

# Jak vytváří statistika obrazy světa a života. I. díl

---

Část III. [odst. 5,1-5,4, 6,1-6,9, 7,1-7,4, 8,1-8,5, 9,1-9,2]

In: Jaroslav Janko (author): Jak vytváří statistika obrazy světa a života. I. díl. (Czech). Praha: Jednota českých matematiků a fyziků, 1942. pp. 64–138.

Persistent URL: <http://dml.cz/dmlcz/403051>

## **Terms of use:**

© 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 III.

(5,1) Teorie náhodného výběru. (Znak alternativní.)  
Hodnota relativní četnosti v základním souboru  
— pravděpodobnost.

Jakmile přecházíme od popisných úkolů k bližšímu vysvětlování pozorovaných jevů hromadných, opíráme se o pojem pravděpodobnosti a věty odvozené počtem pravděpodobnosti. Při t. zv. statistické definici pravděpodobnosti vycházíme od posloupnosti jevů. Procházíme-li zápisy v matrice nějakého většího města, které jsou vedeny časově za sebou třeba po dvacet let a zaznamenáváme porody podle znaku pohlaví, takže označujeme chlapce  $c$ , děvčata  $d$ , dostaneme posloupnost, jejíž členy opatříme pořadovými čísly (v druhém řádku)

$c$	$d$	$c$	$d$	$d$	$c$	$c$	$c$	$d$	$c$	$d$	$d$	$d$	$c$	$c$	$d$	$c$	...
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	...
1	1	2	2	2	3	4	5	5	6	6	6	6	7	8	8	9	...

Abychom nabyli přehledu o četnosti narozených chlapců v určitém úseku posloupnosti, sčítáme v něm písmena  $c$ . Jedná-li se o úsek od začátku až do některého pořadového čísla  $i$ , napíšeme pod něj v třetím řádku zjištěný počet písmen  $c$ , který označujeme jako absolutní četnost  $n_i$ , takže v naší posloupnosti jsou četnosti chlapců postupně

$$n_1 = 1, n_2 = 1, n_3 = 2, \dots, n_7 = 4, n_8 = 5, \dots, n_{16} = 8, n_{17} = 9, \dots$$

Vidíme, že k původní posloupnosti porodů náleží posloupnost absolutních četností znaku  $c$  a tudíž také posloupnost

relativních četností  $f_i = \frac{n_i}{i}$ , která má pro naše další účele

zvláštní význam, neboť jsme viděli, že základní formou podávání výsledků statistického šetření je relativní četnost

znaku v pozorovaném souboru. Je zřejmo, že absolutní četnosti jsou mezi 0 a  $i$ , takže platí nerovnosti  $0 \leq n_i \leq i$  a pro relativní četnosti tedy  $0 \leq f_i \leq 1$ .

Posloupnost relativních četností je v našem případě

$$\frac{1}{1}, \frac{1}{2}, \frac{2}{3}, \frac{2}{4}, \frac{3}{5}, \frac{3}{6}, \frac{4}{7}, \frac{5}{8}, \frac{5}{9}, \frac{6}{10}, \frac{6}{11}, \dots$$

Úsek posloupnosti, který jsme uvedli, je zcela nepatrný. Kdybychom sledovali v dlouhém úseku pěti let, šesti, sedmi, osmi, ... let vývoj čísel  $f_i$ , pozorovali bychom, že se stále blíží určitému číslu na př. 0,51, od něhož se liší na desetinných místech vždy vzdálenějších.

Číslo 0,51 dostáváme pro celý soubor, který představuje posloupnost prvků za celých dvacet let. Tento soubor je vyššího řádu než soubory částečné, jež tvoří posloupnosti pozorované za kratší časové úseky. Považujeme jej za soubor základní.

Relativní četnost v základním souboru nazýváme statistickou pravděpodobností, kterou budeme označovat  $p$ .

Základní soubor si budeme představovat jako soubor, jehož prvky jsou dobře promíchány, což znamená, že ve všech jeho částech je relativní četnost pozorovaného znaku přibližně táž. Budeme nejprve předpokládat, že známe relativní četnost znaku  $c$  v základním souboru čili pravděpodobnost  $p$ . Rozsah základního souboru je  $N$ .

Budeme bráti náhodně ze základního souboru výběry o  $r$  prvcích, tak jako bereme kuličky z osudí. Takové soubory budeme nazývat náhodné výběry. Na vyňatém prvku zjistíme, má-li pozorovaný znak, a zase jej vrátíme do základního souboru, takže se v něm  $p$  nemění. V náhodných výběrech rozsahu  $r$  prvků se bude vyskytovat různý počet prvků s pozorovaným znakem  $c$ , který označíme  $x$ . Budou výběry, v nichž nebude ani jeden prvek se znakem  $c$ , tedy  $x = 0$ , v jiných bude  $x = 1, 2 \dots a$  v některých  $x = r$ . Dělíme-li tento počet prvků se znakem  $c$  rozsahem výběru  $r$ , dostaneme relativní četnost  $\frac{x}{r}$ . Všech možných výběrů

a tedy také hodnot  $x$  bude  $\binom{N}{r}$ , neboť tolika způsoby lze kombinovat  $N$  prvků po  $r$ ; představíme si, že tyto hodnoty tvoří nový soubor. Chceme především stanovit, jaké jsou v něm relativní četnosti jednotlivých hodnot  $x = 0, 1, 2, \dots, \dots, r$ , tedy konkrétních kombinací s  $x$  prvky znaku  $c$ .

Z počtu pravděpodobnosti známe pravděpodobnost, že nastane v souboru  $r$  pokusů právě  $x$ -krát jev, jehož pravděpodobnost je  $p$ . Je dána t. zv. Newtonovou formulí

$$P_r(x) = \binom{r}{x} p^x q^{r-x}, \quad (36)$$

kde  $q = 1 - p$ .

Bude tedy celé rozdělení četností určeno všemi členy  $P_r(x)$ , t. j. pro  $x = 0, 1, 2, \dots, r$ , což jsou jak známo členy binomického rozvoje

$$(q + p)^r = q^r + r p q^{r-1} + \binom{r}{2} p^2 q^{r-2} + \dots + \\ + \binom{r}{x} p^x q^{r-x} + \dots + p^r, \quad (37)$$

jak se odvozuje v počtu pravděpodobnosti [10], [11] pro pravděpodobnosti opakovaných jevů. Jednotlivé členy mají charakter statistických pravděpodobností, jak je patrné z odvození, neboť jsme je nedostali jako výsledky skutečně provedených výběrů.

**(5,2) Binomické rozdělení četností, jeho průměr a rozptyl.** Rozdělení četností, jehož třídní četnosti jsou úměrné členům tohoto rozvoje, se také nazývá rozdělení Bernoulliho.

Jeho důležitost není jenom v tom, že udává nejpravděpodobnější rozdělení výběrů z osudí, nýbrž vystihuje typ rozdělení relativních četností, které dostáváme při nejjednodušších operacích náhodného výběru ve statistice. Tak považuje na př. biolog rozvoj (37) za teoretické rozdělení

relativních četností chlapců v náhodných výběrech o rozsahu  $r$  porodů. Pojistný technik na př. považuje rozvoj (37) za teoretické rozdělení ročních měr úmrtnosti v náhodných výběrech rozsahu  $r$  mužů téhož věku, třeba 25 roků. Při tom nutno zdůrazniti, že tyto výběry jsou brány stále za týchž podmínek, zde ze souboru mužů stále stejného složení vzhledem k znakům, které mohou míti vliv na úmrtnost, tedy vzhledem k povolání, zdravotnímu stavu a pod. Předpoklad stále stejných podmínek čili stálého  $p$  je podkladem základním při odvozování Bernoulliova rozdělení; jinými slovy provádění jednoduchého náhodného výběru předpokládá, že základní pravděpodobnost  $p$  výskytu znaku zůstává konstantní od výběru k výběru, v němž jednotlivé prvky jsou vzájemně nezávislé, t. j. na zahrnutí prvku do výběru nemá významného vlivu zahrnutí prvku předcházejícího.

Nejpravděpodobnější počet  $x'$  prvků, se znakem  $c$  ve výběru rozsahu  $r$  najdeme, utvoříme-li poměr obecného členu rozvoje, k předcházejícímu a pak k následujícímu; tyto dva poměry budou rovny nebo větší než 1.

$$\frac{r-x+1}{x} \frac{p}{q} \geq 1 \quad \text{čili} \quad x \leq pr + p$$

a stejně druhý poměr

$$\frac{x+1}{r-x} \frac{q}{p} \geq 1 \quad \text{čili} \quad x \geq pr - q.$$

Z toho plyne, že pro celá čísla  $x$  je největší hodnota určena nerovnostmi

$$pr - q \leq x' \leq pr + p$$

nebo vzhledem ku  $p + q = 1$

$$pr + p - 1 \leq x' \leq pr + p,$$

takže zanedbáme-li pravý zlomek, nejčetnější hodnota počtu prvků se znakem  $c$  je  $pr$ . Jsou-li  $pr - q$  a  $pr + p$  čísla celá, existují dva stejné největší členy rozvoje. [Ukažte, že první dva členy rozvoje  $(\frac{p}{q} + \frac{q}{p})^r$  jsou stejné.]

Odvodíme si pro toto rozdělení četností první dvě charakteristiky, jimž budeme říkati parametry, ježto jsou to hodnoty v základním souboru, kde relativní četnost znaku pozorovaného je rovna pravděpodobnosti  $p$ . Hodnotu průměru v základním souboru označíme  $\mathfrak{E}(x)$  a dostaneme podle definice, je-li  $x$  hodnota znaku

$$\begin{aligned}\mathfrak{E}(x) &= \sum_{x=0}^r \binom{r}{x} p^x q^{r-x} \cdot x = \\ &= \sum_{x=0}^r \frac{r!}{x! (r-x)!} p^x q^{r-x} \cdot x = \sum_{x=0}^r \frac{r!}{(x-1)! (r-x)!} p^x q^{r-x} = \\ &= rp \sum_{x=1}^r \frac{(r-1)!}{(x-1)! (r-x)!} p^{x-1} q^{r-x} = rp,\end{aligned}$$

neboť

$$\sum_{x=1}^r \binom{r-1}{x-1} p^{x-1} q^{r-x} = (p+q)^{r-1} = 1.$$

Dostáváme tudíž výsledek

$$\mathfrak{E}(x) = rp. \quad (38)$$

Kdyby hodnoty znaku byly  $\frac{x}{r} = f$ , je patrné, že bychom dostali průměr  $\mathfrak{E}(f) = p$ . Abychom odvodili v tomto rozdělení četností hodnotu rozptylu  $\sigma^2(x)$ , utvoříme součet čtverců odchylek od průměru  $\xi = x - rp$ , takže podle definice bude

$$\begin{aligned}\sigma^2(x) &= \sum_{x=0}^r \binom{r}{x} p^x q^{r-x} (x - rp)^2 = \\ &= \sum_{x=0}^r \binom{r}{x} p^x q^{r-x} (x^2 - 2xrp + r^2 p^2).\end{aligned} \quad (39)$$

Místo  $x^2$  uijeme identického výrazu  $x^2 = x + x(x-1)$ , takže první člen můžeme psáti

$$\sum_{x=0}^r \frac{r!}{x!(r-x)!} p^x q^{r-x} + r(r-1) p^2 \sum_{x=2}^r \frac{(r-2)!}{(x-2)!(r-x)!} \times \\ \times p^{x-2} q^{r-x} = rp + r(r-1) p^2,$$

druhý člen je

$$2rp \sum_{x=0}^r \binom{r}{x} p^x q^{r-x} \cdot x = 2r^2 p^2$$

a třetí člen

$$r^2 p^2 \sum_{x=0}^r \binom{r}{x} p^x q^{r-x} = r^2 p^2$$

z čehož plyne, že rozptyl

$$\begin{aligned} \sigma^2(x) &= rp + r(r-1) p^2 - 2r^2 p^2 + r^2 p^2 = \\ &= rp - rp^2 = rp(1-p) = rpq. \end{aligned} \quad (40)$$

Uvažujeme-li odchylky relativní četnosti znaku, od pravděpodobnosti výskytu znaku  $\frac{x}{r} - p$  dostaneme průměr čtverců těchto odchylek, dělíme-li výraz (39) čtvercem  $r^2$ , takže příslušný rozptyl je dán zlomkem  $\sigma^2(f) = \frac{pq}{r}$ . Je tudíž patrné, že rozptyl absolutních četností roste s rostoucím rozsahem  $r$  výběru, kdežto rozptyl relativních četností klesá s rostoucím rozsahem  $r$  výběru.

Podotkneme ještě, že bychom dostali obdobně rozptyl pro případ, t. zv. hypergeometrického rozdělení četností (str. 90), jež vystihuje braní výběrů ze základního souboru tím způsobem, že se prvky nevrací zpět. Rozptyl pak je dán výrazy

$$rpq \left(1 - \frac{r}{N}\right) \text{ resp. } pq \left(\frac{1}{r} - \frac{1}{N}\right).$$

Tento případ přechází v binomický, je-li rozsah základního souboru velmi velký  $N \rightarrow \infty$ . Také tyto výrazy pro rozptyl přecházejí pak na (40) resp. (40').

Budeme dále zkoumati, zda není možno udati pro odchylky  $\frac{x}{r} - p$ , tedy odchylky relativních četností znaku  $c$  ve výběrech  $\frac{x}{r}$  od relativní četnosti v základním souboru  $p$ , takové hranice, v nichž bude většina všech možných výsledků. K tomu cflí si napřed odvodíme důležitou větu.

**(5,3) Věta Bienaymé-Čebyševova.** Představme si, že máme napozorováno množství hodnot statistické proměnné  $x_1, x_2, \dots, x_l$  s relativními četnostmi resp.  $\nu_1, \nu_2, \dots, \nu_l$ , takže  $\nu_1 + \nu_2 + \dots + \nu_l = 1$ . Je-li jejich průměr  $\bar{x}$  a označíme zase odchylky  $x_i - \bar{x} = \xi_i$ , bude rozptyl

$$\sigma_x^2 = \nu_1 \xi_1^2 + \nu_2 \xi_2^2 + \dots + \nu_l \xi_l^2.$$

Rozdělme nyní odchylky  $\xi_i$  na takové, které nedosahují numericky určitého násobku směrodatné odchylky  $\tau\sigma_x$ , přičemž  $\tau > 1$  a na ostatní  $|\xi_i| \geq \tau\sigma_x$ . Relativní četnost prvních odchylek označíme  $P_\tau$ , takže relativní četnost ostatních bude  $1 - P_\tau$ . Můžeme pak psát rovnici pro rozptyl

$$\sigma_x^2 = \sum_{i=1}^l \xi_i^2 \nu_i + \sum_{j=l'+1}^l \xi_j^2 \nu_j.$$

Kde se první součet vztahuje na všechna  $\xi_i$ , která nedosahují  $\tau\sigma_x$  a druhý součet na všechna  $\xi_j$ , která se mu rovnají a převyšují. Poněvadž máme všechny sčítance obou součtů kladné nebo rovny nule, je

$$\sigma_x^2 \geq \sum_{j=l'+1}^l \xi_j^2 \nu_j.$$

Platí tudíž zřejmě nerovnost

$$\sigma_x^2 > \sum_{j=l'+1}^l \tau^2 \sigma_x^2 \nu_j.$$



Vzhledem k tomu, že jsme označili

$$\sum_{i=r+1}^l \nu_i = 1 - P_r$$

je také

$$\sigma_r^2 > \tau^2 \sigma_x^2 (1 - P_r) \text{ čili } \frac{1}{\tau^2} > 1 - P_r$$

a konečně

$$P_r > 1 - \frac{1}{\tau^2}, \quad (41)$$

což znamená, že relativní četnost prvků, jejichž hodnota znaku se bude odchylovati od průměru méně než o  $\tau\sigma_x$  je větší než  $1 - \frac{1}{\tau^2}$ . Tato věta se nazývá kriteriem Bienaymé-Čebyševovým.

Všimneme si ještě, že pravděpodobnosti odchylek podle věty Bienaymé-Čebyševovy mají povahu obecnou, která nezávisí nijak na tvaru rozdělení četností. Za to však jsou tyto pravděpodobnosti určeny v úzkých mezích často nepostačujících, neboť je tu udána dolní mez stanovením, že pravděpodobnost odchylky v mezích  $\tau$ -násobné směrodatné

odchylky je větší než  $1 - \frac{1}{\tau^2}$ . Vzniká zase otázka, jak

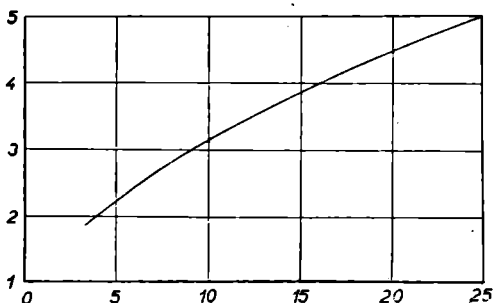
blízko je tato dolní mez skutečné hodnotě pravděpodobnosti. Tato otázka má praktický význam, neboť je-li tato mez značně níže než skutečná pravděpodobnost, musíme provésti k dosažení uspokojivého výsledku zbytečně mnohem více pozorování, než kdybychom znali skutečnou pravděpodobnost. Podle věty Bienaymé-Čebyševovy leží víc než  $1 - \frac{1}{\tau^2}$

z celkového počtu  $r$  prvků souboru v mezích  $\bar{x} \pm \tau\sigma_x$  (kde ovšem  $\tau \geq 1$ ) a tato věta platí pro jakoukoliv množinu konečných čísel, bez ohledu na to, jak byla získána. Pro několik hodnot  $\tau$  si sestavíme přehled:

$\tau$	1	$1\frac{1}{2}$	2	3	4
$1 - \frac{1}{\tau^2}$	0	0,56	0,750	0,889	0,937

Známe-li tedy  $\bar{x}$  a  $\sigma_x$ , můžeme hned říci, že víc než 75% čísel leží v intervalu  $\bar{x} \pm 2\sigma_x$ , čili méně než 25% se liší od  $\bar{x}$  o více než  $2\sigma_x$  atd.

Také vyplývá z věty B.-Č., že při rozsahu  $r = 4$  budou všechny prvky souboru v mezích  $\bar{x} \pm 2\sigma_x$ , neboť jich tam bude víc než  $1 - \frac{1}{4}$ , tedy ze čtyř více než tři čtvrtiny.



Obr. 11. Hodnoty  $\tau$ , pro něž všechny prvky jsou v intervalu  $x \pm \tau\sigma_x$ .

Podobně pro  $r = 10$  vidíme, že budou všechny prvky v intervalu  $\bar{x} \pm 3,16\sigma_x$ , neboť jich tam bude víc než  $1 - \frac{1}{10}$  a pod. Můžeme si graficky znázorniti obory, v nichž jsou podle věty B.-Č. obsaženy všechny prvky souboru; budou vyznačeny křivkou  $\tau^2 = r$  (obr. 11).

