

# Statistika

Zpracování informací ze statistického šetření  
– Třídění statistického souboru

Roman Biskup

(zapálený) statistik ve výslužbě, aktuálně analytik v praxi ;-)  
roman.biskup(at)email.cz

20. února 2012



# Obsah

Třídění dle statistického znaku  
Prosté a intervalové třídění  
Četnosti statistického znaku

Tabulky četností  
Prosté třídění  
Intervalové třídění

Grafická vizualizace rozložení četností  
Polygon četností  
Histogram četnosti  
Výsečový graf



## Třídění dle statistického znaku

- ▶ **Důvody třídění:**
  - ▶ zpřehlednění souboru,
  - ▶ zjištění empirického rozdělení statistického souboru,
  - ▶ **snížení numerické náročnosti výpočtu statistických charakteristik.**
- ▶ **Dle počtu třídících znaků:**
  - ▶ jednostupňové,
  - ▶ dvoustupňové (kontingenční tabulky),
  - ▶ vícestupňové.
- ▶ **Dle typu třídění:**
  - ▶ třídění prosté (malý počet různých hodnot znaku),
  - ▶ třídění intervalové (velký počet různých hodnot znaku, spojitý numerický znak).
- ▶ **Základní zásady při třídění:**
  - ▶ zásada úplnosti (každá jednotka musí někam patřit),
  - ▶ zásada jednoznačnosti (každá jednotka musí mít právě jedno místo při třídění).



## Postup třídění I

- ▶ **Prosté třídění – libovolný statistický znak**
  1. stanovení počtu pozorování různých hodnot znaku (předpokládáme  $k$  různých hodnot)
- ▶ **Intervalové třídění – numerický statistický znak**
  1. stanovení počtu intervalů  $k$ , optimálně  $8 \leq k \leq 20$ 
    - ▶  $k \approx 1 + 3,3 \cdot \log n$  (Sturgesovo pravidlo)
    - ▶  $k \approx \frac{8}{100}(\max x - \min x)$
    - ▶  $k \approx \sqrt{n}$
  2. stanovení délky intervalu  $h$

$$h = \frac{\max x - \min x}{k}$$

3. rozdělení na intervaly  $\langle \bullet; \bullet \rangle, \langle \bullet; \bullet \rangle, \dots, \langle \bullet; \bullet \rangle$

$$\langle \min x + i \cdot h; \min x + (i + 1) \cdot h \rangle, \text{ pro } i = 0, \dots, k - 2$$

$$\langle \min x + (k - 1) \cdot h; \max x \rangle$$

- ▶ Pro popis statistického znaku je vhodné jak délku intervalů, tak hranice intervalů „učesat“, tj. vhodně zaokrouhlit; je však třeba zajistit, aby takto upravené intervaly pokryly všechny hodnoty statistického znaku.



## Postup třídění II

- ▶ Meze jednotlivých intervalů je třeba volit tak, aby nedocházelo k nejasnostem, tj. aby se každé pozorování jednoznačně „spadalo“ do určitého intervalu.
4. stanovení počtu pozorování s hodnotou znaku spadajícího do příslušného intervalu

## Absolutní a relativní četnost II

- ▶ i zde zřejmě:

$$\sum_{i=1}^k p_i = 1 \quad (100\%).$$

## Absolutní a relativní četnost I

Označme sledovaný statistický znak  $x$ , necht' má  $N$  pozorování, pak pro  $i = 1, \dots, k$ :

$n_i$  absolutní četnost

↪ počet pozorování s hodnotou znaku rovnou  $x_i$ , respektive počet pozorování s hodnotou znaku spadající do  $i$ -tého intervalu,

- ▶ zřejmě platí:

$$\sum_{i=1}^k n_i = N.$$

$p_i$  relativní četnost

↪ poměr počtu pozorování s hodnotou znaku rovnou  $x_i$  vzhledem celkovému počtu pozorování, respektive poměr počtu pozorování s hodnotou znaku spadající do  $i$ -tého intervalu vzhledem celkovému počtu pozorování,

$$p_i = \frac{n_i}{N} \quad (p_i \cdot 100\%), \quad i = 1, \dots, k;$$

## Kumulativní četnosti I

$k_{n_i}$  kumulativní (absolutní) četnost

↪ počet pozorování, u nichž je hodnota statistického znaku  $x \leq x_i$ , respektive počet pozorování zařazených díky hodnotě statistického znaku od prvního až do  $i$ -tého intervalu včetně, tj.

$$k_{n_i} = n_1 + n_2 + \dots + n_i = \sum_{j=1}^i n_j.$$

$k_{p_i}$  kumulativní relativní četnost

↪ udává poměr počtu pozorování, u nichž je hodnota statistického znaku  $x \leq x_i$ , vzhledem k celkovému počtu pozorování respektive poměr počtu pozorování zařazených díky hodnotě statistického znaku od prvního až do  $i$ -tého intervalu včetně vzhledem k celkovému počtu pozorování, tj.

