

9 NÁHODNÉ VÝBĚRY A JEJICH ZPRACOVÁNÍ



Čas ke studiu kapitoly: 30 minut



Cíl: Po prostudování tohoto odstavce budete rozumět pojům

- Základní soubor, populace, výběr, výběrové šetření, výběrová statistika a budete znát základní výběrové statistiky pro výběry z normálního rozdělení



Výklad:

Motto:

Chceme-li vědět, jak chutná víno v sudu, nemusíme vypít celý sud. Stačí jenom malý doušek a víme na čem jsme.

Statistika – to je sběr a zpracování dat. V mnoha oborech se setkáme s průzkumy opírajícími se o relativně malý počet zkoumaných jednotek (**výběr**). Statistika pak používá postupy pomocí nichž můžeme, sice s určitým rizikem (předem stanoveným), na základě toho mála usuzovat na chování celku (**populace**). Tomuto zobecňování říkáme **statistická indukce**.

9.1 Statistické zjišťování

Pro většinu statistických souborů, s nimiž se v praxi setkáváme, je typický vysoký **rozsah** (počet zkoumaných jednotek). Jakmile jsme tedy postaveni před úkol provést určité šetření a analyzovat údaje z něj zjištěné, musíme nejprve rozhodnout, zda budeme toto šetření realizovat jako vyčerpávající nebo výběrové.

Vyčerpávající šetření – to je prošetření všech jednotek statistického souboru (populace). Zpravidla se jedná o záležitost velmi nákladnou (personálně, finančně, časově), mnohdy dokonce prakticky nerealizovatelnou (destrukční zkoušky). Pokud však toto šetření proběhne, mezi jeho nesporné výhody patří přesnost zjištěných charakteristik a detailnost formací o každé zkoumané jednotce. Příkladem vyčerpávajícího šetření je například sčítání lidu.

Výběrové šetření – jde o prošetření vybraných jednotek statistického souboru (populace). Z takto pořízených charakteristik pak více či méně usuzujeme na vlastnosti celé populace. Výběrová šetření se používají například při zjišťování jaká je podpora politických stran, při ověřování pevnosti trubek vyráběných určitým podnikem, apod. Mírou objektivnosti informací, které z něho získáme, je kvalita provedení výběrového šetření.

9.2 Typy výběrových šetření

Základní soubory, z nichž vybíráme mohou být buď konečně nebo nekonečně velké. Příkladem konečně velkého základního souboru je dodávka výrobku (např. praček), příkladem nekonečně velkého základního souboru je nepřetržitá pásová výroba (např. praček). Při konstrukci výběrového souboru se snažíme o to, aby výběrový soubor měl stejné vlastnosti jako základní soubor, z něhož výběr pochází.

Mezi druhy výběrových šetření řadíme anketu, metodu základního masivu, záměrný výběr a náhodný výběr.

Anketa oslovuje pouze nesystematicky vybranou část populace (osob, podniků, institucí). Dotazník se k respondentům (dotazovaným) dostává prostřednictvím sdělovacích prostředků (anketa televizních diváků, anketa časopisu Mladí, ...) nebo je zaslán adresně. Návratnost dotazníku je však malá (odhaduje se že 30%). Informace získané anketním šetřením nelze zobecňovat.

Metoda základního masivu se používá v případech, kdy se základní soubor skládá z několika velkých jednotek a z většího počtu jednotek malých. (např. při šetření v oblasti hutnictví se můžeme podle této metody zaměřit na několik „obřích“ společností, tam provést šetření a „malé“ podniky vynechat. Výhody: menší pracnost a menší časová náročnost šetření. Nevýhody: zobecnění poznatků má menší platnost (nevystihuje specifika menších jednotek).

Záměrný výběr spočívá v tom, že skupina odborníků na danou problematiku vybere podle svého nejlepšího uvážení ty jednotky, o nichž se lze domnívat, že ve svém souhrnu nejlépe umožní provést šetření. S tímto typem šetření se setkáme například při průzkumech trhu a při průzkumech veřejného mínění. Nevýhoda: subjektivní přístup k výběru zpochybňuje možnost zobecnění.