**(5,4) Teorém Bernoulliův.** Budeme nyní s hlediska kritéria B.-Č. uvažovati zmíněnou již úlohu, která je jedním z uhelných kamenů moderní statistiky, totiž najíti pravděpodobnost, že odchylka relativních četností  $\left| \frac{x}{r} - p \right|$  bude menší než libo-

volné kladné číslo  $\varepsilon$ . Zvolíme tedy  $\varepsilon = \tau \sigma(f)$ , kde  $\sigma(f) = \sqrt{\frac{pq}{r}}$ , potom platí podle věty Bienaymé-Čebyševovy, že pro

$$|\xi_i| \geq \tau \sigma(f) \text{ bude } P_\tau \leq \frac{1}{\tau^2} \text{ čili } P_\tau \leq \frac{pq}{r\varepsilon^2} \text{ neboť } \frac{1}{\tau} = \frac{\sigma(f)}{\varepsilon}.$$

Může tedy býti pravděpodobnost  $P_\tau$  pro rostoucí  $r$  při určité zvoleném  $\varepsilon$  libovolně malá. Naopak pro pravděpodobnost

$$1 - P_\tau \text{ že } \left| \frac{x}{r} - p \right| < \varepsilon \text{ bude platit } 1 - P_\tau > 1 - \frac{1}{\tau^2} \text{ čili}$$

$$1 - P_\tau > 1 - \frac{pq}{r\varepsilon^2}. \quad (42)$$

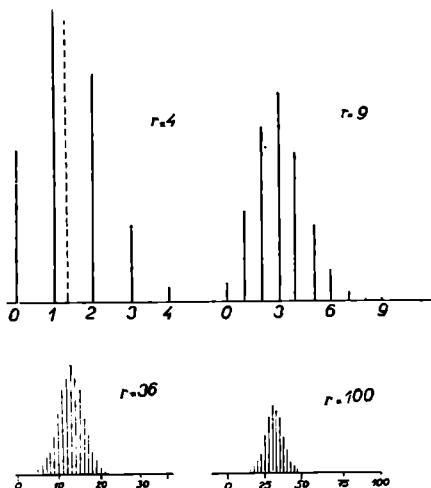
Tato pravděpodobnost se blíží 1, když  $r$  roste nad všechny meze. Odhad na pravé straně (42) můžeme provést nezávisle na určité hodnotě  $p$ , neboť součin  $pq$  nemůže býti větší než  $\frac{1}{4}$ , vzhledem k tomu, že  $p + q = 1$ , takže bude  $1 - P_\tau > 1 - \frac{1}{4r\varepsilon^2}$ . Tento odhad je přirozeně slabší.

Tak jsme dospěli k teorému Bernoulliuvu, který je jedním ze základních pilířů statistiky. Výraz (42) vyjadřuje teorém Bernoulliuv jako větu o mezní hodnotě nejjednodušší formou. Osvětluje otázku, jak se blíží relativní četnost znaku ve výběru o  $r$  prvcích své hodnotě v základním souboru, t. j. konstantní pravděpodobnosti  $p$ , když rozsah  $r$  roste a vyslovíme jej takto:

Je-li  $p$  pravděpodobnost výskytu znaku pro každý prvek náhodného výběru rozsahu  $r$ , pak se pravděpodobnost  $P_\tau$  odchylky  $\frac{x}{r} - p$  relativní četnosti znaku ve výběru od hodnoty  $p$  v základním souboru, která se rovná libovolně malému kladnému číslu  $\varepsilon$ , blíží k nule jakožto limitě, roste-li rozsah náhodného výběru  $r$  nade všechny meze. Pravděpodobnost  $1 - P_\tau$ , že odchylka relativní četnosti znaku ve

výběru  $\frac{x}{r}$  od hodnoty  $p$  v základním souboru bude menší než  $\varepsilon$ , se blíží 1 neboli jistotě.

Způsob, jímž jsme přešli od rozdělení četností absolutních  $x$  k rozdělení relativních četností  $\frac{x}{r} = f$  můžeme považovati za transformaci souřadnic, která sesunuje k sobě úsečky funkce  $P_r(x)$  v poměru  $r : 1$ . Průměr transformovaného rozdělení je konstanta  $p$  a rozptyl  $\sigma^2(f) = \frac{\sigma^2(x)}{r^2} = \frac{pq}{r}$ ,



Obr. 12. Zhušťování binomického rozdělení četností a klesající rozptyl.

neboť rozptyl se mění se čtvercem úseček. Rozptyl tedy klesá k nule s rostoucím  $r$ . Rozdělení  $P_r(f)$  je znázorněno v obr. (12) pro  $p = \frac{1}{3}$ ; pro  $r = 100$  bylo možno zobrazit jeň každou druhou pořadnici. Ubývání rozptylu tu jasně vidíme a současně se jeví, můžeme říci, zhušťování rozdělení četností, čímž je vyjádřena podstata Bernoulliho teorému.

Můžeme jej také formulovati větou:

Relativní četnost nějakého znaku, zjištěná v náhodném výběru rozsahu  $r$  na sobě nezávislých prvků, se blíží hodnotě  $p$  v základním souboru až na odchylku (chybu)  $\varepsilon$  napřed danou s pravděpodobností, která se může zvětšováním rozsahu  $r$  přiblížiti libovolně blízko číslu 1.

V tomto smyslu tudíž reprezentuje náhodný výběr rozsahu  $r$  celý soubor všech prvků, odpovídajících pojmu určujícímu statistickou jednotku, tím lépe, čím je rozsah výběru  $r$  větší. Je to základní věta o větší bezpečnosti delší statistické řady, která tvoří v podstatě obsah t. zv. zákona velkého čísla. Věta o větší bezpečnosti delší statistické řady dává oprávnění principu statistické indukce; mohla by býti označena také jako věta o větší bezpečnosti závěru provedeného statistickou indukcí na základě náhodného výběru o větším rozsahu než na základě výběru o menším rozsahu.

Zákon velkého čísla souvisí přímo s principem stejnotvárnosti přírodního dění, podle něhož za podobných okolností jev probíhá podobně. Předpokládá se tedy, že stejné skupiny (komplexy) příčin mají za následek stejné pochody. Abychom však z teoremu matematicky odvozeného mohli činit závěry na skutečné dění, musíme učinit ještě další krok.

Budeme se dovolávat zkušenosti, že v souborech, které mají povahu našeho základního souboru, pozorujeme skutečně zřídka prvků o znaku s malou relativní četností.

Vyvodíme z toho potom závěr, že velké odchylky  $\frac{x}{r} - p$  se budou u statistických souborů rovněž jen zřídka vyskytovat. To je podstatný obsah věty Cournotovy a jeho formulace zákona velkého čísla.

Na konec nám zbývá určití celkovou relativní četnost všech těch prvků našeho nového (myšlenkového) souboru rozsahu  $\binom{N}{r}$ , u nichž se vyskytuje znak  $c$  nejméně  $(rp - \xi_0)$ -krát a nejvýše  $(rp + \xi_0)$ -krát, neboli u nichž je relativní četnost znaku  $c$  v mezích od  $p - \frac{\xi_0}{r}$  do  $p + \frac{\xi_0}{r}$ . Pro první případ dostaneme hledanou celkovou relativní četnost  $P_r(\bar{x} - \xi_0, \bar{x} + \xi_0)$ , když sečteme všechny relativní četnosti  $P_r(x)$  pro hodnoty  $x$  v uvedeném intervalu. Pro druhý případ obdobně  $P_r(p - z_0, p + z_0)$  dostaneme sečtením přísluš-

ných hodnot  $P_r(f)$ ; při tom je ovšem  $P_r(x) = P_r(f)$ ,  $\frac{\xi_0}{r} = z_0$ . Vzhledem k tomu, že výpočet jednotlivých členů je dosti pracný a tedy také jejich součet

$$P_r(\bar{x} - \xi_0, \bar{x} + \xi_0) = \sum_{\bar{x}-\xi_0}^{\bar{x}+\xi_0} P_r(x) = [P_r(\bar{x} + \xi_0) + P_r(\bar{x} - \xi_0)] + [P_r(\bar{x} + \xi_0 - 1) + P_r(\bar{x} - \xi_0 + 1)] + \dots + [P_r(\bar{x} + 1) + P_r(\bar{x} - 1)] + P_r(\bar{x})$$

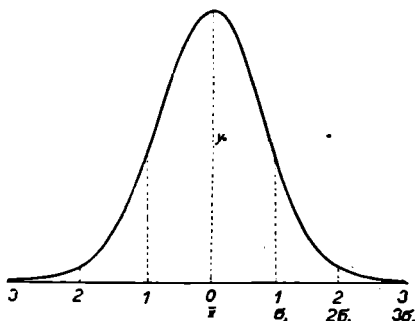
najdeme dále — rovnice (52) nebo (53) — vyhovující řešení přibližné.

Příklad: relativní četnost narozených chlapců v základním souboru je  $p = 0,513$  čili 51,3%. Průměrná absolutní četnost chlapců narozených ročně v místě, kde se rodí ročně  $r = 100$  dětí, je tedy  $rp = 51,3$ ; v místě, kde se rodí ročně 10 000 dětí, bude  $rp = 5130$ . Směrodatná odchylka činí v prvním případě  $\sqrt{rpq} = \sqrt{0,513 \times 0,487r} = 0,5\sqrt{r} = 5$ , kdežto v druhém případě 50. Směrodatná odchylka relativních četností chlapců však klesá, neboť je v prvním případě  $\sqrt{\frac{pq}{r}} = \frac{0,5}{10} = 0,05$ , tedy 5%, kdežto v druhém případě jen 5 promille. Výsledek se považuje za tím přesnější, čím má menší rozptyl a tedy také, čím má menší směrodatnou odchylku. Je z toho zřejmo, že čím jsou náhodné výběry většího rozsahu, tím dávají výsledek bližší hodnotě v základním souboru.

**(6,1) Křivky rozdělení četností.** (Křivka Laplace-Gaussova.)

Lze říci, že Bernoulli začal studovat binomické rozdělení četností a vyjádřil jednu jeho zvláštní vlastnost ve větě po něm pojmenované, která ukazuje, že vytkneme-li libovolně malý interval kolem hodnoty  $p$  a určíme si číslo libovolně blízké jednotce, pak můžeme zvoliti soubor o dosti

velkém počtu prvků, takže relativní četnost pozorovaného znaku padne do zvoleného intervalu s určenou pravděpodobností. Nahradiť binomické rozděléní spojitou křivkou se podařilo Laplaceovi (1812), takže bylo úplně dáno souměrnou zvonovitou křivkou  $e^{-\xi^2}$ , jejíž pořadnice klesají od průměru tak, že se jejich přirozené logaritmy (se záporným znaménkem) chovají jako čtverce vzdálenosti od průměru (viz obraz 13).



Obr. 13. Křivka Laplace-Gaussova.

Odvodili jsme si již z binomického rozděléní četností spojitou křivku Laplace-Gaussovu, čili normální ve zcela zvláštním případě, kde základní relativní četnosti při alternativním znaku byly sobě rovny, tedy  $p = q = \frac{1}{2}$ . Lze však ukázati, že obecné rozděléní binomické  $(p + q)^r$  se blíží pro velká  $r$  křivce normální. Abychom tento postup naznačili, vyjádříme členy binomického rozděléní (37) hodnotami  $y(\xi)$  v jednotkových intervalech tak, že pro rozdíl mezi průměrem  $\bar{x} = rp$  a četností  $x$  znaku ve výběru rozsahu  $r$  zvolíme symbol  $\xi$ . Potom jednotlivé členy rozděléní četností budou

$$y(\xi) = \frac{r!}{(pr + \xi)!(qr - \xi)!} p^{pr + \xi} q^{qr - \xi}. \quad (43)$$

K přibližnému vyjádření faktoriel použijeme Stirlingovy formule

$$n! = n^n e^{-n} \sqrt{2\pi n} \left( 1 + \frac{1}{12n} + \frac{1}{288n^2} + \dots \right).$$

Užijeme-li jen prvního členu této řady, dostaneme přibližnou hodnotu, která se rovná přesné hodnotě dělené nějakým číslem mezi 1 a  $1 + \frac{1}{10n}$ . Stačí tudíž většinou toto přiblížení pro  $n$ , která přicházejí v úvahu. S tímto přiblížením pak dostáváme

$$\begin{aligned} (pr + \xi)! &= (pr + \xi)^{pr + \xi} e^{-(pr + \xi)} \sqrt{2\pi (pr + \xi)} = \\ &= (pr)^{pr + \xi} \left( 1 + \frac{\xi}{pr} \right)^{pr + \xi + \frac{1}{2}} e^{-(pr + \xi)} \sqrt{2\pi pr} \end{aligned}$$

a podobně

$$(qr - \xi)! = (qr)^{qr - \xi} \left( 1 - \frac{\xi}{qr} \right)^{qr - \xi + \frac{1}{2}} e^{-(qr - \xi)} \sqrt{2\pi qr},$$

takže po dosazení do (43) a jednoduché úpravě bude přibližně

$$y(\xi) = \frac{1}{\sqrt{2\pi r p q}} \left( 1 + \frac{\xi}{pr} \right)^{-(pr + \xi + \frac{1}{2})} \left( 1 - \frac{\xi}{qr} \right)^{-(qr - \xi + \frac{1}{2})} \quad (44)$$

K osvětlení, jak se přibližuje tento výraz k (43), srovnáváme odchylky  $\xi$  od průměru se směrodatnou odchylkou  $\sigma_x = \sqrt{r p q}$ , která je řádu  $\sqrt{r}$ , není-li  $p$  ani  $q$  příliš malé. Musíme tedy předpokládati  $r$  tak velké, aby bylo možno zanedbat  $\frac{\xi}{r}$ , ale  $\sqrt{\frac{\xi}{r}}$  musí míti takové konečné hodnoty, jaké se nám vyskytují, když posuzujeme odchylky  $\xi$  srovnáváním se směrodatnou odchylkou.

Můžeme tedy výraz (44), který napíšeme ve tvaru

$$y(\xi) = \frac{1}{\sqrt{2\pi r p q}} A \cdot B,$$



zjednodušiti s uvedenou přibližností, neboť

$$\log A = - \left( rp + \xi + \frac{1}{2} \right) \left[ \frac{\xi}{rp} - \frac{\xi^2}{2r^2p^2} + \frac{\xi^3}{r^3} \Phi_1(\xi) \right]$$

$$\log B = - \left( rq - \xi + \frac{1}{2} \right) \left[ -\frac{\xi}{rq} - \frac{\xi^2}{2r^2q^2} - \frac{\xi^3}{r^3} \Phi_2(\xi) \right],$$

takže

$$\log y(\xi) \sqrt{2\pi r p q} = \frac{(p-q)\xi}{2rpq} - \frac{\xi^2}{2rpq} + \frac{\xi^2}{r^2} \Phi_3(\xi) =$$

$$= \frac{\xi^2}{2rpq} + \frac{\xi}{r} \Phi(\xi),$$

kde  $\Phi_i(\xi)$  a  $\Phi(\xi)$  jsou konečné, neboť představují součty konvergentních řad mocnin zlomku  $\frac{\xi}{r}$ , který je libovolně malý.

Je-li tudíž  $r$  tak velké, že  $\frac{\xi}{r} \Phi(\xi)$  je malé, a tedy zanedbatelné, dostáváme

$$y(\xi) = \frac{1}{\sqrt{2\pi r p q}} e^{-\frac{\xi^2}{2\pi r p q}}.$$

Vzhledem k tomu, že  $\sigma_x^2 = r p q$ , můžeme také psáti

$$y(\xi) = \frac{1}{\sigma_x \sqrt{2\pi}} e^{-\frac{\xi^2}{2\sigma_x^2}}, \quad (45)$$

což je normální křivka rozdělení četností.

Ve směrodatné proměnné  $\frac{x - \bar{x}}{\sigma_x} = \frac{\xi}{\sigma_x} = t$  má pak tvar

$$y(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2}. \quad (46)$$

Křivka je symetrická podle průměru, do něhož jsme položili počátek souřadnic, který je tedy v bodě  $t = 0$  a této

hodnotě odpovídá maximální pořadnice

$$y(0) = \frac{1}{\sqrt{2\pi}} = 0,39894.$$

Uvedeme si přehled několika pořadnic v intervalech  $0,5\sigma_x$ :

$\xi/\sigma_x = 0,5$	1,0	1,5	2,0	2,5	3,0
$y = 0,35207$	0,24117	0,12952	0,05399	0,01753	0,00443
$y/y(0) = 0,88250$	0,60653	0,32465	0,13534	0,04394	0,01111

Podrobnější tabulku poměru pořadnic  $y : y(0)$  možno najíti na př. v [9]. Druhá derivace výrazu (45) je

$$\frac{d^2y}{d\xi^2} = \frac{1}{\sqrt{2\pi}\sigma_x^3} \left( \frac{\xi^2}{\sigma_x^2} - 1 \right) e^{-\frac{\xi^2}{2\sigma_x^2}},$$

z čehož plyne, že křivka má dva inflexní body pro  $\xi = \pm \sigma_x$ , neboť nejbližší derivace v těch bodech od nuly různá je třetí, tedy lichého stupně. Tečny v těchto bodech křivky protínají osu  $\xi$  v bodech  $\xi = \pm 2\sigma_x$ .

Momenty lichého řádu kolem průměru jsou pro symetrickou křivku Laplace-Gaussovu, jako pro každou symetrickou křivku, rovny nule, tedy  $\mu_{x,1} = \mu_{x,3} = \mu_{x,5} = \dots = 0$ . Pro momenty sudého řádu lze odvodit rekurentní vztah [6]

$$\mu_{x,2i} = (2i - 1) \sigma_x^2 \mu_{x,2i-2} \quad (47)$$

takže

$$\mu_{x,2} = \sigma_x^2, \quad \mu_{x,4} = 3\sigma_x^4, \quad \mu_{x,6} = 15\sigma_x^6, \dots \quad (48)$$

Jako jsme viděli u histogramu, že celá jeho plocha představuje rozsah souboru, tak jej také zde znázorňuje plocha ohraničená křivkou a osou  $x$ . Tato plocha je dána při spojitě proměnné integrálem

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-t^2} dt = 1$$

vzhledem k tomu, že

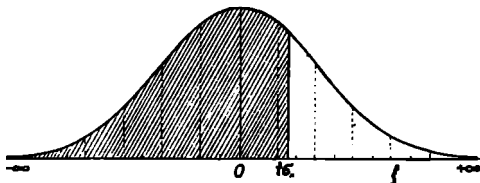
$$\int_{-\infty}^{+\infty} e^{-t^2} dt = \sqrt{2\pi}.$$

Máme-li soubor rozsahu  $r$ , pak je rovnice normální křivky

$$y(\xi) = \frac{r}{\sigma_x \sqrt{2\pi}} e^{-\frac{\xi^2}{2\sigma_x^2}} \quad (49)$$

a maximální pořadnice pro  $\xi = 0$  je  $y_r(0) = \frac{r}{\sigma_x \sqrt{2\pi}}$ .

Vzhledem k souměrnosti křivky, je část ohraničená osou  $\xi$



Obr. 14a. Úseky plochy normální křivky četností.

a křivkou v mezích od  $-\infty$  do 0 rovna (viz obr. 14a) polovině celé plochy, tedy 0,5. Část plochy od  $-\infty$  až do  $\xi = t\sigma_x$ , kde  $t$  je kladné číslo, bude

$$F(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-t^2} dt = 0,5 + \frac{1}{2}\alpha(t), \quad (50)$$

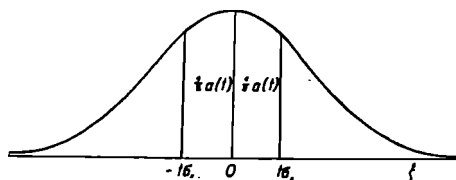
kde jsme zavedli

$$\frac{1}{2}\alpha(t) = \frac{1}{\sqrt{2\pi}} \int_0^t e^{-t^2} dt. \quad (51)$$

Bude tedy plocha pásu (obr. 14b) mezi  $\xi = -t\sigma_x$  a  $\xi = +t\sigma_x$

$$\alpha(t) = \frac{1}{\sqrt{2\pi}} \int_{-t}^t e^{-t^2} dt = \frac{2}{\sqrt{2\pi}} \int_0^t e^{-t^2} dt. \quad (52)$$

Tyto hodnoty můžeme sestaviti do tabulky [6] pro různá  $t$ , t. j. pro různé hodnoty odchylky od průměru vyjádřené ve směrodatné odchylce jako jednotce. Tak je na př.