$$k_{p_i} = p_1 + p_2 + \dots + p_i = \sum_{j=1}^i p_j.$$

- ▶ Je nutné uspořádání znaku  $x$ , tj. má smysl dělat minimálně pro ordinální znak.
- ▶ Nebo ne? Jakou by pak měla kumulativní četnost interpretaci?
- ▶  $k_{n_k} = N$ ,  $k_{p_k} = 1$  (100%)

## Datový soubor – Evidence studijních výsledků LS 2005

Obor	Počet *	Zameškáno	Zápočet	Body	Hodnocení
PUPN	4	0	Ano	4	1
VZ	0	3	Ano	1,5	4
OP	0	2	Rost	4	2
PP	0	0	Biskup	2	4
VZ	0	3	Ano	1	4
OP	0	1	Rost	2	4
ZOO	1	0	Ano	4	2
BT	13	1	Ano	4	2
OP	0	0	Rost	0,5	4
VZ	1	2	Ano	4	2
VZ	0	3	Ne	0	4
VZ	0	2	Ano	1,5	4
ZOO	2	1	Ano	1,5	4
:					
:					

## Tabulka četností – Evidence studijních výsledků LS 2005

Počet bodů získaných z písemné části zkoušky ze statistiky LS 2005 (řádný termín)

$x_i$	$n_i$	$p_i$ (%)	$k_{n_i}$	$k_{p_i}$ (%)
0,0	27	19,42	27	19,42
0,5	11	7,91	38	27,34
1,0	20	14,39	58	41,73
1,5	15	10,79	73	52,52
2,0	14	10,07	87	62,59
2,5	11	7,91	98	70,50
3,0	22	15,83	120	86,33
3,5	8	5,76	128	92,09
4,0	7	5,04	135	97,12
4,5	3	2,16	138	99,28
5,0	1	0,72	139	100,00
5,5	0	0,00	139	100,00
6,0	0	0,00	139	100,00
<b>Σ</b>	<b>139</b>	<b>100,00</b>		

## Přípravné práce – Evidence studijních výsledků LS 2005

Body – počet bodů získaných z písemné části zkoušky ze statistiky LS 2005 (řádný termín)

- $N = 139$ ;  $k = 13$  (0; 0,5; ... ; 6 bodů) stanovení počtu pozorování jednotlivých hodnot znaku ...

## Datový soubor – Splátkový prodej (2004)

Věk	Pohlaví	Stav	Vzdělání	Zaměstnání	Příjem (Kč)	Úvěr (Kč)	Splátek
59	žena	ženatý	základní	důchodce	7 200	5 390	20
27	žena	ženatý	střední	dělník	7 000	7 542	20
50	muž	rozvedený	střední	kuchař	61 000	6 216	10
29	muž	svobodný	vyučený	dělník	10 000	7 002	20
31	muž	ženatý	vyučený	řidič	15 000	8 982	10
19	žena	druhý	základní	mateř dovolená	5 500	6 696	10
22	muž	svobodný	vyučený	malíř, natěrač	10 000	4 621	20
34	muž	ženatý	střední	stát. zam.	15 159	7 624	30
45	žena	ženatý	vyučený	podnikatel	10 000	7 515	20
24	muž	rozvedený	vyučený	technik	12 000	6 680	20
30	muž	rozvedený	vyučený	pekař	12 500	3 228	20
25	muž	svobodný	střední	pol. inspektor	14 000	14 229	30
:							
:							

## Příprava intervalů – Splátkový prodej (2004)

Úvěr – cena zaplacená za celkový spotřebitelský úvěr;

- $N = 737$ ;  $k \approx 1 + 3,3 \cdot \log 737 = 10,463$ , zvolme  $k = 11$ ;
- $\min x = 1584$  a  $\max x = 25164$ ;
- $h = \frac{25164 - 1584}{11} = 2151,273$ , položíme  $h = 2200$  a dolní mez prvního intervalu rovnu 1500 pak:
  - $\langle 1500 ; 3700 \rangle$
  - $\langle 3700 ; 5900 \rangle$
  - $\langle 5900 ; 8100 \rangle$
  - $\vdots$
  - $\langle 21300 ; 23500 \rangle$
  - $\langle 23500 ; 25700 \rangle$
- stanovení počtu pozorování v jednotlivých intervalech ...

## Polygon četností

- vizualizace absolutních četností – prosté třídění
  - na vodorovnou osu se vynášejí hodnoty sledovaného znaku
  - na svislou osu se pak vynášejí absolutní četnosti
  - nad jednotlivými hodnotami znaku jsou vynášeny hodnoty odpovídající příslušným absolutním četnostem
  - jednotlivé hodnoty jsou navíc spojeny lomenou čarou

