

PSY117 2019

Statistická analýza dat v psychologii

Přednáška 2

MÍRY CENTRÁLNÍ TENDENCE A VARIABILITY

He uses statistics as a drunken man uses lampposts – for support rather than illumination.



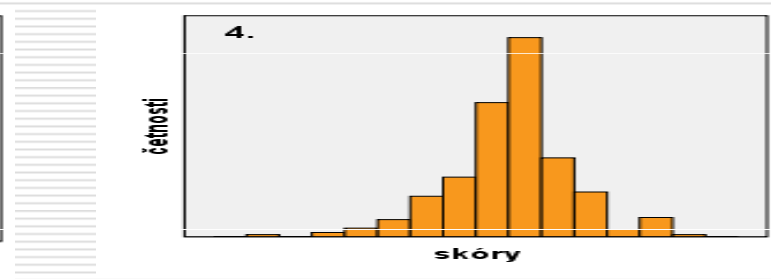
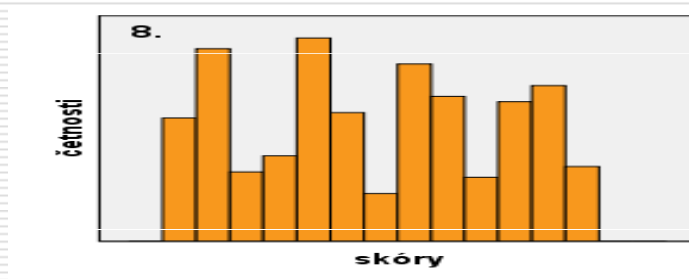
Andrew Lang

Z minula

- První informací (*statistikou*), která nás zajímá je **četnost** výskytu jednotlivých hodnot (resp. hodnot uvnitř jednotlivých intervalů)
- Četnosti – **absolutní, relativní, kumulativní** komunikujeme
 - Tabulkou četností
 - Sloupcovým diagramem, histogramem

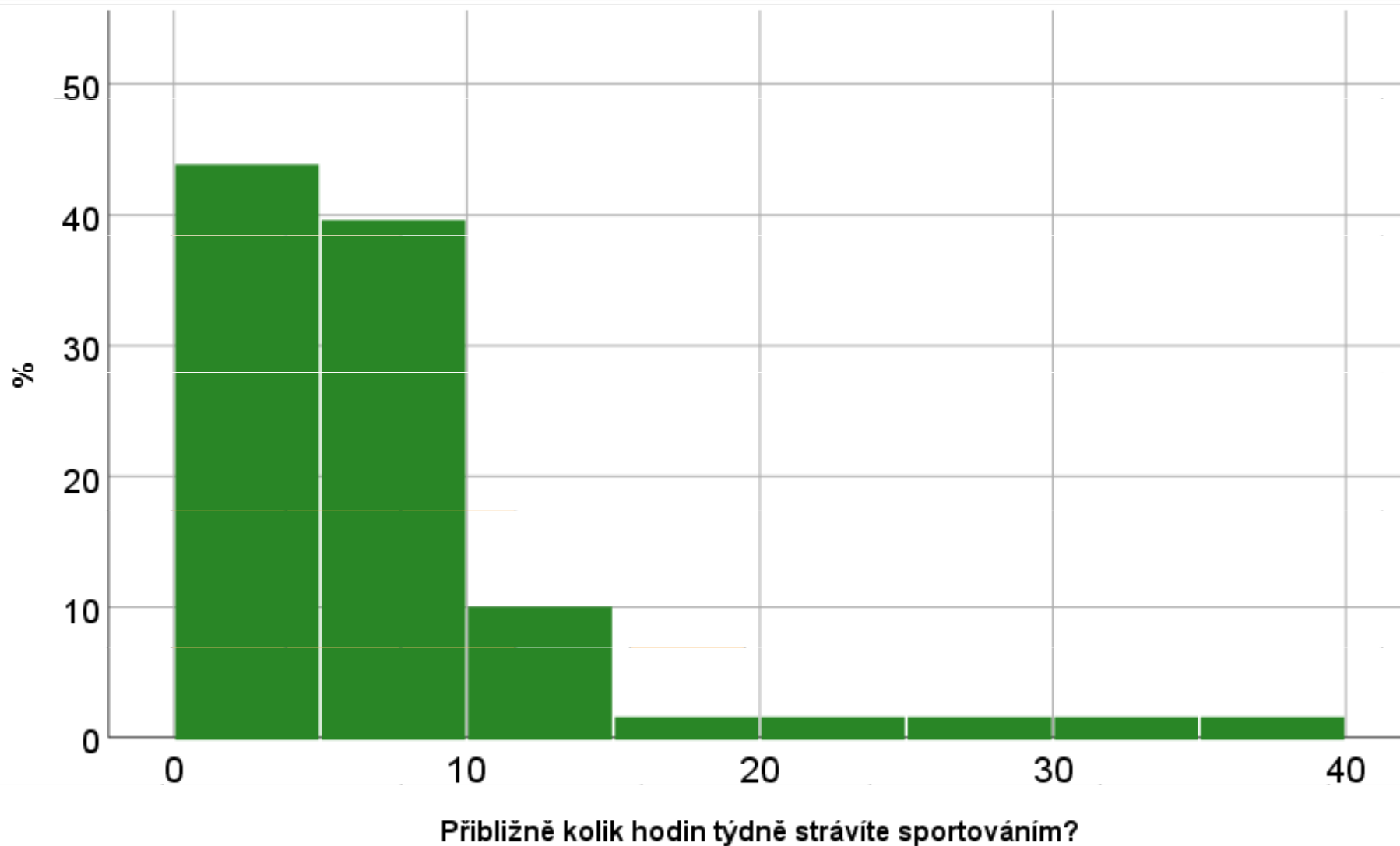
Rozložení *rozdělení, distribuce* četností