Obr. 14b.

$\frac{\xi}{\sigma_x} = t$	0,6745	1	$\sqrt{2}$	2	3
$\alpha(t)$	0,5	0,6827	0,8427	0,9545	0,9973
B.-Č.	0	0,5	0,750	0,889	

Hodnota  $\xi_p = 0,6745\sigma_x$  se nazývá také pravděpodobná chyba. Je patrné, že v mezích  $\bar{x} \pm \xi_p$  je polovina celého souboru a tedy rovněž polovina vně těchto mezí. Je tudíž tato hodnota kvartilovou odchylkou. Hodnota  $\xi_m = \sqrt{2}\sigma_x$  se nazývá modul. Užívá se jí také někdy za jednotku, v níž se vyjadřuje proměnná, takže  $\frac{\xi}{\sigma_x\sqrt{2}} = \gamma$  a dostáváme pak funkci

$$\Phi(\gamma) = \frac{2}{\sqrt{\pi}} \int_0^\gamma e^{-\gamma^2} d\gamma, \quad (53)$$

kteřá bývá tabelována. Přesvědčíme se ovšem snadno, že  $\alpha(t) = \Phi(\gamma)$ , provedeme-li substituci  $t = \gamma\sqrt{2}$ .

Ve třetím řádku byly pro srovnání uvedeny hodnoty vyplývající z teoremu Bienaymé-Čebyševova.

Z uvedených čísel vidíme, jaké procento prvků souboru s normálním rozdělením četností je v určitých mezích odchylek od průměru.

Tak bude prvků:

68,27%	s odchylkou	$ \xi  \leq \sigma_x$ ,	ostatních je tedy	31,73%
95,45%	„	$ \xi  \leq 2\sigma_x$ ,	„ „ „	4,55%
99,73%	„	$ \xi  \leq 3\sigma_x$ ,	„ „ „	0,27%

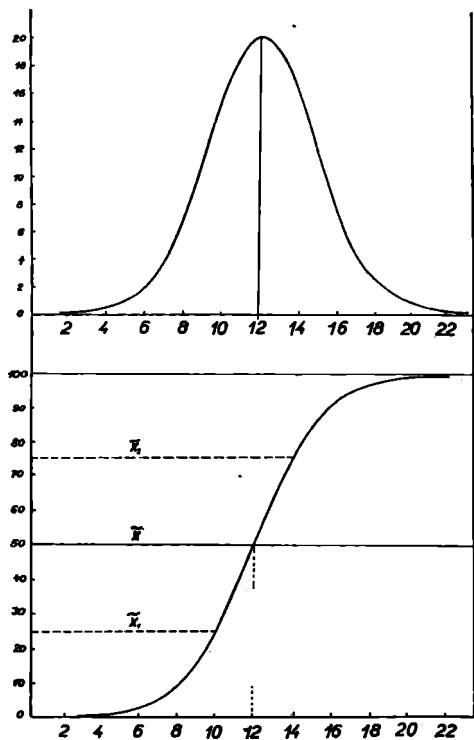
V souboru s normálním rozdělením četností podle toho bude na př. 0,135% prvků s většími hodnotami než  $\bar{x} + 3\sigma_x$  a rovněž tolik s menšími hodnotami než  $\bar{x} - 3\sigma_x$ , čili 0,27% prvků bude mimo interval  $\pm 3\sigma_x$ .

Vlastní praktický význam křivky Laplace-Gaussovy se jeví teprve při těchto daleko důležitějších otázkách, kde potřebujeme součet velkého počtu jednotlivých relativních četností, neboť méně nás zajímá otázka, jaká jest pravděpodobnost, že při 1200 vrzích kostkou padne právě  $x = 180$ krát šestka, jako spíše otázka, jaká je pravděpodobnost, že nebude odchylka od průměru  $\bar{x} = 200$  větší

než  $200 - 180 = 20$ . To vyžaduje zjistiti součet  $\sum_{x=180}^{220} P_r(x)$  čili vypočítati podle Newtonovy formule (36) celkem 41 jednotlivých hodnot pravděpodobností  $P_r(x)$  a sečísti. Integraci křivky Laplace-Gaussovy dosahuje se zde dalekosáhlého zjednodušení. Pro tento t. zv. Laplaceův integrál existují různé tabulky sestojené pro různé argumenty; proto je třeba značné opatrnosti při jejich užívání a především řádného seznámení se s nimi.

Znáznorníme si hodnoty funkce  $F(t)$ , probíhá-li proměnná  $t$  celý obor reálných čísel; dostáváme tak k normálnímu roz-

dělení četností součtovou křivku, která je znázorněna v obr. 15. Její pořadnice je v stupnici pětkrát zmenšené



Obr. 15. Součtová křivka k normálnímu rozdělení četností.

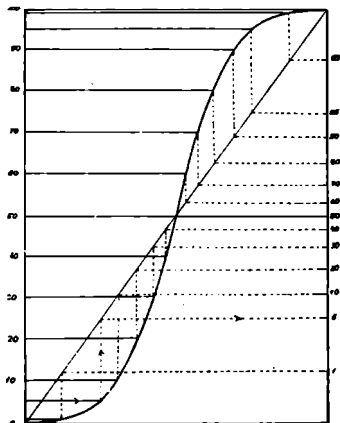
proti pořadnici příslušné normální křivky nahoře. Pomocí součtové křivky se snadno určuje, jak již víme, medián a kvartily.

**(6,2) Normální rozdělení četností kvantitativního znaku.** Odvodili jsme Laplace-Gaussovu křivku normálního rozdělení četností, pomocí náhodných výběrů, z nichž každý vykazuje určitou relativní četnost znaku  $\frac{x}{r} = f$  a pochopili jsme tak vznik této křivky na základě binomického rozdělení četností. Normální rozdělení však vzniká také, provedeme-li mnohonásobné měření kvantitativního znaku na jednom předmětu (řada měření nějaké délky) nebo při měření určitého kvantitativního znaku na různých předmětech, jež jsou prvky jednoho statistického souboru (na př. délka listů určitého stromu). Pro výklad, jak vzniká normální rozdělení, v prvním případě si můžeme představit, že výsledek každého měření závisí na velikém počtu t. zv. elementárních příčin, z nichž každá je s to způsobiti nějakou elementární odchylku od skutečnosti. Tyto odchylky jsou v obou směrech stejně pravděpodobné a vzájemně nezávislé. Je to tedy analogie s náhodnými výběry koulí z osudí se stejnou pravděpodobností pro bílou i černou nebo analogie házení mincí. Takový výklad byl sestrojen původně pro teorii chyb při měření; přenášel se pak také na druhý případ, jímž se zabýváme ve statistice. Je však také jiný výklad, který snad lépe vystihuje skutečnost, takže se ho můžeme přidržeti. Vychází od hypotese, že každá hodnota kvantitativního znaku je součtem množství neznámých a nezávislých sčítanců. Na př. délka nějakého předmětu (listu) se skládá z délky velkého množství nezávislých součástí (buněk). Také tudíž každá jednotlivá odchylka je součtem množství malých neznámých veličin, elementárních odchylek. Rozdělení těchto součtů je blízké normálnímu i když by rozdělení sčítanců nebylo normální. (Tak si můžeme vysvětlit, že se vyskytuje normální rozdělení četností také pro kvantitativní znak ve statistických souborech.)

**(6,3) Pravděpodobnostní stupnice.** Součtovou křivku patřící k normální křivce lze znázorniti přímkou, zvolíme-li

vhodnou stupnicí pro pořadnici. Souvislost pravidelné stupnice relativních četností v procentech se stupnicí t. zv. pravděpodobnostní rovněž v procentech je vyznačena nomograficky v obr. 16, kde jsou patrné body přímky odpovídající bodům součtové křivky.

V této stupnici je znázorněna součtová křivka rozdělení četností na str. 29 (obr. 4b). Podle toho, jak se odchyluje od přímky, můžeme posoudit, že pozorované rozdělení četností se liší od normálního.



Obr. 16. Převod pravidelné stupnice na stupnici pravděpodobnostní.

Normální křivku lze rovněž převést na přímku, zvolíme-li v pravoúhlé soustavě na ose úseček kvadratickou stupnici a na ose pořadnic logaritmickou stupnici [7, str. 19].

**(6,4) Poissonovo rozdělení četností.** (Exponenciální Poissonova.) Abychom z binomického rozdělení četností odvodili ještě jiné křivky rozdělení četností, budeme hledat pro funkci

$$y = \frac{r!}{x! (r-x)!} p^x q^{r-x}$$

vhodný výraz, který by ji vyjádřil přibližně v těch případech, kdy základní pravděpodobnost výskytu pozorovaného znaku  $p$  je malá, ale tak, že  $rp = \lambda$  je číslo konečné pro libovolně veliké  $r$ .

Především je

$$\frac{r!}{(r-x)!} = r(r-1)(r-2) \dots (r-x+1).$$



Dále pišme

$$p = \frac{\lambda}{r} \text{ a tedy } q = 1 - \frac{\lambda}{r},$$

takže bude tedy

$$y = \left(1 - \frac{1}{r}\right) \left(1 - \frac{2}{r}\right) \dots \left(1 - \frac{x-1}{r}\right) \cdot \frac{\lambda^x}{x!} \left(1 - \frac{\lambda}{r}\right)^r q^{-x}.$$

Přibližný výraz dostaneme pro velká  $r$ , zanedbáme-li veličiny řádu  $\frac{1}{r}$ , takže především součin prvních  $x-1$  činitelů v závorkách, který je mezi 1 a  $1 - \frac{x(x-1)}{2r}$ , položíme roven přibližně 1.

Dále můžeme místo  $\left(1 - \frac{\lambda}{r}\right)^r$  klásti přibližně  $e^{-\lambda}$ , což je limita, k níž výraz spěje pro  $r \rightarrow \infty$ . Konečně  $q^{-x}$  spěje pro velká  $r$  k 1, neboť  $q^{-x} = \left[\left(1 - \frac{\lambda}{r}\right)^{-r}\right]^{\frac{x}{r}}$ , což spěje k  $(e^{\lambda})^{\frac{x}{r}}$ , a tedy pro velká  $r$  se  $\frac{x}{r}$  blíží k nule. Z toho všeho tudíž vyplývá, že můžeme přibližně klásti

$$y = \frac{e^{-\lambda} \lambda^x}{x!}; \quad (54)$$

tento výraz se obvykle označuje symbolem  $\psi(x)$  a nazývá se exponenciela Poissonova, udávající pravděpodobnost, že se vyskytuje  $x$ -krát pozorovaný znak, který patří mezi tak zv. řídké jevy, jejichž pravděpodobnost  $p$  je malá. Bortkiewicz jej nazval zákonem malých čísel.

Pravděpodobnosti, že se objeví pozorovaný znak právě 0, 1, 2, ... krát, jsou dány jednotlivými členy řady

$$e^{-\lambda} \left(1 + \lambda + \frac{\lambda^2}{2!} + \frac{\lambda^3}{3!} + \dots\right).$$

Ačkoliv jsme předpokládali při odvozování Poissonovy exponenciely, že  $x$  je malé vzhledem k  $r$ , dostáváme k rozdělení četností, vyjádřenému touto exponencielou, klademe-li za  $x$  všechna celá čísla od  $x=0$  do  $x=r$ , jednoduché a důležité výsledky pro průměr a směrodatnou odchylku. Pro velká  $r$  platí přibližně

$$e^{-\lambda} \left( 1 + \lambda + \frac{\lambda^2}{2!} + \dots + \frac{\lambda^r}{r!} \right) = 1, \quad (55)$$

neboť součet v závorce je přibližně roven  $e^\lambda$ . Jednotlivé členy pravé strany jsou tedy relativní četnosti. Vynásobíme-li každou z nich příslušnou hodnotou znaku 0, 1, 2, ...,  $r$ , dostaneme průměr

$$\begin{aligned} \bar{x} &= e^{-\lambda} \left( 0 + \lambda + \lambda^2 + \frac{\lambda^3}{2!} + \dots + \frac{\lambda^r}{(r-1)!} \right) = \\ &= \lambda e^{-\lambda} \left( 1 + \lambda + \frac{\lambda^2}{2!} + \dots + \frac{\lambda^{r-1}}{(r-1)!} \right) = \lambda, \end{aligned} \quad (56)$$

neboť součet v závorce je pro velká  $r$  přibližně týž jako v rovnici (55).

Podobně dostaneme pro druhý moment obecný

$$\mu'_{x,2} = e^{-\lambda} \left( 0 + \lambda + 2\lambda^2 + \frac{3\lambda^3}{2!} + \dots + \frac{r\lambda^r}{(r-1)!} \right),$$

takže rozptyl

$$\mu_{x,2} = \mu'_{x,2} - \bar{x}^2$$

bude

$$\mu_{x,2} = \lambda e^{-\lambda} \left( 1 + 2\lambda + \frac{3\lambda^2}{2!} + \dots + \frac{r\lambda^{r-1}}{(r-1)!} \right) - \lambda^2,$$

což lze také psáti

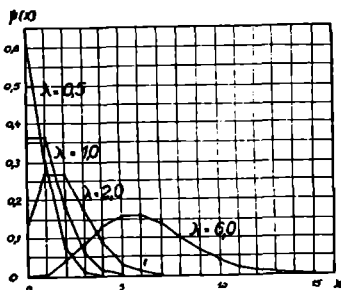
$$\begin{aligned} \mu_{x,2} &= \lambda e^{-\lambda} \left( 1 + \lambda + \frac{\lambda^2}{2!} + \dots + \frac{\lambda^{r-1}}{(r-1)!} \right) + \\ &+ \lambda^2 e^{-\lambda} \left( 1 + \lambda + \frac{\lambda^2}{2!} + \dots + \frac{\lambda^{r-2}}{(r-2)!} \right) - \lambda^2; \end{aligned}$$

vidíme tedy vzhledem k (55)

$$\sigma_x^2 = \lambda + \lambda^2 - \lambda^2 = \lambda. \quad (57)$$

Je tedy rozptyl roven průměru. Vzhledem k tomu, že  $\lambda = rp$ , je to tedy hodnota blízká  $rpq$ , kterou jsme našli pro normální rozdělení četností, neboť  $q$  se liší velmi málo od 1.

Hodnoty Poissonovy exponenciální limity (54) byly tabulovány pro různá  $\lambda$  a  $x$ ; lze je najít na př. v tabulkách [8]. Průběh jejich je znázorněn na obr. 17 pro  $\lambda = 0,5, 1, 2, 6$ .



Obr. 17. Exponenciála Poissonova.

Je jasně viděti, že od úplné nesouměrnosti přecházejí křivky pro rostoucí  $\lambda$  k tvaru stále souměrnějšímu.

**(6,5) Pearsonův systém křivek četností.** Viděli jsme, že lze odvodit z binomického rozdělení četností čili z formule Newtonovy (36) celý systém křivek rozdělení četností. Mohou však býti odvozeny ještě obecnější systémy. Představme

si, že základní soubor konečného rozsahu  $N$  obsahuje  $k$  prvků, majících pozorovaný alternativní znak a  $N - k$  prvků, které jej nemají. Vyjmeme-li z tohoto základního souboru částečné soubory o rozsahu  $r$  prvků, můžeme tak učiniti celkem  $\binom{N}{r}$  různými způsoby, čili můžeme dostati tolik různých výběrů. Každý z těchto výběrů má určitý počet  $x$  prvků s uvažovaným znakem. Kladné celé číslo  $x$  je v intervalu od 0 do  $r$ , když předpokládáme  $k \geq r$ . Abychom stanovili, kolik může býti různých výběrů, jež mají určitý počet  $x$  prvků s pozorovaným znakem, uvědomíme si, že je celkem  $\binom{k}{x}$  skupin, jež obsahují  $x$  různých

prvků z daných  $k$  prvků s uvažovaným znakem v základním souboru, a ke každé z těchto skupin lze přiřaditi  $\binom{N-k}{r-x}$  různých skupin tvořených ze zbývajících  $r-x$  prvků, které nemají uvažovaný znak a doplňují skupinu na celkový rozsah výběru  $r$ . Vidíme, že tedy bude hledaný počet různých výběrů čili absolutní četnost  $\binom{k}{x} \binom{N-k}{r-x}$ . Relativní četnost jejich dostaneme, dělíme-li poslední výraz celkovým počtem různých možných výběrů rozsahu  $r$ , takže bude vyjádřena funkcí

$$f(x) = \frac{1}{\binom{N}{r}} \binom{k}{x} \binom{N-k}{r-x} \quad (58)$$

pro hodnoty  $x = 0, 1, 2, \dots, r$ ; jsou to postupně za sebou jdoucí členy konečné řady hypergeometrické. Z této funkce vyšel K. Pearson, aby odvodil t. zv. Pearsonův systém křivek rozdělení četností, v němž jsou křivky (45) a (54) zahrnuty jako zvláštní případy [6], neboť také binomická funkce (36) je zvláštním případem hypergeometrické, která v ni přechází pro nekonečný rozsah základního souboru, takže  $N = \infty$ ,  $k = \infty$ , ale jejich poměr  $\frac{k}{N} = p$  je konstantní a konečný.

Vhodnou volbou typu křivky je pak možno s postačujícím přiblížením vyjádřiti statisticky pozorovaná rozdělení četností. Tato volba je tu usnadněna tím, že bylo odvozeno — pomocí momentů — kritérium, které umožňuje rozhodnouti se mezi možnými typy pro vhodnější.

**(6,6) Pólyovo výběrové schema pro jevy vázané.** K určité funkci hypergeometrické jsme vedeni, provádíme-li ze základního souboru o  $N$  prvcích, z nichž má  $k$  prvků pozorovaný znak, výběr rozsahu  $r$  tak, že vyjmeme prvek a zjistíme, má-li pozorovaný znak. V kladném případě se

počet prvků s pozorovaným znakem v základním souboru zvětší o  $1 + \Delta$ ; neměl-li prvek pozorovaný znak, zvětší se o  $1 + \Delta$  počet těchto druhých prvků v základním souboru. V okamžiku, když jsme vyňali  $r$  prvků, bude mít základní soubor celkem  $N + r\Delta$  prvků.

Bylo-li mezi nimi  $x$  prvků s pozorovaným znakem a tudíž  $r - x$  ostatních, je v základním souboru  $k + x\Delta$  prvků s pozorovaným znakem a  $N - k + (r - x)\Delta$  ostatních. Pravděpodobnost výskytu pozorovaného znaku je na začátku v základním souboru  $\frac{k}{N} = p$  a opačná  $\frac{N - k}{N} = q$ ; mění se po vynětí každého prvku do výběru, takže po vynětí  $r$ -tého je  $\frac{k + x\Delta}{N + r\Delta}$  resp.  $\frac{N - k + (r - x)\Delta}{N + r\Delta}$ .