## Tabulka četností – Splátkový prodej (2004)

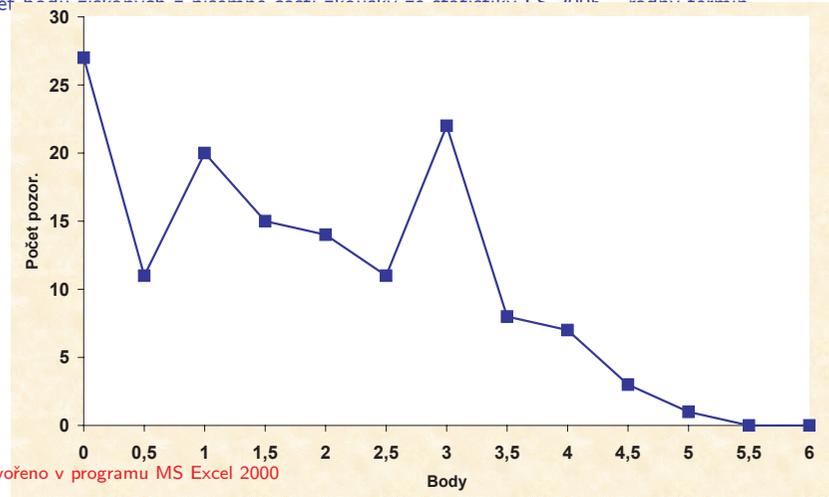
Cena zaplacená za celkový spotřebitelský úvěr

		Tabulka četností: Celková výše úvěru			
OD	DO	Četnost	Kumulativní četnost	Rel. četnost	Kumulativní rel. četnost
1 500 Kč	$\leq x < 3 700$ Kč	69	69	9,36228	9,3623
3 700 Kč	$\leq x < 5 900$ Kč	217	286	29,44369	38,8060
5 900 Kč	$\leq x < 8 100$ Kč	218	504	29,57938	68,3853
8 100 Kč	$\leq x < 10 300$ Kč	104	608	14,11126	82,4966
10 300 Kč	$\leq x < 12 500$ Kč	55	663	7,46269	89,9593
12 500 Kč	$\leq x < 14 700$ Kč	54	717	7,32700	97,2863
14 700 Kč	$\leq x < 16 900$ Kč	15	732	2,03528	99,3216
16 900 Kč	$\leq x < 19 100$ Kč	3	735	0,40706	99,7286
19 100 Kč	$\leq x < 21 300$ Kč	0	735	0,00000	99,7286
21 300 Kč	$\leq x < 23 500$ Kč	1	736	0,13569	99,8643
23 500 Kč	$\leq x < 25 700$ Kč	1	737	0,13569	100,0000
<b>ČSD</b>		<b>737</b>	<b>737</b>	<b>0,00000</b>	<b>100,0000</b>

Vytvořeno v programu STATISTICA komplet 6.1.02

## Polygon četností

Počet bodů získaných na ústřední části zkoušky ze statistiky LS 2005 – každý bod má



Vytvořeno v programu MS Excel 2000

## Histogram četnosti

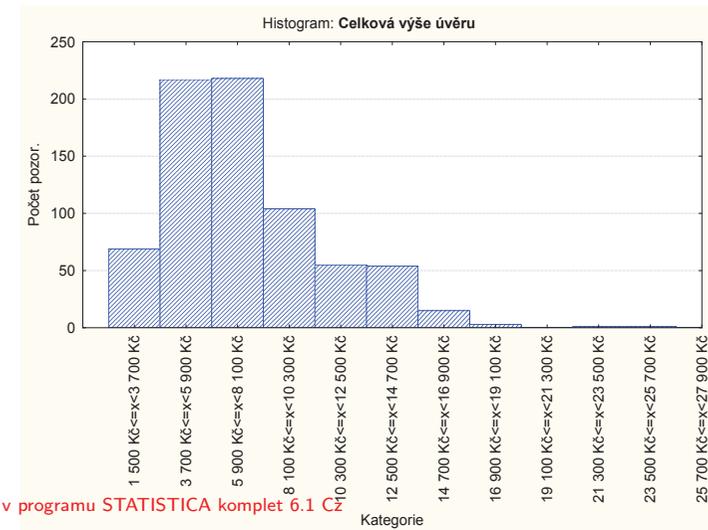
- ▶ vizualizace absolutních četností – intervalového třídění
  - ▶ na vodorovnou osu se vynášejí meze intervalů
  - ▶ na svislou osu pak absolutní četnosti
  - ▶ nad jednotlivými intervaly jsou vykresleny sloupce s podstavou šířky intervalu a výškou absolutní četnosti
  - ▶ někdy jsou hodnoty vynášené na svislou osu modifikovány tak, aby celková plocha sloupců byla rovná jedné
- ▶ vše pochopitelně v měřítku ;-)

## Výšečový (koláčový) graf

- ▶ vizualizace relativních četností
  - ▶ plocha grafu je dělena na kruhové výšeče v poměru, který je dán relativní četností, tj.
    - ▶  $|\angle_i| = 360^\circ \cdot p_i$ ,
    - ▶ zřejmě platí:  $\sum_{i=1}^k \angle_i = 360^\circ$ .
- ▶ Graf je obvykle doplněn o legendu a relativní četnosti v procentech

## Histogram četnosti

Cena zaplacená za celkový spotřebitelský úvěr



## Výšečový graf

Výsledné známky ze Statistiky 2004/05 – LS