Prostý náhodný výběr je základním a v praxi nejpoužívanějším typem výběru. Jde o výběr, při němž mají všechny jednotky základního souboru stejnou pravděpodobnost, že do výběru budou zařazeny.

Nyní se náhodným výběrem budeme zabývat podrobněji (a formálněji).

9.3 Náhodný výběr

Náhodný výběr (\underline{X}) je speciální náhodný vektor, jehož složky jsou nezávislé náhodné veličiny se stejným rozdělením pravděpodobnosti.

Pokusíme se uvedenou definici vysvětlit. Opakujeme-li n -krát nezávisle pokus (pozorování, měření), jehož výsledek je náhodná veličina X s distribuční funkcí $F(x)$, pozorujeme vlastně náhodný vektor $\underline{X} = (X_1, \dots, X_n)^T$, $X_i \approx F(x)$, jehož složky jsou vzájemně nezávislé náhodné veličiny s touž distribuční funkcí $F(x)$. Tento vektor nazýváme **náhodný výběr** z rozdělení $F(x)$ nebo **náhodný výběr ze základního souboru** (nebo **populace**) s rozdělením $F(x)$. Číslo n se nazývá **rozsah** náhodného výběru.

Podle rozsahu obvykle rozdělujeme náhodné výběry na **malé** ($n \leq 30$) a **velké** ($n > 30$).

Náhodný výběr má zřejmě **sruženou distribuční funkci** $F(\underline{x})$:

$$\begin{aligned} F(\underline{x}) &= F(x_1, \dots, x_n) = P(X_1 < x_1; \dots; X_n < x_n) = P(X_1 < x) \cdot \dots \cdot P(X_n < x_n) = \\ &= F(x_1) \cdot \dots \cdot F(x_n) = \prod_{i=1}^n F(x_i) \end{aligned}$$

a podobně i **sruženou hustotu pravděpodobnosti**:

$$f(\underline{x}) = \prod_{i=1}^n f(x_i)$$

Číselný vektor (x_1, \dots, x_n) , který získáme při realizaci náhodného výběru (X_1, \dots, X_n) , nazýváme **statistický soubor** nebo vzorek o rozsahu n . Množina všech těchto vektorů se nazývá **výběrový prostor**. Je to zřejmě podmnožina množiny R^n .

Řadu informací o posuzované náhodné veličině X poskytují její číselné charakteristiky, např. EX , DX , σ_X atd. Při statistické indukci jsme při určování jejich hodnot odkázáni na realizace náhodných výběrů, tedy na statistické soubory. Užíváme přitom následující pojmy:

Funkci náhodného výběru $\underline{X} = (X_1, \dots, X_n)^T$, k jejímuž určení není třeba znát konkrétní hodnoty parametrů příslušného rozdělení, nazýváme **statistika** nebo **výběrová charakteristika** a značíme ji $T(\underline{X}) = T(X_1, \dots, X_n)$. Je to obecně náhodná veličina.

Její hodnotu $t = T(x_1, \dots, x_n)$, kterou nabývá na statistickém souboru $(x_1, \dots, x_n)^T$, nazýváme **pozorovaná hodnota statistiky T** nebo empirická charakteristika.

Používáme zejména následující statistiky:

$$1. \quad T_1(\underline{X}) = \bar{X} = \frac{\sum_{i=1}^n X_i}{n} \quad \dots \quad \text{výběrový průměr}$$

$$E(T_1(\underline{X})) = E\bar{X} = \frac{1}{n} \cdot \sum_{i=1}^n EX_i = \frac{EX_i}{n} \cdot \sum_{i=1}^n 1 = EX_i$$

$$2. \quad T_2(\underline{X}) = S^2 = \frac{1}{n-1} \cdot \sum_{i=1}^n (X_i - \bar{X})^2 \quad \dots \quad \text{výběrový rozptyl}$$

Není těžké ukázat, že $ES^2 = DX_i$

$$3. \quad T_3(\underline{X}) = \sqrt{S^2} = S \quad \dots \quad \text{výběrová směrodatná odchylka}$$