Pravděpodobnost, že prvních  $x$  prvků bude mít pozorovaný znak ve výběru rozsahu  $r$ , bude jako složená pravděpodobnost dána součinem

$$\frac{k}{N} \cdot \frac{k + \Delta}{N + \Delta} \cdot \dots \cdot \frac{k + (x - 1)\Delta}{N + (x - 1)\Delta} \cdot \frac{N - k}{N + x\Delta} \cdot \frac{N - k + \Delta}{N + (x + 1)\Delta} \cdot \dots \cdot \frac{N - k + (r - x - 1)\Delta}{N + (r - 1)\Delta}.$$

Zavedeme-li označení  $\frac{\Delta}{N} = \delta$ , přechází poslední výraz na tvar

$$\frac{p}{1} \cdot \frac{p + \delta}{1 + \delta} \cdot \dots \cdot \frac{p + (x - 1)\delta}{1 + (x - 1)\delta} \cdot \frac{q}{1 + x\delta} \cdot \frac{q + \delta}{1 + (x + 1)\delta} \cdot \dots \cdot \frac{q + (r - x - 1)\delta}{1 + (r - 1)\delta}.$$

Pravděpodobnost, že bude ve výběru téhož rozsahu  $r$  jiných  $x$  prvků se znakem pozorovaným, bude dána tímž výrazem, jen pořadí jednotlivých faktorů bude jiné.

Kombinací, v nichž se může vyskytnouti mezi  $r$  prvky  $x$  s pozorovaným znakem je  $\binom{r}{x}$ , takže celkem pravděpodob-

nost, že mezi  $r$  prvky výběru provedeného ze základního souboru, který se uvedeným způsobem mění, bude  $x$  prvků s pozorovaným znakem, je

$$f(x, r) = \binom{r}{x} \cdot \frac{p(p+\delta) \dots [p+(x-1)\delta] q[q+\delta] \dots [q+(r-x-1)\delta]}{[1+\delta][1+2\delta] \dots [1+(r-1)\delta]}$$

Je-li pravděpodobnost  $p$  malá, ale pro velká  $r$  je  $rp = \lambda$  konečné číslo, a při kladném  $\delta$  označíme  $r\delta = d > 0$ , platí přibližně

$$f(x, r) = \frac{1}{x!} \lambda(\lambda+d)(\lambda+2d) \dots (\lambda+x-1d) (1+d)^{-\frac{\lambda}{d}-x} \quad (59)$$

což se nazývá zákonem Pólyovým a je zobecněním exponenciely Poissonovy (54), která z něho vyplývá jako limita pro  $d = 0$ .

Pro uvedené schema výběrové to znamená, že  $\Delta = 0$  čili redukuje se na případ schematu Bernoulliho o konstantní pravděpodobnosti  $p$ .

Jiný případ dostaneme pro  $\Delta = -1$ , který znamená, že prvek se vyjme do výběru a rozsah základního souboru se tím vždy o jeden prvek zmenšuje, což je případ Pearsonův, odpovídající konečnému souboru základnímu, do něhož se vyňatý prvek nevrací zpět.

Zákon Pólyův uvádíme vzhledem k jeho obecnosti a také proto, že se osvědčil k vystižení případů, kde se nejedná o jevy nezávislé, nýbrž nějakým způsobem vázané, jako je případ úmrtnosti vlivem nakažlivých nemocí, nebo smrti cestujících následkem neštěstí na drahách a pod.

**(6,7) Rozvoje v řady.** (Řada Brunsova.) Praktický problém vyjádření pozorovaného rozdělení četností analyticky je také velmi obecně řešen pomocí řady, jejímž prvním členem je funkce Laplace-Gaussova jako vytvářející a dalšími členy její derivace.

Omezíme-li se jen na první dva členy, od nuly různé, dostáváme vyjádření

$$f(\xi) = \varphi(\xi) - \varphi(\xi) \frac{\mu_{x,3}}{3! \sigma_x^3} \left( \frac{3\xi}{\sigma_x} - \frac{\xi^3}{\sigma_x^3} \right), \quad (60)$$

kde

$$\varphi(\xi) = \frac{1}{\sigma_x \sqrt{2\pi}} e^{-\frac{\xi^2}{2\sigma_x^2}}, \quad \varphi(\xi) \left( \frac{3\xi}{\sigma_x} - \frac{\xi^3}{\sigma_x^3} \right) = \varphi'''(\xi),$$

takže z pozorovaného rozdělení četností musíme stanovit první tři momenty, abychom určili potřebné tři konstanty, jež se ve výrazu vyskytují  $\bar{x}$ ,  $\sigma_x$ ,  $\mu_{x,3}$ .

(Řada Poisson-Charlierova.) Také jiné funkce mohou sloužiti k podobným rozvojem v řadu. Zvláště jednoduchý a pro vystižení nesymetrických rozdělení četností vhodný je rozvoj pomocí exponenciely Poissonovy, který uvedeme také bez odvozování

$$f(x) = \psi(x) + \frac{1}{2}(\mu_{x,2} - \lambda) \Delta^2 \psi(x). \quad (61)$$

$\psi(x)$  značí exponenciely Poissonovu (54) a druhá diference je

$$\Delta^2 \psi(x) = \psi(x) - 2\psi(x-1) + \psi(x-2).$$

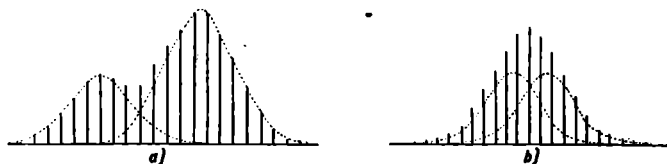
Na rozdíl od Pearsonova systému křivek nemáme zde kriteria, které by nám pomohlo rozhodnouti, zda máme použiti řady Brunsovy nebo Poisson-Charlierovy pro rozpojitou proměnnou  $x$ , takže musí rozhodnouti statistik sám podle vhodnosti a účelnosti.

Praktický význam rozvojem vytvořených pomocí jiných funkcí je omezen požadavkem rychlé konvergence řady, aby bylo možno se omeziti jen na několik málo členů.

**(6,8) Vícevrcholová rozdělení četností.** Někdy se vyskytují soubory, jejichž rozdělení četností má dva vrcholy (čili dvě maxima) jako v obr. 18a nebo více vrcholů. Vznik takového rozdělení četností se vysvětluje tím, že soubor zahrnuje prvky nestejnorodé podle některého znaku, takže bychom dostali dvě různá rozdělení četností, kdybychom

podle něho soubor roztrídili. Představujeme si tedy, že výsledná křivka rozdělení četností vznikla superposicí dvou jednoduchých křivek; při tom by ovšem mohla vzniknouti také křivka jednovrcholová (obr. 18b).

Za účelem oddělení obou jednoduchých křivek je možno použití pro některé tvary dvojevrcholových rozdělení čet-



Obr. 18a, b. Dvojevrcholové rozdělení četností.

ností křivek normálních, takže pak dané rozdělení je vyjádřeno rovnicí

$$r f(x) = \frac{r_1}{1\sigma_x\sqrt{2\pi}} e^{-\frac{(x-\bar{x}_1)^2}{2_1\sigma_x^2}} + \frac{r_2}{2\sigma_x\sqrt{2\pi}} e^{-\frac{(x-\bar{x}_2)^2}{2_2\sigma_x^2}},$$

kde čísla  $r_1$  a  $r_2$  udávají, v jakém poměru se vyskytují v celkovém rozsahu  $r$  prvky prvního a druhého souboru složkového. Konstanty vypočítáme pomocí momentů celkového rozdělení četností.

Úloha: Vypočítejte konstanty pro jednoduchý případ, kdy oba vrcholy spadají do téhož místa, takže  $\bar{x}_1 = \bar{x}_2 = \mu'_{x,1}$ .

V tomto případě je rozdělení symetrické, takže momenty lichého stupně kolem průměru se rovnají nule, tedy  $\mu_{x,1} = \mu_{x,3} = \mu_{x,5} = 0$  a ostatní jsou vzhledem k rovnicím (48) postupně

$$\begin{aligned} r &= r_1 + r_2 \\ r\mu_{x,2} &= r_{11}\sigma_x^2 + r_{21}\sigma_x^2 \\ r\mu_{x,4} &= 3(r_{11}\sigma_x^4 + r_{22}\sigma_x^4) \\ r\mu_{x,6} &= 15(r_{11}\sigma_x^6 + r_{22}\sigma_x^6). \end{aligned}$$



Řešením těchto čtyř rovnic je možno určit  $r_1, r_2, {}_1\sigma_x, {}_2\sigma_x$ , neboť  ${}_1\sigma_x^2$  a  ${}_2\sigma_x^2$  dostaneme jako dva kořeny jedné rovnice druhého stupně.

### (6,9) Příklady.

1. (Normální rozdělení četností.) Vyjádříme skupinové rozdělení četností dané v sloupci (1) a (2) pomocí křivky Laplace-Gaussovy. Použijeme k tomu částí její plochy, daných

výrazem  $F(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-t^2} dt$ , kde  $t = \frac{x - \bar{x}}{\sigma_x} = \frac{u - \bar{u}}{\sigma_u}$  vzhle-

dem k (15) a (17).

$x_i$	$n_i$	$u_i$	$u_i n_i$	$u_i^2 n_i$	$u_i - \bar{u}$
(1)	(2)	(3)	(4)	(5)	(6)
17	14	-3	-42	126	-2,095
22	121	-2	-242	484	-1,095
27	335	-1	-335	335	-0,095
32	349	0	0	0	0,905
37	150	1	150	150	1,905
42	29	2	58	116	2,905
47	2	3	6	18	$\infty$
$\Sigma$	1000		-405	1229	

$t = \frac{u_i - \bar{u}}{\sigma_u}$	$F(t)$	$\Delta F(t)$	$r \cdot \Delta F(t)$
(7)	(8)	(9)	(10)
-2,1140	1-0,9827	0,0173	17,3
-1,1049	1-0,8654	0,1173	117,3
-0,0959	1-0,5382	0,3272	327,2
0,9132	0,8195	0,3577	357,7
1,9223	0,9727	0,1532	153,2
2,9314	0,9983	0,0256	25,6
$\infty$	1,0000	0,0017	1,7
		1,0000	1000,0

$$\begin{aligned} \bar{u} &= -0,405 & \bar{x} &= 29,975 \\ \mu'_{u,2} &= 1,229 & {}^a\mu_{u,2} &= 0,982 \\ \mu_{u,2} &= 1,065 & {}^o\sigma_u &= 0,991 \end{aligned}$$

Výsledek v sl. 10. dává t. zv. teoretické rozdělení četností.

2. (Poissonova exponenciála.) Počet úmrtí žen starších 85 let, pozorovaný denně v období tří let je uveden v (1) a (2) sloupci tabulky rozdělení četností. Vzhledem k jeho nesymetrickému tvaru se pokusíme o vyjádření exponenciálou Poissonovou. Potřebujeme k tomu cíli zjistit průměr rozdělení  $\lambda$ . Celkový rozsah je  $r = 1086$ .

Počet úmrtí denně	Počet dní				
$x_i$	$n_i$	$x_i n_i$	$x_i^2 n_i$	$\psi(x)$	$r \psi(x)$
(1)	(2)	(3)	(4)	(5)	(6)
0	364	0	0	0,30360	329,7
1	376	376	376	0,36179	392,9
2	218	436	872	0,21568	234,2
3	89	267	801	0,08576	93,1
4	33	132	528	0,02559	27,8
5	13	65	325	0,00611	6,7
6	2	12	72	0,00122	1,3
7	1	7	49	0,00025	0,3
$\Sigma$	1086	1295	3023	1,00000	1086,0

$$\lambda = \bar{x} = 1,1924$$

$$\mu'_{x,2} = 2,7836$$

$$\mu_{x,2} = 1,3618$$

3. (Řada Poisson-Charlierova.) Pryskeřík (*Ranunculus*) je rostlina s korunou zpravidla pětičetnou, vyskytují se však i případy s korunou vícečetnou. Pozorováním 222 květů bylo zjištěno dole uvedené rozdělení četností. Pro jeho analytické vyjádření použijeme řady Poisson-Charlierovy (61). Potřebujeme vyčíslit koeficient druhého členu a další postup je patrný z tabulky výpočtů.

$$\lambda = \bar{u} = 0,631$$

$$\mu'_{u,2} = 1,315$$

$$\mu_{u,2} = 0,917$$

$$\frac{1}{2}(\mu_{x,2} - \lambda) = c = 0,143$$

Počet lístků v koruně	Čet- nost				
$x_i$	$n_i$	$u_i$	$u_i n_i$	$u_i^2 n_i$	$\psi(u)$
(1)	(2)	(3)	(4)	(5)	(6)
5	133	0	0	0	0,53262
6	55	1	55	55	0,33497
7	23	2	46	92	0,10588
8	7	3	21	63	0,02243
9	2	4	8	32	0,00359
10	2	5	10	50	0,00051
$\Sigma$	222		140	292	1,00000

$\Delta \psi(u)$	$\Delta^2 \psi(u)$	$c \Delta^2 \psi(u)$	$y$	$r \cdot y$
(7)	(8)	(9)	(10)	(11)
0,53262	0,53262	0,07616	0,6088	135,2
-0,19765	-0,73027	-0,10443	0,2305	51,2
-0,22909	-0,03144	-0,00450	0,1014	22,5
-0,08345	0,14564	0,02083	0,0433	9,6
-0,01884	0,06461	0,00924	0,0128	2,8
-0,00308	0,01576	0,00225	0,0028	0,7
				<u>222,0</u>

Ve sloupci (6) určena hodnota  $\psi(u)$  pro znak  $u \geq 5$ , aby se docílilo součtu relativních četností 1; difference pak byly počítány tak, jakoby to byla hodnota  $\psi(5)$ , neboť vliv této úpravy je zanedbatelný.

Úloha: Vyjádřete pomocí řady Poisson-Charlierovy rozdělení četností použité v předchozím příkladu.

### (7,1) Aplikace a zobecnění Bernoulliova teorému. (Od

Bernoulliova teorému k závěrům o skutečném průběhu jevů.) Laplaceovým integrálem jsme získali důležitý prostředek k řešení některých úloh praktické statistiky, jež se často opakují. Proto umožňuje statistikovi ohromnou úsporu práce a času. Nesmíme však nikdy zapomínati, že pro konečný rozsah souboru  $r$  znamená jen přibližnou formuli, jejíž meze chyb nelze obyčejně ani dosti přesně odhadnout. Především můžeme vhodně použít Laplaceova integrálu k takové formulaci Bernoulliova teorému, jež by usnadnila jeho praktickou aplikaci. Základní problém

Bernoulliův jest v určení pravděpodobnosti  $P_r(x)$ , že v náhodném výběru rozsahu  $r$  bude právě  $x$  prvků s pozorovaným znakem alternativním je-li  $p$  jeho relativní četnost v základním souboru. Tato pravděpodobnost je určena Newtonovou formulí (36). Můžeme pak snadno určit pravděpodobnost, že četnost  $x$  znaku ve výběru ze základního souboru o konstantní relativní četnosti  $p$  se odchýlí od průměru  $rp$  nejvýše o  $\pm \xi_0$ . Pro dosti velká  $r$  je hledaná pravděpodobnost

$$P_r(\bar{x} - \xi_0, \bar{x} + \xi_0)$$

dána Laplaceovým integrálem

$$\Phi(\gamma_0) = \frac{2}{\sqrt{\pi}} \int_0^{\gamma_0} e^{-\gamma^2} d\gamma \quad \text{pro} \quad \gamma_0 = \frac{\xi_0}{\sqrt{2rpq}}$$

Bernoulliův teorém dostaneme z toho malou změnou proměnné. Uvažujeme místo četnosti  $x$  znaku jeho relativní četnost  $f = \frac{x}{r}$  a ptáme se, jaká je pravděpodobnost, že relativní četnost znaku ve výběru má určitou hodnotu  $f$ . Označíme-li tuto pravděpodobnost  $P_r(f)$ , bude zřejmé  $P_r(f) = P_r(x)$ . Dále je  $\frac{\bar{x} - \xi_0}{r} = \frac{rp - \xi_0}{r} = p - z_0$  píšeme-li  $z_0 = \frac{\xi_0}{r}$ . Můžeme tedy udati pravděpodobnost, že relativní četnost znaku ve výběru se bude odchylovati od relativní četnosti v základním souboru nejvýše o  $z_0$ , neboť je opět dána Laplaceovým integrálem

$$P_r\left(\frac{\bar{x} - \xi_0}{r}, \frac{\bar{x} + \xi_0}{r}\right) = P_r(p - z_0, p + z_0) = \Phi\left(\frac{rz_0}{\sqrt{2rpq}}\right)$$

čili

$$P_r(p - z_0, p + z_0) = \Phi\left(z_0 \sqrt{\frac{r}{2pq}}\right). \quad (62)$$

Na pravé straně této rovnice je funkce  $z_0$ , která s rostoucím  $r$  spěje k 1 při každé hodnotě  $z_0$ , neboť  $\Phi(\infty) = 1$ . Platí tedy pro pevné  $z_0$

$$\lim_{r \rightarrow \infty} P_r(p - z_0, p + z_0) = 1, \quad (63)$$

což je výrazem Bernoulliova teorému, který znovu vyslovíme: Je-li rozsah  $r$  náhodného výběru dosti velký, je pravděpodobnost, že relativní četnost alternativního znaku v něm se odchýlí od své relativní četnosti v základním souboru o méně než  $z_0$  libovolně blízka 1, ať je  $z_0$  jakkoliv malé.

Pravděpodobnosti velmi blízké 1 se také říká méně přesně „skoro-jistota“. Potom může předcházející věta zníti: Je skoro jisto, že relativní četnost bude libovolně blízko statistické pravděpodobnosti, je-li jen  $r$  dosti velké.

Je důležité si uvědomiti, že vývody až potud byly provedeny jen matematickými úvahami z oboru t. zv. kombinatoriky. Proto nemůžeme za tohoto stavu nic říci o tom, jaká četnost znaku  $c$  by se ve skutečnosti objevila, kdybychom vzali  $r$  prvků ze základního souboru rozsahu  $N$ .

Mohli bychom dostati relativní četnost  $\frac{r}{r} = 1$ , kdybychom vzali prvky výběru z jedné části základního souboru, která má jen prvky se znakem  $c$ . Kdyby však byly všechny prvky se znakem  $c$  v jiné části základního souboru, dostali bychom při téže relativní četnosti  $p$  v základním souboru výsledek  $\frac{0}{r} = 0$  při platnosti všech formulí, jež jsme si odvodili.

Abychom mohli se svými vývody pokročiti k nějakým závěrům o vztahu mezi  $\frac{x}{r}$  a  $p$  museli jsme udělati dodatečný předpoklad, že základní soubor rozsahu  $N$  je dobře promíchán čili prvky se znakem  $c$  jsou v něm více méně stejnoměrně rozděleny. Tomu promíchání jsme rozuměli technicky, tedy asi tak jako vznikne beton pečlivým promícháním cementu, šterku, písku a vody. Podati pro pojem dobrého

promíchání ryze matematickou definicí je ovšem úkol zcela jiný. Byly takové definice sestrojeny a založen na nich celý počet pravděpodobnosti. Tak na př. Misesova definice vychází od základního souboru nekonečného rozsahu zvaného kolektiv, který má tři vlastnosti: 1. Prvky tvoří posloupnost nepravidelnou. 2. Relativní četnost  $f_i$  znaku  $c$  spěje při neomezeném počtu prvků k pevné mezní hodnotě  $p$ ; předpokládá se tedy, že existuje limita  $\lim_{i=\infty} f_i = p$ , zvaná pravděpodobnost. Říkáme, že relativní četnost  $f_i$  spěje ku  $p$  stochasticky a tato konvergence ve smyslu teorie pravděpodobnosti čili stochastická je charakterisována větou Bienaymé-Čebyševovou a Bernoulliovou. 3. V každé posloupnosti libovolně odvozené ze základní, je též mezní hodnota relativní četnosti.