- Četnosti hodnot ordinálních a vyšších proměnných obvykle nebývají **distribuvány** nahodile – jejich rozložení zobrazené histogramem má popsateľný tvar.

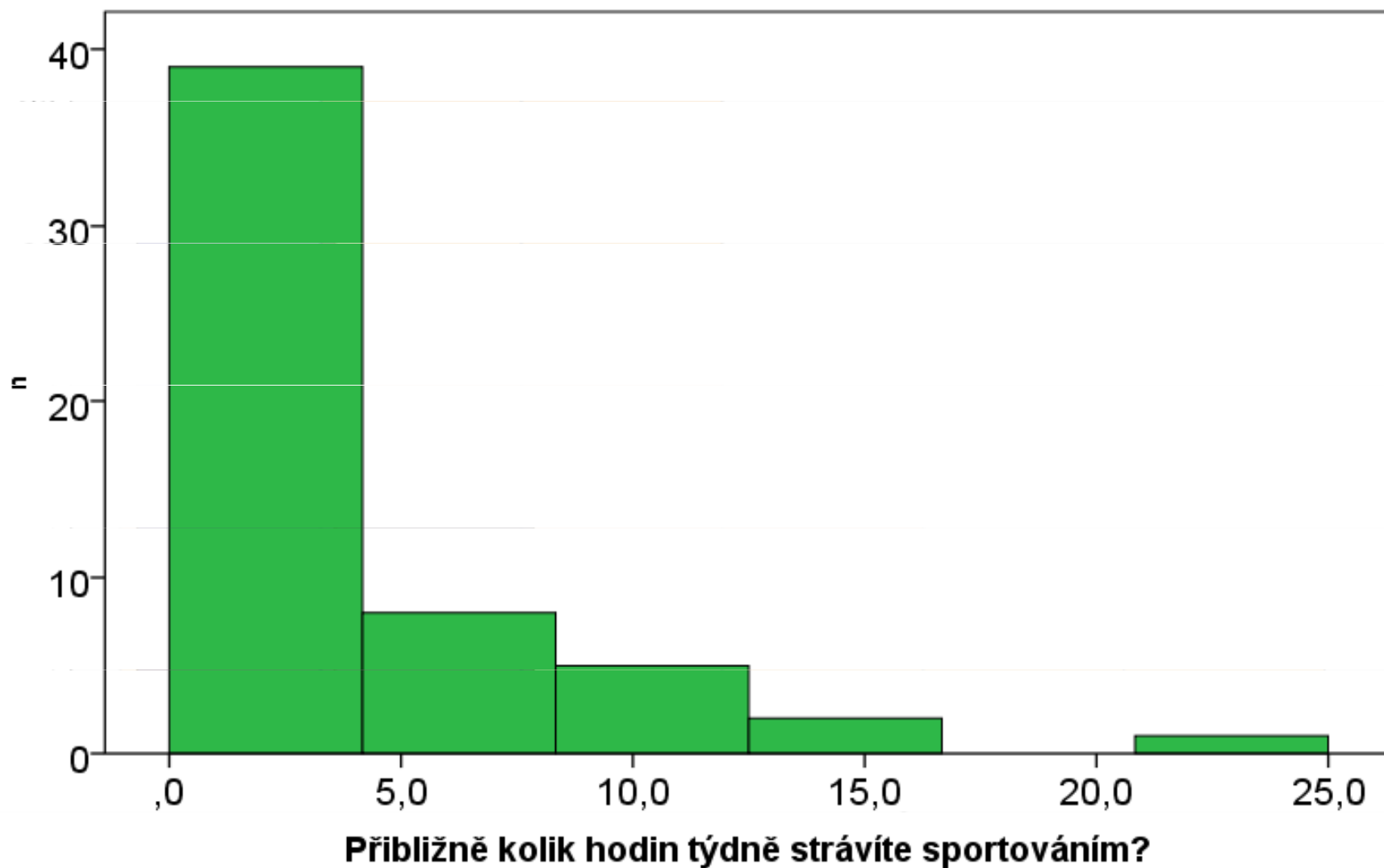


- **Empirické rozložení (rozdělení)** četností je tedy to, kolik relativně (či absolutně) máme kterých hodnot měřené proměnné.