Nechť v daném výběru je počet prvků se sledovanou vlastností x_v , pak:

$$4. \quad T_4(\underline{X}) = p = \frac{x_v}{n} \quad \dots \quad \text{výběrová relativní četnost (výběrový podíl)}$$

9.4 Výběrová rozdělení – rozdělení statistik či výběrových charakteristik

Předpokládejme, že daný náhodný výběr pochází z normálního rozdělení:

$$\underline{X} = (X_1, \dots, X_n)^T, \quad X_i \rightarrow N(\mu, \sigma^2)$$

$$1. \quad \bar{X} = \frac{\sum_{i=1}^n X_i}{n} \rightarrow N\left(\mu, \frac{\sigma^2}{n}\right)$$

(plyne jednak z centrální limitní věty pro velká n , ale dá se také ukázat na základě odvození rozdělení součtu náhodných veličin)

$$2. \quad Z = \frac{\bar{X} - \mu}{\sigma} \cdot \sqrt{n} \rightarrow N(0,1)$$

(plyne z transformace předešlého rozdělení)

$$3. \chi = \frac{S^2}{\sigma^2} \cdot (n-1) \rightarrow \chi^2(n-1)$$

(bylo vysvětleno při diskusi rozdělení χ^2)

$$4. T = \frac{\bar{X} - \mu}{S} \cdot \sqrt{n} \rightarrow t_{n-1}$$

(odvozeno při diskusi o použití Studentova rozdělení)

$$5. P_1 = \frac{p - \pi}{\sqrt{\pi(1-\pi)}} \cdot \sqrt{n} \rightarrow N(0;1)$$

(odvozeno při diskusi aplikaci centrální limitní věty – kap. 7.5.1)

Nyní předpokládáme dva výběry z normálních rozdělení:

$$\underline{X} = (X_1, \dots, X_{n_x})^T, X_i \rightarrow N(\mu_x, \sigma_x^2), \underline{Y} = (Y_1, \dots, Y_{n_y})^T, Y_j \rightarrow N(\mu_y, \sigma_y^2).$$

Potom platí:

$$6. Z_2 = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}} \rightarrow N(0,1)$$

$$\bar{X} \rightarrow N\left(\mu_x; \frac{\sigma_x^2}{n_x}\right); \quad \bar{Y} \rightarrow N\left(\mu_y; \frac{\sigma_y^2}{n_y}\right); \quad (\bar{X} - \bar{Y}) \rightarrow N\left(\mu_x - \mu_y; \frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}\right)$$

$$7. F = \frac{\frac{S_x^2}{\sigma_x^2} \cdot (n_x - 1)}{\frac{S_y^2}{\sigma_y^2} \cdot (n_y - 1)} = \frac{\frac{S_x^2}{n_x - 1}}{\frac{S_y^2}{n_y - 1}} \rightarrow F_{n_x-1, n_y-1}$$

(zdůvodněno v souvislosti s F -rozdělením)

Předpokládejme speciální případ, že **rozptily jsou neznámé avšak stejné**: $\sigma_x^2 = \sigma_y^2$. Potom se dá ukázat, že platí:

$$8. T_2 = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sqrt{S_x^2(n_x - 1) + S_y^2(n_y - 1)}} \cdot \sqrt{\frac{n_x \cdot n_y \cdot (n_x + n_y - 2)}{n_x + n_y}} \rightarrow t_{n_x + n_y - 2}$$

Nechť mají dané výběrové soubory počty prvků se sledovanou vlastností x_V a y_V . Pak výběrové relativní četnosti určíme jako:

$$p_X = \frac{x_V}{n_x}; \quad p_Y = \frac{y_V}{n_y}$$

Pak platí:

$$9. \quad P_2 = \frac{(p_X - p_Y) - (\pi_X - \pi_Y)}{\sqrt{p(1-p)\left(\frac{1}{n_X} + \frac{1}{n_Y}\right)}} \rightarrow N(0;1), \text{ kde } p = \frac{x_V + y_V}{n_X + n_Y}$$



Shrnutí:

Statistika pak používá postupy pomocí nichž můžeme, sice s určitým rizikem (předem stanoveným), na základě toho mála usuzovat na chování celku (**populace**). Tomuto zobecnování říkáme **statistická indukce**.