Tím, že v definici pravděpodobnosti předpokládáme nějakou limitu relativní četnosti, idealisujeme pozorovanou skutečnost k účelu definice. V některých směrech je to analogická idealisace jako přímka nebo koule v geometrii či hmota a síla ve fyzice.

Přesvědčili jsme se, že lze jen matematickou cestou dospěti k závěru, že relativní četnost výběrů čili jednotlivých kombinací  $r$ -té třídy bude tím menší, čím je počet  $x$  prvků se znakem  $c$  vzdálenější od průměru  $\bar{x} = rp$ . Dokonce můžeme snadno pomocí integrálu Laplaceova uvést čísla udávající, že celková relativní četnost výběrů, u nichž je odchylka  $\xi = x - rp$  v mezích  $\pm 2\sigma_x$  je rovna 0,9545 a tedy celková relativní četnost těch případů, v nichž  $\xi$  je menší než  $-2\sigma_x$  a v nichž je větší než  $+2\sigma_x$  je rovna  $1 - 0,9545 = 0,0455$ . Relativní četnost případů, v nichž je

$|\xi| > 3\sigma_x$  je 0,0027, pro  $|\xi| > 4\sigma_x$  je 0,000063 atd.

Od těchto ryze matematických výsledků nás převádí ke skutečnostem světa nás obklopujícího věta, která je výsledkem našich zkušeností a možno říci skoro axiomem denního života, t. zv. věta Cournotova. Tato věta konsta-

tuje, že zřídka se stane, abychom vyňali náhodně (t. j. bez zvláštního vybírání a hledání) z promíchaného základního souboru prvek se znakem, jehož relativní četnost v něm, čili pravděpodobnost, je velmi malá. Kdyby v množství 10 000 zrněk hrachu bylo jedno černé, pak zřídka vyjmeme náhodně z promíchaného množství právě to černé zrnko.

Odvodili jsme pak, že větší odchylky četnosti  $x$  (resp. relativní četnosti  $\frac{x}{r}$ ) v náhodných výběrech od její hodnoty  $rp$  (resp.  $p$ ) v základním souboru mají při dostatečně velkém  $r$  velmi malou relativní četnost v souboru rozsahu  $\binom{N}{r}$  všech možných výběrů rozsahu  $r$  ze základního souboru rozsahu  $N$ . Soudíme tudíž, že ve statistických souborech dostatečně velkého rozsahu se vyskytují také ve skutečnosti jen velmi zřídka větší odchylky relativních četností od odpovídající jim statistické pravděpodobnosti. Tato věta bývá nazývána zákonem velkých čísel a je jednou z těch, které vyjadřují princip velkých čísel.

Naše vědění spočívající na principu velkých čísel není prosto jisté subjektivní libovůle. Vidíme to, chceme-li říci, do které hodnoty máme považovati relativní četnost za „velmi malou“, takže pozorovaný znak se ve skutečnosti objeví jen „velmi zřídka“ nebo jakou hodnotu relativní četnosti máme nejvýše připustiti, abychom mohli očekávati, že se pozorovaný znak „prakticky neobjeví“. Rozhodnutí nezávisí jen na absolutní velikosti pravděpodobnosti a názoru badatele, nýbrž také na stupni důležitosti, jakou by pro něho subjektivně mohl mít nepravděpodobný jev, kdyby přece nastal.

V praxi se ustálila zvyklost, že za mez pravděpodobností, na které se ještě bere zřetel, se volí celková pravděpodobnost odchylek, které přesahují  $\pm 3\sigma_x$ . Našli jsme, že pravděpodobnost odchylek větších než trojnásobek směrodatné

odchylky je dána číslem  $0,0027 = \frac{1}{3685}$ . Toto „pravidlo tří sigma“ je nyní velmi populární. Kdyby ovšem závisel náš vlastní život na vyskytnutí se jevu, který má tuto pravděpodobnost 0,3%, nezdála by se nám jistě úplně zanedbatelnou. Zavádí se také v poslední době jistý decimální systém mezních pravděpodobností a sice 5% pro optimisty, 2%, a 1% pro pesimisty.

**(7,2) Poissonovo zobecnění teorému Bernoulliiova.** Uvažujme výběr prvků se znakem alternativním, které mají různé základní pravděpodobnosti. Máme tedy  $r$  základních souborů, které mají relativní četnosti pozorovaného znaku  $c$  a doplňky na jednotku postupně  $p_1, q_1; p_2, q_2; \dots; p_r, q_r$ . Zobecnění vztahu (62), které podal Poisson (1837), spočívá v tom, že se vezme z každého z těchto základních souborů jeden prvek a určí se pravděpodobnost  $P_r(x)$ , že výběr bude mít  $x$  prvků se znakem  $c$ . Klademe-li zase

$$x = rf, \text{ je } P_r(x) = P_r(f),$$

kde  $P_r(f)$  značí pravděpodobnost, že dostaneme výběr rozsahu  $r$  s relativní četností  $f$ . Poisson dokázal, že také pro toto rozdělení  $P_r(f)$  platí vztah (63) o mezní hodnotě, v němž  $p$  musí býti průměr čísel  $p_1, p_2, \dots, p_r$ , takže v Laplaceově integrálu bude hranice

$$\gamma_0 = \frac{z_0^r}{\sqrt{2 \sum_{k=1}^r p_k (1 - p_k)}}$$

**(7,3) Průměr a rozptyl rozdělení četností vzniklého tvořením součtů z několika náhodných proměnných.**

— (Bernoulliův problém jako zvláštní případ.) V jednom základním souboru jsou hodnoty, jichž nabývá kvantitativní znak, označený čísly  $x_1, x_2, \dots, x_l$ , jimž odpovídá rozdělení relativních četností  $p_1(x_1), p_1(x_2), \dots, p_1(x_l)$ , takže  $\sum_{i=1}^l p_1(x_i) = 1$ . Tak se stává znak náhodnou proměn-



nou. Průměr je podle definice

$$x_1 \cdot p_1(x_1) + x_2(p_1(x_2) + \dots + x_l p_1(x_2) = \sum_{i=1}^l x_i p_1(x_i). \quad (64)$$

Tato hodnota, jakožto parametr základního souboru se označuje často zvláštním symbolem  $\mathfrak{E}(x)$ , jehož jsme již užili; je obdobným na př. znaménku integračnímu a odlišuje se tím jasně od průměru jako charakteristiky výběrové. Budeme ji nazývatí očekávaná hodnota; vyskytují se v počtu pravděpodobnosti také názvy střední hodnota nebo matematická naděje. V druhém základním souboru buďtež hodnoty kvantitativního znaku  $y_1, y_2, \dots, y_m$  s rozdělením četností  $p_2(y_1), p_2(y_2), \dots, p_2(y_m)$ . Očekávaná hodnota této náhodné proměnné je

$$\mathfrak{E}(y) = \sum_{i=1}^m y_i p_2(y_i). \quad (65)$$

Odvodíme nyní očekávanou hodnotu součtu dvou náhodných proměnných. Vezmeme náhodně z prvního základního souboru jeden prvek. Pravděpodobnost, že hodnota znaku bude  $x_i$  je  $p_1(x_i)$ . Obdobně bude  $p_2(y_k)$  pravděpodobnost, že z druhého základního souboru vezmeme náhodně  $y_k$ . Pravděpodobnost, že se současně vyskytne znak  $x_i$  a  $y_k$ , je podle pravidla o složené pravděpodobnosti dána součinem  $p_1(x_i) p_2(y_k) = p_{ik}$ , a to je tedy také pravděpodobnost, že dostaneme určitý součet  $x_i + y_k$ .

Pro součet obou náhodných proměnných chceme najítí očekávanou hodnotu jako průměr. Sestavíme si tedy hodnoty pravděpodobností čili relativních četností v nově vzniklém základním souboru, pro jednotlivé možné páry hodnot  $x_i$  a  $y_k$  do této tabulky

	$x_1$	$x_2$	$\dots$	$x_l$
$y_1$	$p_{11}$	$p_{21}$	$\dots$	$p_{l1}$
$y_2$	$p_{12}$	$p_{22}$	$\dots$	$p_{l2}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$y_m$	$p_{1m}$	$p_{2m}$	$\dots$	$p_{lm}$



$$\begin{aligned} \sigma^2(x+y) &= \mathfrak{E}(\xi+\eta)^2 = (\xi_1+\eta_1)^2 p_{11} + \dots + (\xi_l+\eta_m)^2 p_{lm} = \\ &= (\xi_1^2 + 2\xi_1\eta_1 + \eta_1^2) p_{11} + (\xi_1^2 + 2\xi_1\eta_2 + \eta_2^2) p_{12} + \dots \\ &\quad \dots + (\xi_1^2 + 2\xi_1\eta_m + \eta_m^2) p_{1m} + \\ &\quad \dots \\ &\quad + (\xi_l^2 + 2\xi_l\eta_1 + \eta_1^2) p_{l1} + (\xi_l^2 + 2\xi_l\eta_2 + \eta_2^2) p_{l2} + \dots \\ &\quad \dots + (\xi_l^2 + 2\xi_l\eta_m + \eta_m^2) p_{lm}. \end{aligned}$$

Sečteme opět čtverce  $\xi^2$  v každém řádku a čtverce  $\eta^2$  v každém sloupci a dostaneme

$$\begin{aligned} \sigma^2(x+y) &= \xi_1^2 p_1(x_1) + \dots + \xi_l^2 p_1(x_l) + \\ &\quad + \eta_1^2 p_2(y_1) + \dots + \eta_m^2 p_2(y_m), \end{aligned}$$

čili

$$\sigma^2(x+y) = \sigma^2(x) + \sigma^2(y),$$

neboť součet všech součinů  $\sum \xi_i \eta_k p_{ik}$  se rovná nule. O tom se snadno přesvědčíme, ježto jej dostaneme provedením součinu součtů

$$\sum_{i=1}^l \xi_i p_1(x_i) \cdot \sum_{k=1}^m \eta_k p_2(y_k)$$

a každý z těchto součtů je roven nule, neboť je to první moment kolem aritmetického průměru.

Také zde platí obecně

$$\sigma^2(x+y+z+\dots) = \sigma^2(x) + \sigma^2(y) + \sigma^2(z) + \dots \quad (67)$$

Podobně bychom odvodili

$$\sigma^2(x-y) = \sigma^2(x) + \sigma^2(y). \quad (67')$$

Bernoulliův problém se jeví jako nejjednodušší zvláštní případ tvoření součtů. Vzniká, když se pravděpodobnosti  $p_1(x_i)$ ,  $p_2(y_k)$ ,  $p_3(z_j)$ , ... vztahují na alternativu, takže se hodnoty každého znaku redukuje na dvě, jež označíme 1, 0. Potom je

$$p_1(1) = p_2(1) = \dots = p, \quad p_1(0) = p_2(0) = \dots = q.$$

Pravděpodobnost, že v náhodném výběru bude  $x$  prvků s pozorovaným znakem  $c$  je táž, jako pravděpodobnost, že

součet jednotek bude  $x$  a tedy počet nul  $r - x$ . Výsledky, jež jsme našli, odpovídají právě odvozeným větám, neboť bylo  $\mathfrak{E}(x) = rp$ ,  $\sigma^2(x) = rpq$  pro rozdělení četností  $P_r(x)$ .

Pro očekávané hodnoty platí ještě další věty, které stačí uvést:

$$\alpha) \quad \mathfrak{E}(a) = a,$$

kde  $a$  je konstanta. Z toho důvodu také

$$\beta) \quad \mathfrak{E}[\mathfrak{E}(x)] = \mathfrak{E}(x),$$

$$\gamma) \quad \mathfrak{E}(ax) = a \mathfrak{E}(x).$$

Očekávaná hodnota součinu dvou náhodných proměnných na sobě nezávislých se rovná součinu jejich očekávaných hodnot.

$$\delta) \quad \mathfrak{E}(xy) = \mathfrak{E}(x) \mathfrak{E}(y).$$

Dvě náhodné proměnné jsou na sobě nezávislé, zůstává-li rozdělení četností jedné proměnné stále totéž, ať druhá proměnná nabývá kterékoli hodnoty. Říká se také, že jsou stochasticky nezávislé.

Platí dále analogická věta jako (5) mezi obecným druhým momentem a druhým momentem kolem aritmetického průměru

$$\varepsilon) \quad \mathfrak{E}(x^2) = \mathfrak{E}(\xi^2) + [\mathfrak{E}(x)]^2.$$

**(7,4) Zákon velkých čísel.** Můžeme nyní odvodit podle Misesa další obecnou větu, která zahrnuje jako zvláštní případy teorém Bernoulliův i Poissonovo zobecnění tvořící součást vět vyjadřujících princip velkých čísel.

Odvodili jsme si očekávanou hodnotu a rozptýl rozdělení pravděpodobnosti vzniklého tvořením součtů náhodných proměnných. Hledejme nyní tyto parametry nikoliv pro součet náhodných proměnných, nýbrž pro jejich průměr. Vydeme od  $r$  základních souborů a vezmeme z každého z nich jeden prvek; budeme na nich sledovati (pro jednoduchost) znak alternativní, který bude vyznačen 1 a 0.

Celkový součet hodnot znaků bude tedy součtem jednotek na př.  $x$ . Budeme tvořit průměry tím, že součty  $x$  dělíme počtem prvků  $r$ , tedy  $f = \frac{x}{r}$ . Přejít od původních základních souborů k novému se znakem  $f$  označujeme jako tvoření průměrů.

Hledáme pravděpodobnost  $P_r(f)$ , že z  $r$  prvků bude  $x$  prvků s pozorovaným znakem, takže dostaneme z nich průměr  $f$ . Tato pravděpodobnost průměru souvisí vztahem

$$P_r(f) = P_r(x) = P_r(rf)$$

s pravděpodobností  $P_r(x)$  součtu  $x$ , jak jsme již konstatovali (str. 98).

Průměr rozdělení pravděpodobností  $P_r(f)$  je

$$\mathfrak{E}(f) = \sum_f f P_r(f) = \sum_x \frac{x}{r} P_r(x) = \frac{\mathfrak{E}(x)}{r}. \quad (68)$$

Rozptyl

$$\begin{aligned} \sigma^2(f) &= \sum_f (f - \mathfrak{E}(f))^2 P_r(f) = \sum_x \left( \frac{x - \mathfrak{E}(x)}{r} \right)^2 P_r(x) = \\ &= \frac{\sigma^2(x)}{r^2}. \end{aligned} \quad (69)$$

Značí tedy přechod od  $P_r(x)$  ku  $P_r(f)$  sesunutí úseček (obr. 12) v poměru  $r:1$ . Poněvadž pro  $r$  základních souborů jsou očekávané hodnoty  $\mathfrak{E}(x_1), \dots, \mathfrak{E}(x_r)$  a tedy podle (66)  $\mathfrak{E}(x) = \mathfrak{E}(x_1) + \mathfrak{E}(x_2) + \dots + \mathfrak{E}(x_r)$ , bude vzhledem k (68)

$$\mathfrak{E}(f) = \frac{\mathfrak{E}(x_1) + \mathfrak{E}(x_2) + \dots + \mathfrak{E}(x_r)}{r}.$$

Podobně můžeme psát vzhledem k (67) a (69) rozptyl

$$\sigma^2(f) = \frac{\sigma^2(x_1) + \sigma^2(x_2) + \dots + \sigma^2(x_r)}{r^2}.$$

Zavedeme-li předpoklad, že rozptyly  $\sigma^2(x_i)$  těch jednotlivých rozdělení četností mají horní hranici  $\sigma^2$ , že tedy  $\sigma^2(x_i) \leq \sigma^2$

pro  $i = 1, 2, \dots, r$  pak z poslední rovnice plyne, že

$$\sigma^2(f) \leq \frac{\sigma^2}{r}$$

čili  $\lim_{r \rightarrow \infty} \sigma^2(f) = 0$ .

Rozptyl rozdělení pravděpodobností  $P_r(f)$  spěje s rostoucím  $r$  k nule právě jako v případě Bernoulliově.

Pravděpodobnost, že  $f$  bude v mezích  $\pm z_0$  kolem průměru, bude vymezena zase nerovninou

$$P_r(\mathfrak{E}(f) - z_0, \mathfrak{E}(f) + z_0) \geq 1 - \frac{\sigma^2}{rz_0^2}$$

čili

$$\lim_{r \rightarrow \infty} P_r(\mathfrak{E}(f) - z_0, \mathfrak{E}(f) + z_0) = 1.$$

Můžeme tedy vyslovit větu:

Pravděpodobnost, že průměr  $r$  veličin, z nichž každá podléhá nějakému libovolnému rozdělení pravděpodobností, leží v libovolně malém intervalu u své očekávané hodnoty, je libovolně blízka 1, když  $r$  je dosti velké. Předpokladem je, že rozptyly  $\sigma^2(x_i)$  jednotlivých rozdělení mají určitou horní hranici, nebo jejich součet roste slaběji než  $r^2$ . Lze také říci stručněji: Při velkém  $r$  je skoro jisto, že průměr čísel, která podléhají nějakým  $r$  rozdělením, bude přibližně roven své očekávané hodnotě.

### **(8,1) Odhad parametrů základního souboru podle příslušných charakteristik výběrových.**

Dosud jsme se zabývali hlavně otázkou, co můžeme říci o relativní četnosti  $f$  pozorovaného znaku v náhodných výběrech, známe-li jeho relativní četnost  $p$  v příslušném základním souboru, z něhož byly vzaty. Odvodili jsme velmi užitečné věty o rozptylu alternativního znaku v náhodných výběrech.

Při statistické praxi však je častěji třeba usuzování směrem obráceným. Ze znalosti charakteristiky v jednom

nebo několika pozorovaných výběrech máme odhadnouti neznámou hodnotu příslušného parametru v základním souboru. K tomu cíli hledáme odpověď hlavně na tyto čtyři typy otázek:

1. Jaká je pravděpodobnost určité hodnoty neznámého parametru?

2. Jaký je tudíž rozptyl jeho hodnot?

3. Kterou hodnotu máme podle pozorování určitého náhodného výběru považovati za nejbližší a tedy nejlepší hodnotu neznámého parametru?

4. Lze považovati dva nebo několik souborů za náhodné výběry z téhož základního souboru?

Statistickým úkolem tedy je především, udati na základě pozorovaného výběru meze, v nichž je neznámý parametr základního souboru, čili stanoviti jeho rozptyl a najíti, kterou hodnotu lze pro tento parametr pokládati za nejlepší.

**(8,2) Meze základní relativní četnosti.** Poněvadž se v tomto oddílu zabýváme jen znakem alternativním, budeme řešiti naznačené úkoly pro relativní četnost  $f$  a ji odpovídající parametr  $p$ .

