \delta^2}. \quad (61)$$

Můžeme-li považovati prakticky rozsah základního souboru

za nekonečný, je modul vyjádřen výrazem $\frac{\sigma(x)\sqrt{2}}{\sqrt{r}}$, takže pak dostaneme přímo výběrový rozsah

$$r = \frac{2\gamma^2\sigma^2(x)}{\delta^2\bar{x}^2}$$

Vidíme, že výběrový koeficient závisí na veličinách \bar{x} , $\sigma(x)$, γ , δ , N , z nichž první dvě jsou parametry základního souboru, které musíme nějak odhadnouti, což je možno jen podle zkušeností z oboru, jehož se šetření týká, nebo podle dříve provedených analogických šetření.

Poznámka: Dobrým cvičením jest provedení úvah a odvození výběrových rovnic pro odchylky rovnající se nějakému násobku směrodatné odchylky.

Výběrová rovnice nám dala odhad rozsahu pro výběr, který můžeme považovati za reprezentativní po stránce kvantitativní. Reprezentativní síla výběru se však jeví ještě stránkou kvalitativní, což značí, že výběr má představovat pro všechny zkoumané znaky přesné a tedy co nejvěrnější zobrazení základního souboru, který má určitý stupeň homogenity, velkou či malou rozmanitost jednotlivých znaků, která závisí na tom, kolika hodnot může dotýčný znak nabývat a určitý rozsah. Pod těmito třemi hledisky se nám jeví povaha souboru, kterou musíme vždy podrobně uvažovati, abychom zjistili, zda odpovídá povaze souborů, které byly základem odvozené teorie. Soubor s malým rozsahem může poskytovat při úplné homogenitě a malé rozmanitosti přibližně reprezentativní četnost malého počtu znaků ve výběru. Ovšem v rozsáhlém souboru se projeví jednotlivé znaky reprezentativně s jistotou mnohem větší. Rozmanitost přizpůsobujeme tím, že se vzdáme jistých výsledků a zkoumáme jen ty znaky, které vystupují ve velkých četnostech, pro něž tudíž existuje velká pravděpodobnost, že se objeví ve výběru. Do jaké míry je takové přizpůsobení možné, závisí na možnosti omezení po případě

pozměnění účelu šetření. Ve všech případech, v nichž nelze připustiti toto zjednodušení poznatků, musí býti provedeno místo reprezentativního šetření vyčerpávající šetření.

Oblastní výběr vzniká postupem, který rozdělí základní soubor rozsahu N na několik menších souborů rozsahu N_i (obecně různého rozsahu) t. zv. oblastí (strata). Z každé oblasti pak vezmeme technikou náhodného výběru určitý počet prvků r_i zpravidla tak, že poměr $r_i : N_i$ je pro všechna i stejný, což však není nutné. Z těchto částečných výběrů oblastních dostaneme shrnutím celkový výběr. Tento postup se nazývá metodou stratifikovaného čili oblastního výběru. Také můžeme užívatí typu oblastního výběru skupinového [1].

(4,5) Záměrný výběr. Podle principu záměrného, také uváženého čili systematického výběru se rozdělí základní soubor na skupiny známých rozsahů a základních parametrů (průměry, směrodatné odchylky). Z těchto skupin se snažíme vybrat záměrně takové, které dohromady poskytnou pro určitou charakteristiku (průměr) týž výsledek jako základní soubor [1]. Matematické podmínky, které jsou podkladem této metody, jsou dosti omezující a teoretická i praktická zkoumání naznačují, že obecně nelze přikládati získaným výsledkům tolik spolehlivosti, jako získaným pomocí náhodného výběru. Výsledky získané reprezentativní metodou se přenášejí na neznámý základní soubor s určitou rezervou. Zvláště v oboru hospodářských statistik se vyskytují často reprezentativní šetření, která nejsou vždy provedena průhlednými metodami záměrného výběru a výsledky se bez důkazu vydávají za reprezentativní. Metody náhodného výběru se v tomto oboru neuzívá ještě takovou měrou, jak by to dovoľoval nynější stav teorie. Větší pronikání lze pozorovati ve statistice výroby zemědělské a průmyslové.