Jakmile jsme postavení před úkol provést určité šetření a analyzovat údaje z něj zjištěné, musíme nejprve rozhodnout, zda budeme toto **šetření** realizovat jako **vyčerpávající nebo výběrové**.

Vyčerpávající šetření – to je prošetření všech jednotek statistického souboru (populace).

Výběrové šetření – jde o prošetření vybraných jednotek statistického souboru (populace).

Mezi druhy výběrových šetření řadíme **anketu, metodu základního masivu, záměrný výběr a náhodný výběr**.

Náhodný výběr (\underline{X}) je speciální náhodný vektor, jehož složky jsou nezávislé náhodné veličiny se stejným rozdělením pravděpodobnosti.

Číselný vektor $(x_1, \dots, x_n)'$, který získáme při realizaci náhodného výběru $(X_1, \dots, X_n)'$, nazýváme **statistický soubor** nebo vzorek o rozsahu n . Množina všech těchto vektorů se nazývá **výběrový prostor**.

Funkci náhodného výběru $\underline{X}=(X_1, \dots, X_n)'$, k jejímuž určení není třeba znát konkrétní hodnoty parametrů příslušného rozdělení, nazýváme **statistika** nebo **výběrová charakteristika** a značíme ji $T(\underline{X})$. Její hodnotu $t=T(x_1, \dots, x_n)$, kterou nabývá na statistickém souboru $(x_1, \dots, x_n)'$, nazýváme **pozorovaná hodnota statistiky T** nebo empirická charakteristika.

Používáme zejména následující statistiky: **výběrový průměr \bar{x} , výběrový rozptyl s^2 a výběrovou směrodatnou odchylku s a výběrový podíl p** .

Za předpokladu, že náhodný výběr pochází z normálního rozdělení pravděpodobnosti, se dají z daného náhodného výběru odvodit další významné statistiky se známým rozdělením:

Výběrová charakteristika	Rozdělení výběrové charakteristiky
$Z = \frac{\bar{X} - \mu}{\sigma} \cdot \sqrt{n}$	$N(0;1)$
$T_{n-1} = \frac{\bar{X} - \mu}{S} \cdot \sqrt{n}$	t_{n-1}
$\chi = \frac{(n-1)S^2}{\sigma^2}$	χ_{n-1}^2
$P_1 = \frac{p - \pi}{\sqrt{\pi(1-\pi)}} \cdot \sqrt{n}$	$N(0;1)$

Máme-li k dispozici dva výběry z normálního rozdělení, setkáváme s následujícími výběrovými statistikami:

Výběrová charakteristika	Rozdělení výběrové charakteristiky
$Z_2 = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\frac{\sigma_X^2}{n_X} + \frac{\sigma_Y^2}{n_Y}}}$	N(0;1)
$T_2 = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{S_p \cdot \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}}}$	$t_{n_X + n_Y - 2}$
$F = \frac{\frac{S_X^2}{\sigma_X^2}}{\frac{S_Y^2}{\sigma_Y^2}}$	$F_{n_X - 1, n_Y - 1}$
$P_2 = \frac{(p_X - p_Y) - (\pi_X - \pi_Y)}{\sqrt{p(1-p)\left(\frac{1}{n_X} + \frac{1}{n_Y}\right)}}$	N(0;1)



Otázky

1. Co je statistická indukce ?
2. Charakterizujte pojmy náhodný výběr a statistický soubor.
3. Co jsou výběrové charakteristiky a které z nich se nejčastěji používají ?
4. Vyjmenujte některé z dalších výběrových statistik, tj. statistik odvozených z náhodného výběru z normálního rozdělení a u některých z takto vyjmenovaných statistik se pokuste zdůvodnit rozdělení pravděpodobnosti.