Řešení nám zase usnadní Laplaceův integrál, který udává pravděpodobnost  $\alpha(t)$ , že odchylka četnosti  $x$  od průměru  $\bar{x} = rp$  bude v mezích  $\pm t \sigma(x)$ , čili s pravděpodobností

$$\alpha(t) = \frac{2}{\sqrt{2\pi}} \int_0^t e^{-\frac{1}{2}\tau^2} d\tau \quad (70)$$

platí nerovnosti

$$-t \sigma(x) \leq x - rp \leq +t \sigma(x). \quad (71)$$

Znamená to, že v souboru, který má za prvky všechny kombinace  $r$  prvků z celkového počtu  $N$ , a má tedy rozsah  $\binom{N}{r}$ , existuje zcela určitá relativní četnost takových kombi-

nací, v nichž počet prvků s pozorovaným znakem se neodchyluje od  $rp$  více než o  $t\sigma(x)$  dolů nebo nahoru. Určitými nerovnostmi (71) je v daném souboru stanovena pravděpodobnost (70); také obráceně, předepíšeme-li si určitou pravděpodobnost, (70) plynou z ní přímo určité nerovnosti (71); Tyto nerovnosti můžeme psát také v jiném tvaru, přičteme-li na každé straně  $rp$

$$rp - t\sigma(x) \leq x \leq rp + t\sigma(x)$$

nebo

$$p - \frac{t\sigma(x)}{r} \leq \frac{x}{r} \leq p + \frac{t\sigma(x)}{r} \quad (72)$$

$$p - t\sqrt{\frac{pq}{r}} \leq f \leq p + t\sqrt{\frac{pq}{r}}.$$

Tím je tedy relativní četnost  $\frac{x}{r} = f$  sevřena do určitých mezí při daném  $p, t, r, N$ , neboť směrodatná odchylka  $\sigma(x)$  je buď  $\sqrt{r pq}$ , nebo  $\sqrt{r pq \left(1 - \frac{r}{N}\right)}$ , nevrací-li se při provádění výběru prvky do základního souboru konečného rozsahu  $N$ .

Jedná se nám nyní o to, abychom odvodili přípustné meze, v nichž musí býti  $p$  při určitém, daném  $\frac{x}{r} = f$ .

Dolní hranice (72) je  $-\frac{t\sigma(x)}{r} = f - p$  a horní hranice  $+\frac{t\sigma(x)}{r} = f - p$ . Jejich čtverec je týž, a dosadíme-li v něm za  $\sigma(x)$  druhý obecnější výraz, máme

$$t^2 p (1 - p) \left( \frac{1}{r} - \frac{1}{N} \right) = (f - p)^2.$$

To je rovnice druhého stupně pro  $p$  a jejím řešením dostáváme dva kořeny



$$p = f + \left\{ t^2 \left( \frac{1}{2} - f \right) \left( \frac{1}{r} - \frac{1}{N} \right) \pm t \right. \quad (73)$$

$$\left. \cdot \sqrt{f(1-f) \left( \frac{1}{r} - \frac{1}{N} \right) + \frac{t^2}{4} \left( \frac{1}{r} - \frac{1}{N} \right)^2} \right\} : \left\{ 1 + t^2 \left( \frac{1}{r} - \frac{1}{N} \right) \right\},$$

keré tvoří horní a dolní mez pro  $p$ . Tento výsledek může býti zjednodušen především tím, že klademe výraz v děliteli

$$\text{přibližně roven jedné, neboť vzhledem k } t = \frac{\xi}{\sqrt{r p q \left( 1 - \frac{r}{N} \right)}}$$

bude

$$1 + t^2 \left( \frac{1}{r} - \frac{1}{N} \right) = 1 + \frac{\xi^2 \frac{1}{r} \left( 1 - \frac{r}{N} \right)}{r p q \left( 1 - \frac{r}{N} \right)} = 1 + \frac{\xi^2}{r^2 p q}$$

a veličiny řádu  $\frac{\xi^2}{r^2 p q}$  jsme při odvozování křivky Gaussovy zanedbávali, takže také zde můžeme zůstat v obdobných mezích přibližnosti.

Dostali jsme tak pro relativní četnost v základním souboru nerovnosti, jimiž je sevřena při známé relativní četnosti výběrové  $f$

$$f + t^2 \left( \frac{1}{2} - f \right) \left( \frac{1}{r} + \frac{1}{N} \right) - \quad (74)$$

$$- t e \leq p \leq f + t^2 \left( \frac{1}{2} - f \right) \left( \frac{1}{r} - \frac{1}{N} \right) + t e,$$

kde

$$e = \sqrt{f(1-f) \left( \frac{1}{r} - \frac{1}{N} \right) + \frac{t^2}{4} \left( \frac{1}{r} - \frac{1}{N} \right)^2}.$$

Je-li rozsah  $r$  tak velký, že stačí přihlížeti k veličinám řádu  $\frac{1}{\sqrt{r}}$  a zanedbat veličiny řádu  $\frac{1}{r}$ , dostaneme přibližné ne-

rovnosti

$$f - t \sqrt{f(1-f) \left( \frac{1}{r} - \frac{1}{N} \right)} \leq p \leq f + t \sqrt{f(1-f) \left( \frac{1}{r} - \frac{1}{N} \right)} \quad (75)$$

a pro základní soubor nekonečného rozsahu  $N = \infty$  čili pro případ výběru s vracením prvků

$$f - t \sqrt{\frac{f(1-f)}{r}} \leq p \leq f + t \sqrt{\frac{f(1-f)}{r}}. \quad (76)$$

Je zřejmo, že nerovnosti (76) jsou inverzí nerovností (72), neboť  $p$  a  $f$  si vyměnily místo. Nerovnosti (76) tedy udávají hranice, v nichž je sevřena pravděpodobnost znaku  $p$  při dané relativní četnosti  $f$  a určité zvoleném  $t$  s pravděpodobností  $\alpha(t)$ . Velký praktický význam této inverse je v tom, že dostáváme i při neznámém  $p$  dobré přiblížení pro  $\alpha(t)$  z tabulky Laplaceova integrálu, nahradíme-li  $p$  ve výrazech pro směrodatnou odchylku hodnotou  $f$ , kterou jsme stanovili z výběru. Použijeme pak zase věty Cournotovy, abychom přešli od matematických výsledků k závěrům o skutečnosti. Nejprve si stanovíme určitou nejmenší hranici pro pravděpodobnosti, na něž ještě chceme bráti zřetel. Potom považujeme hodnoty znaku nebo odchylky, jejichž celková pravděpodobnost je menší, za „velmi zřídka se vyskytující“ nebo „prakticky se nevyskytující“. Tyto nejmenší hranice se v literatuře nazývají také „fiduciální meze“, nebo „interval konfidence“. Rozhodneme se na příklad, že nebudeme přihlížeti k pravděpodobnostem  $0,0027 = 1 - \alpha(t)$ ; tato hranice odpovídá hodnotě  $t = 3$ . Tím říkáme, že odchylky od průměru větší než  $\pm 3\sigma_x$  pozorujeme v náhodném výběru z dobře promíchaného základního souboru „velmi zřídka“ nebo „prakticky nikdy“. Potom určíme pomocí této hodnoty  $t = 3$  meze pro  $p$  v nerovnostech (74), (75) nebo (76). Konečně pak vyvodíme závěr, který je obrácením Cournotovy formulace zákona velkých čísel a odpovídá na otázku,

v jakých mezích je neznámý parametr  $p$  takto: Příhází se „velmi zřídka“ nebo „prakticky nikdy“, aby pravděpodobnost pozorovaného znaku byla vně určených hranic (pro  $t = 3$ ), byl-li vzat náhodný výběr rozsahu  $r$  ze základního souboru dobře promíchaného. Předpokladem je, že  $r$  je tak velké, že užití Laplaceovy formule je přípustné.

**(8,3) Přibližná hodnota parametru  $p$ .** Nyní máme dáti odpověď na druhou otázku: Kterou přibližnou hodnotu máme nejlépe přijmouti pro parametr  $p$ . Obyčejně se považuje za nejlepší přibližnou hodnotu pro  $p$  relativní četnost  $f$  nejčastěji pozorovaná, t. j. ta, která má v základním souboru největší relativní četnost, nebo průměr všech hodnot  $f$ , které se vyskytují. Obě cesty, o nichž se blíže zmíníme až v druhém díle, vedou zde k téměř výsledku, že za nejlepší přibližnou hodnotu parametru  $p$  bereme pozorovanou relativní četnost  $f$ . Trochu jiný výsledek dostaneme, když vyjdeme od nerovností (74), neboť tam vidíme, že se odchylky nepočítají od  $f$  nýbrž od

$$f + t^2 \left( \frac{1}{2} - f \right) \left( \frac{1}{r} - \frac{1}{N} \right),$$

což nás může vésti k hodnotě opravené druhým členem, který vymizí pro  $f = \frac{1}{2}$  a je tím větší, čím je  $f$  vzdálenější od  $\frac{1}{2}$  a čím je větší  $t$ . Tato oprava posunuje vždy přibližnou hodnotu  $f$  blíže k  $\frac{1}{2}$ . Můžeme si uvésti pro  $N = \infty$  několik čísel pro ilustraci. Pro  $r = 100$ ,  $t = \sqrt{10}$  dostáváme při pozorovaném  $f$

0,10	0,20	0,30	0,40	0,50	0,60	0,70	0,80	0,90
opravenou přibližnou hodnotu pro parametr $p$								
0,14	0,23	0,32	0,41	0,50	0,59	0,68	0,77	0,86.

Třebaže nemůžeme stanoviti jednoznačně přibližnou hodnotu pro  $p$ , poněvadž stojí v cestě obtíže vyplývající z povahy problému, přece je opravdovým úspěchem matema-

tické teorie, že můžeme dosáci za určitých podmínek velmi cenných odhadů parametru tím, že lze zkoumati a odhadnouti rozptyl resp. směrodatnou odchylku hodnot, z nichž jsme jednu zjistili náhodným výběrem.

Očekávaná hodnota  $\mathfrak{E}(f)$  relativní četnosti  $f = \frac{x}{r}$ , která je průměrem jejích hodnot, zjištěných ve všech možných výběrech rozsahu  $r$  je rovna příslušnému parametru  $p$  v základním souboru. Jako přibližnou jeho hodnotu dostáváme relativní četnost  $\frac{x}{r}$  z pozorovaného náhodného výběru, která je mu tím bližší, čím je  $r$  větší. Potřebujeme tedy vyjádřit očekávanou hodnotu rozptylu  $\mathfrak{E}(\sigma^2)$  pomocí pozorovaných hodnot přibližných. Víme již (str. 69), že očekávaná hodnota rozptylu relativních četností je  $\frac{pq}{r}$ , jakožto hodnota rozptylu v základním souboru. Známe však pro  $p$  a  $q$  jen přibližné hodnoty  $f$  a  $1 - f$ . Nemůžeme vzíti za přibližnou hodnotu rozptylu jednoduše  $\frac{1}{r} f(1 - f)$ , která by vyplývala, kdybychom kladli za  $p$  přibližnou hodnotu  $f$ . O tom se přesvědčíme, když si vypočítáme, jaká by byla očekávaná hodnota součinu  $\frac{x}{r} \left(1 - \frac{x}{r}\right)$  daného výsledkem pozorování. Stanovíme tedy očekávanou hodnotu výrazu  $\frac{x}{r} - \frac{x^2}{r^2} = \frac{x}{r} \left(1 - \frac{x}{r}\right)$ . Očekávaná hodnota  $\frac{x}{r}$  je  $p$ ; očekávaná hodnota  $\frac{x^2}{r^2}$  je podle věty ( $\gamma$ ) rovna  $\frac{pq}{r} + p^2$ , neboť druhý moment kolem průměru je  $\frac{pq}{r}$  a čtverec očekávaného průměru, který je totožný s průměrem v základním souboru je  $p^2$ .

Bude tudíž celková očekávaná hodnota uvažovaného součinu podle (66)

$$p - \frac{pq}{r} - p^2 = p(1-p) - \frac{pq}{r} = pq \left(1 - \frac{1}{r}\right)$$

a je odlišná od součinu  $pq$ . Kdybychom přijali  $\frac{x}{r} \left(1 - \frac{x}{r}\right)$  za přibližnou hodnotu očekávané hodnoty  $pq$ , dopouštěli bychom se tedy jednak chyby systematické, jež se jeví v součiniteli  $\frac{r-1}{r}$ , jednak druhé chyby v tom, že existuje

odchylka mezi zvláštní pozorovanou hodnotou  $\frac{x}{r}$  náhodné proměnné a její očekávanou hodnotou  $p$ . Systematickou chybu můžeme opravit tím, že vezmeme za přibližnou hodnotu  $\frac{r}{r-1} \frac{x}{r} \left(1 - \frac{x}{r}\right)$ , neboť její očekávaná hodnota je pak právě  $pq$ . Z toho tedy vyplývá, že

$$\frac{1}{r} \frac{x}{r} \left(1 - \frac{x}{r}\right) = \frac{r-1}{r} \frac{pq}{r}. \quad (77)$$

Další otázkou, kdy lze považovati dva nebo více souborů za náhodné výběry z téhož základního souboru, budeme se zabývat v druhém díle. Zde se omezíme jen na konvenci, která se ujala dnes ve statistice. Dostaneme-li pro jednu charakteristiku, v našem případě pro relativní četnost ze dvou různých výběrů hodnoty  $f_1$  a  $f_2$ , považujeme souhlas mezi nimi za dobrý, když rozdíl  $|f_1 - f_2|$  je menší než směrodatná odchylka této difference podle (67') tedy  $\sqrt{\sigma_{f_1}^2 + \sigma_{f_2}^2}$ , a za uspokojivý, je-li menší než dvojnásobek, někdy i trojnásobek jeho směrodatné odchylky.

Přesahuje-li rozdíl trojnásobek směrodatné odchylky, nepovažuje se souhlas za uspokojivý a vzniká domněnka, že lze najíti vysvětlení této odchylky zvláštní příčinou, nikoliv náhodným výběrem.

(8,4) Pearsonovo kritérium  $\chi^2$ . Podle výsledků, jež jsme dosud odvodili, můžeme stanovit meze, v nichž je relativní četnost v základním souboru  $p$  sevřena, zvolíme-li si za přípustný interval odchylek délku  $\pm 3\sigma_x$ . Tak pro tabulku I. našeho příkladu (str. 29) máme pro hodnotu třídního znaku  $x_i = 75$  relativní četnost  $f_i = 0,141$ , takže směrodatná odchylka  $\sigma_i = 0,021$  a tudíž meze jsou  $f_i \pm 0,063$ . Tak si můžeme vypočítati meze pro parametr  $p_i$  každé třídy z pozorovaného rozdělení četností.

Klademe si však dále otázku, jak bychom vystihli, do jaké míry se liší rozdělení pozorovaného souboru jako celek od základního souboru, nikoliv jak se liší jednotlivé četnosti od příslušných parametrů. V odpověď na to sestrojil K. Pearson t. zv. kritérium  $\chi^2$ .

V základním souboru jsou statistické pravděpodobnosti hodnot znaku kvantitativního resp. třídních hodnot znaku  $p_1, p_2, \dots, p_i$ . Výběr rozsahu  $r$ , který by měl tytéž relativní četnosti, by vykazoval třídní četnosti  $v_1 = rp_1, \dots, v_i = rp_i$ . Tyto četnosti porovnáváme s pozorovanými  $rf_i$  tak, že tvoříme jejich rozdíly; čtverce rozdílů pak vyjádříme v poměru k teoretickým četnostem  $v_i$  a sečteme. Tak dostaneme výraz

$$\chi^2 = \sum_{i=1}^l \frac{(rf_i - rp_i)^2}{rp_i}$$

Všechny čtverce rozdílů se sčítají a je zřejmo, že čím jsou rozdíly obojích četností větší, tím je větší  $\chi^2$ ; jsou-li obě rozdělení shodná, je  $\chi^2 = 0$ . Uvedený výraz můžeme také psáti v tvaru

$$\chi^2 = \sum_{i=1}^l r \frac{(f_i - p_i)^2}{p_i} = r \left( \sum_{i=1}^l \frac{f_i^2}{p_i} - 1 \right)$$

neboť

$$\sum_{i=1}^l r \left( \frac{f_i^2}{p_i} - 2f_i + p_i \right) = r \sum_{i=1}^l \frac{f_i^2}{p_i} - 2r + r$$

vzhledem k tomu, že

$$\sum_{i=1}^l f_i = \sum_{i=1}^l p_i = 1.$$

Vidíme, že je to charakteristika, vztahující se k určitému výběru rozsahu  $r$ , která nám podává zhuštěnou informaci o tom, jak se tento výběr v celku liší svým rozdělením četností od očekávaného. Pro každý výběr bychom dostali pravděpodobně jinou hodnotu, takže ze všech  $\binom{N}{r}$  hodnot bude utvořeno rozdělení četností této charakteristiky. K tomuto rozdělení četností se utvoří součtová křivka  $F_1(\chi^2)$  integrací obdobně jako jsme dostali Laplaceův integrál (53) nebo (50), která však závisí ještě na druhé veličině  $l - 1$ , kde  $l$  je počet tříd.

Z ní se tedy dovíme, jaká je pravděpodobnost, že při určitém  $r$  a daných hodnotách  $p_i$  dostaneme větší hodnotu pro  $\chi^2$ , než je pozorovaná. Je to tedy pravděpodobnost, s níž můžeme očekávat horší souhlas s teoretickým rozdělením, než je pozorovaný.

Presvědčíme se, jak vystihuje v příkladu 3. str. 97 teoretické rozdělení četností pomocí dvou členů řady Poisson-Charlierovy rozdělení pozorované tím, že vypočítáme charakteristiku  $\chi^2$ .

$x_i$	$n_i$	$rp_i$	$ n_i - rp_i $	$(n_i - rp_i)^2$	$\frac{(n_i - rp_i)^2}{rp_i}$
5	133	135,2	2,2	4,84	0,04
6	55	51,2	3,8	14,44	0,28
7	23	22,5	0,5	0,25	0,01
8	7	9,6	2,6	6,76	0,70
9	2	2,8	0,8	0,64	0,23
10	2	0,7	1,3	1,69	2,41
$\Sigma$	222	222,0			3,67

Vidíme, že  $\chi^2 = 3,67$  a podle příslušné tabulky Eldertovy mu odpovídá pro  $l - 1 = 5$   $F_1(\chi^2) = 0,60$ , což je pravděpodobnost, že dostaneme v náhodných výběrech větší hodnoty  $\chi^2$ , než je pozorovaná; bylo by to tedy přibližně v 60 případech ze 100. Takové vystižení není příliš dobré. Ovšem v tomto případě se uplatňuje příliš vliv posledních dvou málo obsazených tříd; příslušná čísla  $rp_i$  ve jmenovateli pak příliš zvyšují hodnotu  $\chi^2$ , jak vidíme na poslední třídě a není to tedy jen vlivem rozdílů. Proto a také z důvodů spočívajících v odvození, jež předpokládá, že odchylky od očekávaných četností vyhovují normální křivce, se obyčejně krajní třídy málo obsazené spojují dohromady, aby četnost byla aspoň 5.

Spojíme-li tedy poslední dvě třídy, bude pak  $\chi^2 = 1,1$  a jemu odpovídá pro  $l - 1 = 4$  pravděpodobnost  $F_1(\chi^2) = 0,78$ . Z toho můžeme usuzovati, že bychom dostali přibližně v 78 případech náhodných výběrů ze sta řadu pozorovaných četností, jež dává skupinu odchylek od teoretického rozdělení tedy  $\chi^2$ , jež je méně pravděpodobné než pozorované; měli bychom tedy očekávat zhruba v každém stu náhodných výběrů 78krát horší souhlas s teoretickým, než je pozorovaný, vyjádřený charakteristikou  $\chi^2 = 1,1$ .

**(8,5) Příklad.** 1. Roční míra úmrtnosti 60letých osob byla v nějakém rozsáhlém souboru zjištěna a uvedena v tabulce úmrtnosti  $q_{60} = 0,0287$ . Jaká je pravděpodobnost náhodné odchylky menší než  $\pm z_0 = 0,01$  roční míry pozorované v souboru rozsahu  $r = 2500$  a menší než  $\pm z_0 = 0,005$  v souboru rozsahu  $r = 10\ 000$ .

Tuto pravděpodobnost udává Laplaceův integrál, jehož horní mez určíme ze vztahu  $\gamma_0 = z_0 \sqrt{\frac{r}{2pq}} = z_0 \sqrt{\frac{r}{2f(1-f)}}$ , kde bude  $z_0 = 0,01$ ,  $r = 2500$ ,  $f = 0,0287$ ,  $1 - f = 0,9713$ . Tak dostáváme  $\gamma_0 = 1,339$  a tedy  $\Phi(\gamma_0) = 0,997$  a stejnou hodnotu máme v druhém případě.

2. Mezi 1 359 671 narozenými chlapci bylo 58 744 mrtvě narozených a mezi 1 285 086 děvčaty 44 224 mrtvě narozených. Vypočítejte podle pohlaví procento mrtvě narozených, které je





Gaussovy je překročen s pravděpodobností 0,0027, nepovažuje se v praxi za odchylku nahodilou.

6. Ve sklárně se zjistí, že automat na výrobu lahví dal při přejímací zkoušce 2% vadných lahví ze 4000 kusů udělaných při zkoušce. Jaké bude asi mezní procento výmětů při plynulé výrobě? Vezmeme tedy za přibližnou hodnotu očekávané hodnoty  $f = 0,98$ ,  $1 - f = 0,02$ , takže

$$\sigma_f = \sqrt{0,02 \times 0,98 : 4000} = 0,223 \text{ procent,}$$

takže mez pro výměty je pravděpodobně  $2 + 3 \times 0,223 = 2,7$  procent.

7. Určité křížení hrachu dalo 5321 žlutých a 1804 zelená zrnka. Podle hypotese Mendelovy je očekávaný počet zelených zrněk 25%. Lze považovati tuto odchylku od očekávané hodnoty za vzniklou jen náhodným výběrem?

Odchylka pozorovaného výsledku od očekávaného je  $\xi = 23$ . Směrodatná odchylka  $\sigma_f = \sqrt{0,25 \times 0,75 \times 7125} = 36,6$ . Poněvadž odchylka  $\xi$  je jen asi  $0,6\sigma_f$ , mohla vzniknouti zcela dobře jen náhodným výběrem.

8. Za víceleté období se objevil ve statistice dětských sebevražd roční průměr  $\bar{x} = 1,96$ . Můžeme uvést jako příklad použití Poissonovy exponentiely tento jev a nepotřebujeme znáti explicitně  $r$  a  $p$ . Tak dostaneme pravděpodobnost, že se v jednom roce nevyskytne ani jedna sebevražda, pak že se vyskytne v jednom roce jedna, dvě, atd. Z rovnice (54) dostáváme  $\psi(0) = 0,141$ ,  $\psi(1) = 0,276$ ,  $\psi(2) = 0,271$ ,  $\psi(3) = 0,177$ ,  $\psi(4) = 0,087$ . Součet těchto pravděpodobností je 0,952, takže na všechny ostatní dohromady zbývá 4,8%, což je pravděpodobnost více než čtyř případů dětských sebevražd v roce. Největší pravděpodobnosti mají případy  $x = 1$  a  $x = 2$ , mezi nimiž leží průměr, a to blíže k  $x = 2$  hlavně v důsledku toho, že  $\psi(3)$  je větší než  $\psi(0)$ .

**(9,1) Lexisova teorie.** Všimli jsme si, že pro teorém Bernoulliův a teorii s ním související až na zobecnění Poissonovo je podstatným znakem, že pravděpodobnost  $p$ , která je podkladem relativních četností získaných pozorováním, je konstantní. Pozorované statistické soubory bývají složeny z prvků mnohotvárnějších a složitějších než odpovídá schématu Bernoulliovu. Zakladatelé matematické statistiky, i Laplace, považovali totožnost pozorované statistické

řady s řadou Bernoulliiovou za samozřejmou. Teprve Lexis ukázal nepostačitelnost dosavadních úvah a podal jasnější pohled na povahu statistických řad. Používání směrodatné odchylky (40) pro rozbor pozorovaných řad dává příliš hrubé výsledky, které jsou tím vzdálenější od skutečnosti, čím její podklad se více liší od podkladu typické binomické řady. Kolísání numerických hodnot pozorovaného znaku na prvcích souboru se neřídí jednoduchými zákony jako schema Bernoulliovo, působí-li na statisticky studovaný jev rušivé vnější vlivy, a proto potřebujeme míru k hodnocení zjištěných rozdílů. Tuto míru dává Lexisova teorie řad. Lexis a současně Dormoy, formuloval otázku, jak určit míry podobnosti nebo rozdílu mezi strukturou statistické řady pozorované a příslušné binomické.

Tomuto určení slouží srovnávání rozptylů resp. směrodatných odchylek řad, s nimiž se potkává statistická praxe; k němu užívá Lexisova teorie tří typů statistických řad jako norem. Metoda rozboru pak spočívá v tom, že pozorovaný soubor se rozloží na částečné soubory, v nichž by mohly býti zkoumány změny relativní četnosti znaku. Hledisko pro odvození těchto částečných souborů není dáno jen všeobecnými zásadami, nýbrž uplatněním statistických zkušeností a znalostí vědního oboru, do něhož spadá studium pozorovaného souboru, jakož i podrobné znalosti původního materiálu a jeho pramenů. K objevení a vysvětlení podstatných změn lze pak proniknouti především statistickým uměním, které pomáhá zvoliti vhodný vědecký postup. Tyto všeobecné úvahy dále objasníme na pozorovaném materiálu. Nyní se seznámíme s uvedenými třemi typy řad, odpovídajícími jednoduchým schematům náhodných her. Srovnáváním s nimi je osvětlována náhodná stránka ve statistickém dění.

1. S prvním typem řad jsme se již seznámili. Je představován řadou, jejíž základní pravděpodobnost  $p$  výskytu pozorovaného znaku je konstantní a nazývá se řadou Bernoulliiovou. Její očekávaná hodnota průměru podle (38)

je  $\mathfrak{E}(x) = rp$  a očekávaná hodnota směrodatné odchylky (40) teoretického rozdělení četností byla odvozena ve výrazu  $\sigma(x) = (rpq)^{\frac{1}{2}}$ ; očekávaná hodnota směrodatné odchylky příslušného rozdělení relativních četností je dána výrazem  $\left(\frac{pq}{r}\right)^{\frac{1}{2}}$  a očekávaná hodnota průměru je  $p$ .

Uvažme nyní, že neznáme hodnotu pravděpodobnosti  $p$ , nýbrž jen pozorované hodnoty relativních četností  $f_i = \frac{x_i}{r}$  z  $n$  výběrů rozsahu  $r$  prvků.

Musíme pak vzít za přibližnou hodnotu parametru  $p$  zlomek, který je průměrem pozorovaných hodnot  $f_i$

$$\bar{f} = \frac{1}{n} (f_1 + f_2 + \dots + f_n) = \frac{1}{rn} (x_1 + x_2 + \dots + x_n)$$

a za této hypotézy je očekávaná hodnota  $\mathfrak{E}(\bar{f}) = p$  a také očekávaná hodnota každé jednotlivé relativní četnosti  $\mathfrak{E}(f_i) = p$ . Dále očekávaná hodnota

$$\mathfrak{E}(f_i - p)^2 = \frac{pq}{r}, \quad (78)$$

neboť je to průměr čtverců odchylek relativních četností z výběru rozsahu  $r$  od jejich průměru, čili rozptyl vyjádřený pomocí hodnot základního souboru.

Dále je

$$\mathfrak{E}(\bar{f} - p)^2 = \frac{1}{n^2} \mathfrak{E} \left[ \sum_{i=1}^n (f_i - p) \right]^2 = \frac{pq}{r \cdot n} \quad (79)$$

vzhledem k tomu, že

$$\begin{aligned} \bar{f} - p &= \frac{1}{n} [(f_1 - p) + (f_2 - p) + \dots + (f_n - p)] = \\ &= \frac{1}{n} \sum_{i=1}^n (f_i - p) \end{aligned} \quad (80)$$

a očekávané hodnoty součinů  $\mathfrak{E}(f_i - p)(f_j - p)$ , kde  $i \neq j$

jsou rovny nule, ježto očekávaná hodnota každého činitele se rovná nule. Zbývá tedy  $n$  čtverců a očekávaná hodnota každého z nich je podle (78) rovna  $\frac{pq}{r}$ . Z toho je patrné, že každá relativní četnost  $f_i$  je přibližnou hodnotou s rozptylem  $\frac{pq}{r}$  a jejich průměr  $\bar{f}$  je přibližnou hodnotou, která je bližší ve smyslu teorie pravděpodobnosti s menším rozptylem  $\frac{pq}{r} \cdot \frac{1}{n}$ .

Statistická řada pozorovaných relativních četností má tedy projevovat rozptyl  $\frac{pq}{r}$  kolem hodnoty základního souboru  $p$ .

Známe však jen přibližnou hodnotu  $\bar{f}$  parametru  $p$ , takže musíme zkoumati rozptyl pozorovaných relativních četností  $f_i$  kolem jejich průměru  $\bar{f}_i$ ; při tom musíme mít na paměti, že budou hráti svoji roli uvedené již dvě odchylky (str. 115).

Abychom stanovili očekávanou hodnotu tohoto rozptylu, vypočítáme si nejprve očekávanou hodnotu čtverců odchylek od přibližného průměru  $\bar{f}$ , takže

$$\begin{aligned} \mathfrak{E}(f_i - \bar{f})^2 &= \mathfrak{E}[(f_i - p) - (\bar{f} - p)]^2 = \\ &= \frac{pq}{r} + \frac{pq}{rn} - 2\mathfrak{E}(f_i - p)(\bar{f} - p), \end{aligned}$$

neboť očekávané hodnoty čtverců známe podle rovnic (78) a (79) a očekávanou hodnotu posledního součinu stanovíme za předpokladu, že výběry jsou na sobě nezávislé, takže dosadíme-li tam z rovnice (80) vidíme, že  $n - 1$  očekávaných hodnot součinů, kde  $i \neq j$ , je rovno nule a zůstává jediný  $\frac{1}{n} (f_i - p)^2$ , jehož očekávaná hodnota je  $\frac{pq}{rn}$ .

Je tedy

$$\mathbb{E}(f_i - \bar{f})^2 = \frac{pq}{r} \left(1 - \frac{1}{n}\right).$$

Poněvadž rozptyl řady empirických hodnot kolem jejich průměru je  $\frac{1}{n} \sum_{i=1}^n (f_i - \bar{f})^2$ , bude také

$$\mathbb{E} \frac{1}{n} \sum_{i=1}^n (f_i - \bar{f})^2 = \frac{pq}{r} \left(1 - \frac{1}{n}\right). \quad (81)$$

Vzhledem k (78) je zřejmě  $\frac{pq}{r}$  očekávanou hodnotou průměru čtverců odchylek pozorovaných relativních četností od  $p$ .

$$\mathbb{E} \frac{1}{n} \sum_{i=1}^n (f_i - p)^2 = \frac{pq}{r}.$$

To je rozptyl řady Bernoulliovy, který budeme označovat  $\sigma_B^2$ . Pro přibližnou hodnotu  $\sigma_B^2$  tedy stačí vzhledem k rovnici (81) brátí výraz

$$\frac{1}{n-1} \sum_{i=1}^n (f_i - \bar{f})^2,$$

jehož očekávaná hodnota je právě  $\sigma_B^2$ .

Máme tudíž dva výrazy pro přibližnou hodnotu rozptylu  $\sigma_B^2$ . Jednak tvoříme t. zv. hodnotu počítanou  $\frac{\bar{f}(1-\bar{f})}{r} = \sigma_f^2$ , jednak t. zv. hodnotu měřenou  $\sigma_f^2$ . Za předpokladu stálého složení v základním souboru a nezávislosti prvků náhodně vybíraných, mohou se tyto dvě hodnoty málo lišit, takže jejich podíl (nebo jeho odmocnina)

$$Q^2 = \frac{\sigma_f^2}{\sigma_f^2}$$

musí býti blízký jednotce. Říkáme pak, že statistická řada

má normální rozptyl, je-li  $Q = 1$ . Pozorování je pak reprezentováno schematem Bernoulliovým, jestliže seskupení relativních četností  $f_i$  kolem jejich průměru  $\bar{f}$  odpovídá binomickému rozdělení.

Můžeme jej také vyjádřit podrobněji

$$Q^2 = \frac{1}{n-1} \sum_{i=1}^n (f_i - \bar{f})^2 : \frac{\bar{f}(1-\bar{f})}{r} = \sigma_f^2 : \frac{n-1}{n} \frac{\bar{f}(1-\bar{f})}{r}, \quad (82)$$

kde

$$\sigma_f^2 = \frac{\sum_{i=1}^n (f_i - \bar{f})^2}{n}.$$

2. Druhý typ si přiblížíme představou  $n$  zalidněných okresů, v nichž pozorujeme úmrtnost  $x$ -letých (na př. 30letých) mužů; tato pravděpodobnost je v každém okresu jiná, ale konstantní. Příklad si znázorníme modelem, sestrojeným z  $n$  osudí  $O_1, O_2, \dots, O_n$ . Stálá pravděpodobnost vytažení černé kuličky z osudí  $O_1$  budiž  $p_1, \dots$ , a z  $O_n$  budiž  $p_n$ .

$$\begin{array}{c|c} O_1 & \overbrace{p_1 \ p_1 \ \dots \ p_1}^r \\ O_2 & p_2 \ p_2 \ \dots \ p_2 \\ \vdots & \vdots \ \vdots \ \vdots \\ O_n & p_n \ p_n \ \dots \ p_n \end{array}$$

Z každého osudí vytáhneme  $r$  kuliček; očekávaný průměr počtu černých kuliček z  $i$ -tého osudí je tedy  $rp_i$ . Označme průměr pravděpodobností  $p = \frac{p_1 + p_2 + \dots + p_n}{n}$ .

Vezmeme-li z každého osudí náhodný výběr  $r$  kuliček, bude celkový očekávaný průměr počtu černých kuliček mezi  $nr$  vytaženými  $rp_1 + rp_2 + \dots + rp_n = nrp$ .

Jestliže je  $nrp$  očekávaný průměr počtu černých kuliček v  $nr$  tazích, je  $rp$  očekávaný průměr v  $r$  tazích, jež učiníme vždy z jednoho osudí náhodně vybraného. Tato hodnota

očekávaného průměru je totožná s očekávaným průměrem počtu černých kuliček ve výběrech  $r$  kuliček řady Bernoulliovy s konstantní pravděpodobností  $p$ .

Uvažujme nyní, jak velký bude rozptyl. Rozptyl ve výběru  $r$  kuliček z osudí  $O_i$ , kde pravděpodobnost černé je  $p_i$ , je dán výrazem  $rp_iq_i$ . Je to průměr čtverců odchylek od průměru výběrového z osudí  $O_i$ , jenž je  $rp_i$ . Hledáme však průměrnou čtvercovou odchylku od hodnoty  $rp$  místo od výběrového průměru  $rp_i$ . Chceme tedy stanovit obecný druhý moment kolem počátku v  $rp$ , který se podle rovnice (5) rovná druhému momentu kolem aritmetického průměru zvětšenému o čtverec rozdílu mezi průměrem a zvoleným počátkem. Je tudíž dán výrazem  $rp_iq_i + (rp_i - rp)^2$ .

Kdybychom vzali z jednoho osudí  $O_i$  takových náhodných výběrů na př.  $N$ , byl by ovšem očekávaný průměr součtu čtverců odchylek od  $rp$  větší  $N$ -krát, tedy

$$Nrp_iq_i + Nr^2(p_i - p)^2. \quad (83)$$

Utvoříme součet výrazů (83) pro všechna osudí, pak dostaneme

$$Nr \sum_{i=1}^n p_iq_i + Nr^2 \sum_{i=1}^n (p_i - p)^2, \quad (84)$$

což je očekávaný průměr součtu čtverců odchylek od  $rp$  pro  $n$  osudí, z každého z nichž jsme vzali  $N$  náhodných výběrů rozsahu  $r$  kuliček.

Celkem máme  $Nn$  výběrů a poněvadž hledáme průměrnou čtvercovou odchylku od hodnoty  $rp$ , připadající na jeden výběr, kterou označíme  $S_L^2$ , musíme dělit jejich počtem poslední výraz (84), čímž dostaneme

$$S_L^2 = \frac{r}{n} \sum_{i=1}^n p_iq_i + \frac{r^2}{n} \sum_{i=1}^n (p_i - p)^2.$$

Součet v prvním členu na pravé straně rovnice však můžeme upravit položíme-li  $p_i = p + (p_i - p)$ , a vzhledem



k  $p_i + q_i = 1$  tedy  $q_i = q - (p_i - p)$ . Potom součin

$$p_i q_i = pq - (p_i - p)(p - q) - (p_i - p)^2$$

a součet jejich

$$\sum_{i=1}^n p_i q_i = npq - \sum_{i=1}^n (p_i - p)^2, \quad (85)$$

ježto

$$(p - q) \sum_{i=1}^n (p_i - p) = 0,$$

vzhledem k tomu, že

$$\sum_{i=1}^n (p_i - p) = 0,$$

neboť

$$p_1 + p_2 + \dots + p_n = np.$$

Na základě (85) bude tedy

$$S_L^2 = rpq + \frac{r^2 - r}{n} \sum_{i=1}^n (p_i - p)^2. \quad (86)$$

Označíme-li  $S_B^2$  rozptyl výběru o rozsahu  $r$  z hypotetického souboru spočívajícího na schematu Bernoulliově s konstantní pravděpodobností  $p$ , která se rovná průměru daných pravděpodobností  $p_1 + p_2 + \dots + p_n$ , můžeme poslední rovnici psát

$$S_L^2 = S_B^2 + \frac{r^2 + r}{n} \sum_{i=1}^n (p_i - p)^2$$

a vidíme, že rozptyl Lexisovy řady je větší než řady Bernoulliovy, spočívající na pravděpodobnosti  $p$ .

Příslušný výraz pro Lexisovy řady relativních četností dostaneme dělením pravé strany rovnice (86) čtvercem rozsahu výběru  $r^2$ , takže potom

$$\sigma_L^2 = \sigma_B^2 + \frac{1 - \frac{1}{r}}{n} \sum_{i=1}^n (p_i - p)^2;$$

pro velká  $r$  pak se užívá přibližně

$$\sigma_L'^2 = \sigma_B'^2 + \frac{1}{n} \sum_{i=1}^n (p_i - p)^2 \quad (87)$$

a píšeme-li

$$\frac{1}{n} \sum_{i=1}^n (p_i - p)^2 = \sigma_p^2,$$

bude

$$\sigma_L^2 = \sigma_B^2 + \sigma_p^2. \quad (88)$$

Druhý člen na pravé straně této rovnice se často nazývá podstatnou komponentou kolísání. Jiný způsob výkladu podává analýza rozptylu, o níž pojednáme později.

Při pozorovaných statistických řadách bude tedy směrodatná odchylka větší než směrodatná odchylka vypočtená z průměrné relativní četnosti znaku, která vystihuje náhodné kolísání bez vlivů rušivých.

Schema osudí s různým složením nám tedy zobrazilo rozptyl řad, který se tím více liší od normálního, čím se základní pravděpodobnosti těchto osudí od sebe více liší.

Rozptyl  $\sigma_L^2$  však nemůžeme podle rovnice (87) počítati, ježto neznáme pravé hodnoty  $p_i$  resp.  $p$ , které se v ní vyskytují a musíme užiti přibližné hodnoty rozptylu

$$\sigma_f^2 = \frac{\sum_{i=1}^n (f_i - \bar{f})^2}{n}.$$

Stanovme tedy její očekávanou hodnotu.

Průměru  $p = \frac{1}{n} \sum_{i=1}^n p_i$  odpovídá empiricky stanovený

průměr  $\bar{f} = \frac{1}{n} \sum_{i=1}^n f_i$ . Zavedeme si k řešení našeho úkolu identitu

$$f_i - \bar{f} = (f_i - p_i) + (p_i - p) - (\bar{f} - p),$$

takže její čtverec bude

$$(f_i - \bar{f})^2 = (f_i - p_i)^2 + (p_i - p)^2 + (\bar{f} - p)^2 + \\ + 2(f_i - p_i)(p_i - p) - 2(f_i - p_i)(\bar{f} - p) - \\ - 2(p_i - p)(\bar{f} - p).$$

Abychom stanovili očekávanou hodnotu  $\mathfrak{E}(f_i - \bar{f})^2$  uvědomíme si, že

$\mathfrak{E}(f_i) = p_i$ ,  $\mathfrak{E}(f_i - p_i) = 0$ ,  $\mathfrak{E}(\bar{f}) = p$ ,  $\mathfrak{E}(\bar{f} - p) = 0$ ,  
podle rovnice (78) je

$$\mathfrak{E}(f_i - p_i)^2 = \frac{p_i q_i}{r}.$$

Poněvadž  $\bar{f} - p = \frac{1}{n} \sum_{i=1}^n (f_i - p_i)$ , je dále

$$\mathfrak{E}(f_i - p_i)(\bar{f} - p) = \mathfrak{E}\left[\frac{1}{n} (f_i - p_i) \sum_{i=1}^n (f_i - p_i)\right] = \\ = \mathfrak{E}\frac{1}{n} (f_i - p_i)^2 = \frac{p_i q_i}{rn},$$

neboť očekávaná hodnota ostatních  $n - 1$  členů, v nichž se indexy liší, se rovná nule.

$$\mathfrak{E}(\bar{f} - p)^2 = \frac{1}{n^2} \mathfrak{E}\left[\sum_{i=1}^n (f_i - p_i)\right]^2 = \\ = \frac{1}{n^2} \left[ \mathfrak{E} \sum_{i=1}^n (f_i - p_i)^2 + 2 \mathfrak{E} \sum_{i=1, j=1}^n (f_i - p_i)(f_j - p_j) \right],$$

kde  $i \neq j$ ; členy, kde by bylo  $i = j$ , se v tomto druhém součtu nevyskytují.

Očekávaná hodnota jednotlivých členů v druhém součtu se rovná nule a tedy i celého součtu, takže

$$\mathfrak{E}(\bar{f} - p)^2 = \frac{1}{n^2} \mathfrak{E} \sum_{i=1}^n (f_i - p_i)^2 = \frac{1}{n^2} \sum_{i=1}^n \frac{p_i q_i}{r}.$$

Můžeme tudíž psáti celkem

$$\begin{aligned} \mathfrak{E}(f_i - \bar{f})^2 &= \frac{p_i q_i}{r} + (p_i - p)^2 + \frac{1}{n^2} \sum_{i=1}^n \frac{p_i q_i}{r} - 2 \frac{p_i q_i}{r}, \\ \mathfrak{E} \sum_{i=1}^n (f_i - \bar{f})^2 &= \frac{\sum_{i=1}^n p_i q_i}{r} + \sum_{i=1}^n (p_i - p)^2 + \\ &+ \frac{n}{n^2} \frac{\sum_{i=1}^n p_i q_i}{r} - \frac{2 \sum_{i=1}^n p_i q_i}{rn} = \\ &= \frac{n-1}{n} \frac{\sum_{i=1}^n p_i q_i}{r} + \sum_{i=1}^n (p_i - p)^2. \end{aligned}$$

Tento výsledek ještě můžeme upravit, píšeme-li  $p_i = p + (p_i - p)$ , takže potom  $q_i = q - (p_i - p)$ , a součet je jako na str. 127.

$$\sum_{i=1}^n p_i q_i = npq - (p - q) \sum_{i=1}^n (p_i - p) - \sum_{i=1}^n (p_i - p)^2,$$

ježto prostřední člen se rovná nule, bude

$$\sum_{i=1}^n p_i q_i = npq - n\sigma_p^2$$

a konečný výsledek tedy je

$$\mathfrak{E} \frac{1}{n} \sum_{i=1}^n (f_i - \bar{f})^2 = \frac{n-1}{n} \frac{pq}{r} + \frac{n(r-1) + 1}{nr} \sigma_p^2,$$

a při dostatečně velkém  $r$  se užívá obyčejně výrazu

$$\mathfrak{E}(\sigma_f^2) = \frac{n-1}{n} \frac{pq}{r} + \sigma_p^2. \quad (89)$$

Jsou-li všechny pravděpodobnosti v základním souboru sobě rovny  $p_1 = p_2 = \dots = p_n = p$ , pak je  $\sigma_p^2 = 0$  a dostáváme již známý výsledek

$$\mathbb{E}(\sigma_f^2) = \frac{n-1}{n} \frac{pq}{r}$$

Obě hodnoty směrodatných odchylek se srovnávají utvořením podflu. Označme  $\sigma_f$  směrodatnou odchylku řady relativních četností pozorovaných ve studovaném souboru. Za předpokladu konstantní pravděpodobnosti  $p$  je pro příslušné rozdělení Bernoulliovo t. zv. teoretická hodnota směrodatné odchylky  $\sigma_B = \left(\frac{pq}{r}\right)^{\frac{1}{2}}$ .

Podíl  $L = \frac{\sigma_f}{\sigma_B}$  se nazývá Lexisův poměr nebo koeficient.

Také se nazývá v teoreticko-statistické literatuře koeficient divergence (podle Dormoye). Místo směrodatných odchylek rozdělení relativních četností bychom mohli použítí směrodatných odchylek rozdělení absolutních četností, neboť

$$\sigma_x = r\sigma_f \text{ a také } S_B = r\sigma_B.$$

Lexisův poměr je tím větší, čím se více odchyluje (diverguje) statisticky zkoumaný jev od dění náhodného.

Říká se, že řada pozorovaných relativních četností má rozptyl normální, je-li  $L = 1$ , nadnormální (super-normální), je-li  $L > 1$  a podnormální (subnormální), je-li  $L < 1$ .

Vzhledem k (89) musíme tedy při zkoumání rozptylu srovnávati empirickou hodnotu  $\sigma_f^2$  s výrazem  $\frac{n-1}{n} \frac{pq}{r}$ , při čemž za  $\frac{pq}{r}$  musíme vzítí přibližnou hodnotu  $\frac{\bar{f}(1-\bar{f})}{r}$ .

Lexisův koeficient pak bude

$$Q^2 = 1 + \frac{n(r-1)+1}{nr} \sigma_p^2 : \frac{n-1}{n} \frac{pq}{r}, \quad (90)$$

nebo přibližně

$$Q^2 = 1 + \sigma_p^2 : \frac{n-1}{n} \frac{pq}{r}. \quad (91)$$

Máme tedy v koeficientu divergence důležitý prostředek k řešení nejvýznamnější úlohy statistiky, spočívající v zjištění, zda můžeme souditi na přítomnost změn v základních podmínkách výskytu znaku nebo na stále stejné, tedy konstantní, působení a složení základních podmínek.

Není sice absolutním kriteriem, ale dobrým vodítkem k posouzení kolísání výskytu znaku, jak u hromadných jevů fyzikálních, tak sociálních.

V praktické statistice se velmi často vyskytují řady, které daleko přesahují míru očekávaného rozptylu. Za příklad si zvolíme statistiku úmrtí s nadnormálním rozptylem, jejíž rozbor provedeme podle Misesa k objasnění uvedené teorie.

Ve státě se 45 miliony obyvatelů byla na př. cifra úmrtnosti obyvatelstva, t. j. počet úmrtí připadající na 1000 obyvatelů v desítiletém období, v němž stejnoměrnost životní úrovně nebyla rušena nějakými pozoruhodnými vnějšími jevy, zaznamenána v těchto promilech 28,0, 27,8, 27,2, 27,5, 26,9, 27,2, 27,3, 27,4, 27,2, 27,6.

Tyto relativní četnosti naplňují údivem svou stálostí toho, kdo na ně pohlíží bez znalostí matematické teorie statistiky. Skutečně dřívější statistikové byli v úžasu nad mimořádnou stabilitou lidských poměrů, jevíci se ve statistice. Dojdeme však ke zcela jinému závěru, vypočítáme-li skutečný rozptyl a srovnáme jej s očekávaným podle Lexisovy teorie.

Průměr uvedených deseti čísel je 27,41 promile, tudíž  $\bar{f} = 0,02741$ . Rozptyl pak dostaneme  $\sigma_L^2 = 0,000\ 000\ 0949$ . Očekávaná hodnota rozptylu řady Bernoulliho bude  $\sigma^2 =$

$$= \frac{\bar{f}(1 - \bar{f})}{r} \cdot \frac{n - 1}{n},$$
 kde koeficient  $\frac{n - 1}{n}$  vyplývá z teorie náhodného výběru podle (82), a hodnota  $p$  je nahrazena přibližnou hodnotou z pozorování, takže pro  $r = 45\,000\,000$  (průměrný počet obyvatelstva v uvažovaném desetiletí)  $\bar{f} = 0,02741$ ,  $n = 10$  dostaneme  $\sigma^2 = 0,000\,000\,000\,533$  a Lexisův poměr je  $L = 13,34$ . Přesahuje tedy skutečně pozorovaná směrodatná odchylka očekávanou teoretickou víc než 13krát.

Naznačme, jak možno provéstí rozbor tohoto výsledku. Lexisova teorie tu srovnává průběh roční úmrtnosti s deseti výběry, z nichž každý vznikl provedením 45 milionů tahů z osudí, v němž je stále mezi 100 000 kuličkami 2741 černých a 97 259 bílých. Kdyby na začátku každého z uvažovaných roků přišel každý obyvatel státu před toto osudí a vytáhl z něho svůj los života nebo smrti, museli bychom očekávat, že úmrtnost v tomto období vykáže rozptyl  $\sigma_B^2$ , který je 178krát menší než skutečně pozorovaný. Tento obraz nevystihuje hru o životě a smrti přiléhavě, neboť ze zkušenosti víme, že mnohé příčiny smrti působí současně na řadu lidí, jako na př. nepříznivý vývoj povětrnosti v nějakém zimním nebo letním měsíci, endemické onemocnění atd. Vzhledem k tomu bychom se přiblížili skutečnosti lépe, kdybychom předpokládali, že za celý soubor přijde k osudí menší část a každý se otáže po osudu celé skupiny, kterou zastupuje.

Je zřejmo, že podle vzorce  $\frac{pq}{r} \cdot \frac{n - 1}{n}$  bude tato očekávaná hodnota rozptylu tolikrát větší, kolikrát bude počet nezávislých jednotlivých případů  $r$  menší. Kdybychom tedy předpokládali v našem případě, že pro každých 178 obyvatelů bude tažen společný los, který rozhodne o životě nebo smrti celé jejich skupiny, dostali bychom úplný souhlas mezi pozorováním a očekáváním. Zda lze v konkrétním případě považovati vysvětlení silně nadnormálního rozptylu solidaritou jevů za případné, je třeba dále zkoumati. Bylo

by nutné, aby rozsah skupiny solidarity (178) zůstal zachován, když pozorujeme jiné analogické řady, na př. z jiných desetiletí. Kdyby to nebylo v dostačující míře splněno, bylo by třeba hledati jiné teoretické vysvětlení. V tomto případě lze je podati pomocí podstatné komponenty kolísání. Poněvadž se pravděpodobnost úmrtí rok od roku mění, jedná se o druhý typ Lexisovy řady, čili poměr černých a bílých kuliček je každý rok jiný. Potom, jak víme, je očekávaná hodnota rozptylu dána výrazem (89), čili k číslu nahore vypočítaného rozptylu přistupuje další složka, která nezávisí na  $r$ , nýbrž jen na kolísání pravděpodobnosti od jednoho roku ke druhému. V tom, že podstatná komponenta kolísání nezávisí na  $r$ , v našem případě na počtu obyvatelstva, máme kontrolu teorie, neboť při kolísání pravděpodobnosti úmrtí, následkem hospodářských nebo klimatických poměrů v celém státě, musí se tato komponenta vyskytovat v přibližně stejné výši v jednotlivých větších oblastech státu.

3. Třetím typem jsou řady Poissonovy. Jejich schema si představíme tak, že náhodný výběr rozsahu  $r$  se skládá z prvků, z nichž každý byl vzat z osudí jiného složení; pravděpodobnost výskytu pozorovaného znaku je tedy u každého prvku výběru jiná, takže schema můžeme napsati takto

$$\frac{O_1 \ O_2 \ \dots \ O_r}{p_1 \ p_2 \ \dots \ p_r}$$

$$p_1 \ p_2 \ \dots \ p_r$$

$$\dots \dots \dots$$

$$p_1 \ p_2 \ \dots \ p_r$$

kde  $p$  je pravděpodobnost vytažení černé kuličky z osudí  $O_k$ .

Označíme-li průměr těchto pravděpodobností  $p$ , píšeme  $p = \frac{p_1 + p_2 + \dots + p_r}{r}$ . Očekávaný průměr počtu černých

kuliček ve výběru rozsahu  $r$ , jehož každý prvek je z jiného osudí, je  $rp$  a rovná se očekávanému průměru počtu černých



kuliček, bereme-li náhodný výběr rozsahu  $r$  z jednoho osudí o konstantní pravděpodobnosti  $p$ .

Odvodíme nyní rozptyl počtu černých kuliček v řadě Poissonově. Rozptyl pro osudí  $O_k$  je dán výrazem  $S_k^2 = rp_kq_k$ , kdyby byl celý výběr z něho vzat. Vezmeme-li jen jeden prvek z něho, položíme  $r = 1$ . Jsou-li pravděpodobnosti  $p_1, p_2, \dots, p_r$  na sobě nezávislé, pak platí věta o sčítání rozptylů (67), takže celkový rozptyl

$$S_P^2 = S_1^2 + S_2^2 + \dots + S_r^2. \quad (92)$$

Dostáváme tudíž pro náš náhodný výběr

$$S_P^2 = p_1q_1 + p_2q_2 + \dots + p_rq_r = \sum_{k=1}^r p_kq_k$$

a tento výraz opět upravíme tak, že položíme

$$\begin{aligned} p_k &= p + (p_k - p) \\ q_k &= q - (p_k - p), \end{aligned}$$

takže

$$p_kq_k = pq - (p_k - p)(p - q) - (p_k - p)^2$$

a tudíž součet

$$\sum_{k=1}^r p_kq_k = rpq - \sum_{k=1}^r (p_k - p)^2,$$

neboť

$$\sum_{k=1}^r (p_k - p) = 0.$$

Pro rozptyl teoretického rozdělení počtu černých kuliček ve výběrech rozsahu  $r$  podle schematu Poissonova dostáváme tedy

$$S_P^2 = rpq - \sum_{k=1}^r (p_k - p)^2.$$

Je tedy rozptyl řady Poissonovy menší než rozptyl příslušné řady Bernoulliovy s konstantní pravděpodobností

rovnou průměru proměnné pravděpodobnosti:

$$S_P^2 = S_B^2 - \sum_{k=1}^r (p_k - p)^2. \quad (93)$$

Obdobnou rovnici pro rozdělení relativních četností lze snadno napsati jako v případě řad Lexisových

$$\sigma_P^2 = \sigma_B^2 - \frac{1}{r^2} \sum_{k=1}^r (p_k - p)^2. \quad (94)$$

**(9,2) Koefficient nestálosti.** Vedle Lexisova koeficientu zavedl Charlier koefficient nestálosti nebo disturbační, který rovněž měří vnější vlivy působící na změnu pravděpodobnosti v základním souboru. Definuje jej

$$e = \frac{\sqrt{\sigma_L^2 - \sigma_B^2}}{p}. \quad (95)$$

Jeho přibližnou hodnotu dostáváme, klademe-li

místo  $\sigma_L^2$  přibližnou hodnotu  $\frac{\sum (f_i - \bar{f})^2}{n - 1}$ ,

místo  $\sigma_B^2$  „ „  $\frac{\bar{f}(1 - \bar{f})}{r}$

a místo  $p$  „ „  $\bar{f}$ .

Jako příklad si zvolíme poměr pohlaví živě-narozených dětí, což je velmi probádaným předmětem statistického šetření. Bylo mínění, že řady těchto čísel odpovídají poměrům náhodné hry o konstantní pravděpodobnosti, tedy řadě Bernoulliově. Přesto máme výsledky konkrétních šetření, kde se objevuje rozptyl podnormální, což je v praxi statistické řídkým případem. Tak byly na př. ve Vídni (Wien) pozorovány ve 24 měsících let 1908 a 1909 tyto relativní četnosti chlapců mezi celkovým počtem živě narozených:

0,5223	0,5125	0,5141	0,5246	0,5126	0,5136
0,5187	0,5213	0,5105	0,5203	0,5124	0,5141
0,5143	0,5093	0,4904	0,5097	0,5140	0,5089
0,5129	0,5275	0,5178	0,5130	0,5177	0,5027

Jejich průměr je  $\bar{f} = 0,514$  a rozptyl  $\sigma_f^2 = 0,0000533$ . Celkem se tam narodilo v té době 93 661 dětí, takže průměrně připadá na jeden měsíc  $r = 3903$  dětí. Ve smyslu Lexisovy teorie se nyní ptáme, jaký je očekávaný rozptyl při konstantní pravděpodobnosti  $p$ , který by odpovídal 24 výběrům rozsahu 3903 z osudí stálého složení, kde mezi každým tisícem losů je 514 označeno znakem  $c$  (chlapec).

Vypočítáme tedy tento rozptyl podle formule  $\frac{pq}{r} \cdot \frac{n-1}{n}$ ,

kde  $n = 24$ ,  $r = 3903$ ,  $p = \bar{f} = 0,514$  a dostaneme  $\sigma^2 = 0,0000613$ , takže pro Lexisův poměr dostáváme  $L = 0,93$ , tedy rozptyl je podnormální. To nasvědčuje tomu, že případ neodpovídá osudí téhož složení, nýbrž jsou tu části obyvatelstva, jimž přísluší rozmanité pravděpodobnosti porodu se znakem  $c$ .

V podobném statistickém šetření, provedeném na př. ve Švédsku, byl zjištěn rozptyl poněkud nadnormální, kdežto na př. pro počet dvojčat v poměru k jednotlivým porodům se projevil rozptyl silně podnormální. Když byla pomocí čísla  $L$  konstatována existence rušivých vlivů na statistické řady, je pak úkolem statistickým, pátrati po příčině poruch. Obecnou metodu k tomu dává teorie korelace. Na základě uvažování předložené řady lze dospěti jen k určitým závěrům o povaze rušivých vlivů, jimž je vystaven statisticky pozorovaný jev. Podle teorie Lexisovy vznikají tyto poruchy tím, že pravděpodobnost pro výskyt znaku se mění.

Úloha. Máme schema deseti osudí, z nichž každé obsahuje 15 kuliček, ale má postupně mezi nimi 3, 4, 5, 6, 7, 8, 9, 10, 11, 12 bílých. Průměr pravděpodobností táhnouti bílou kuličku

je 0,5. Tvořme pokusem náhodné výběry po 10 prvcích tak, že z každého osudí vytáhneme jednu kuličku, pak je zase vrátíme do osudí a vytáhneme nových deset kuliček. Porovnejme potom čísla  $L$  pro 200, 500 po případě 1000 výběrů, abychom si pomocí tohoto kriteria zjistili stupeň shody, po případě neshody pozorování s tím, co očekáváme podle teorie.

---