 - Typicky lze přibližně popsat slovy, např.: vyskytlo se hodně středních hodnot a relativně málo extrémních hodnot.
 - Toto **rozložení** jevů na měřené škále je nejlépe vidět na grafech.
 - Obvykle nějaké rozložení očekáváme – **teoretické** rozložení.
 - Mluvíme o „rozložení/rozdělení proměnné“

Histogram s relativními četnostmi (%)

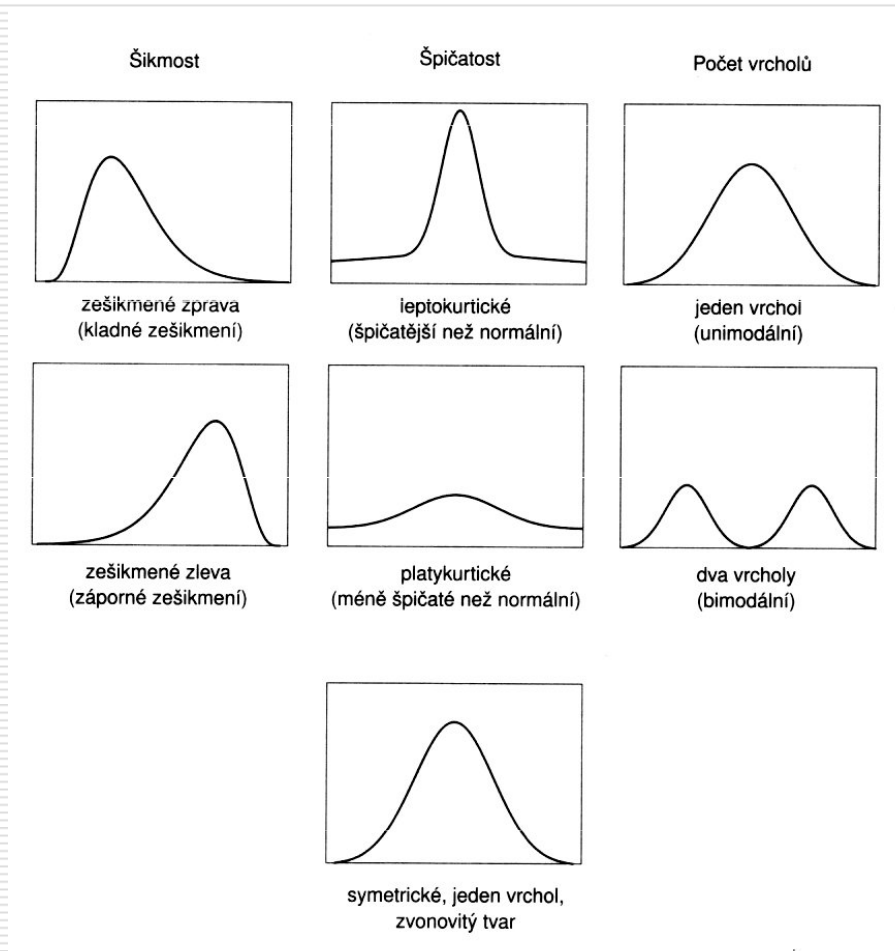


Histogram s širšími intervaly

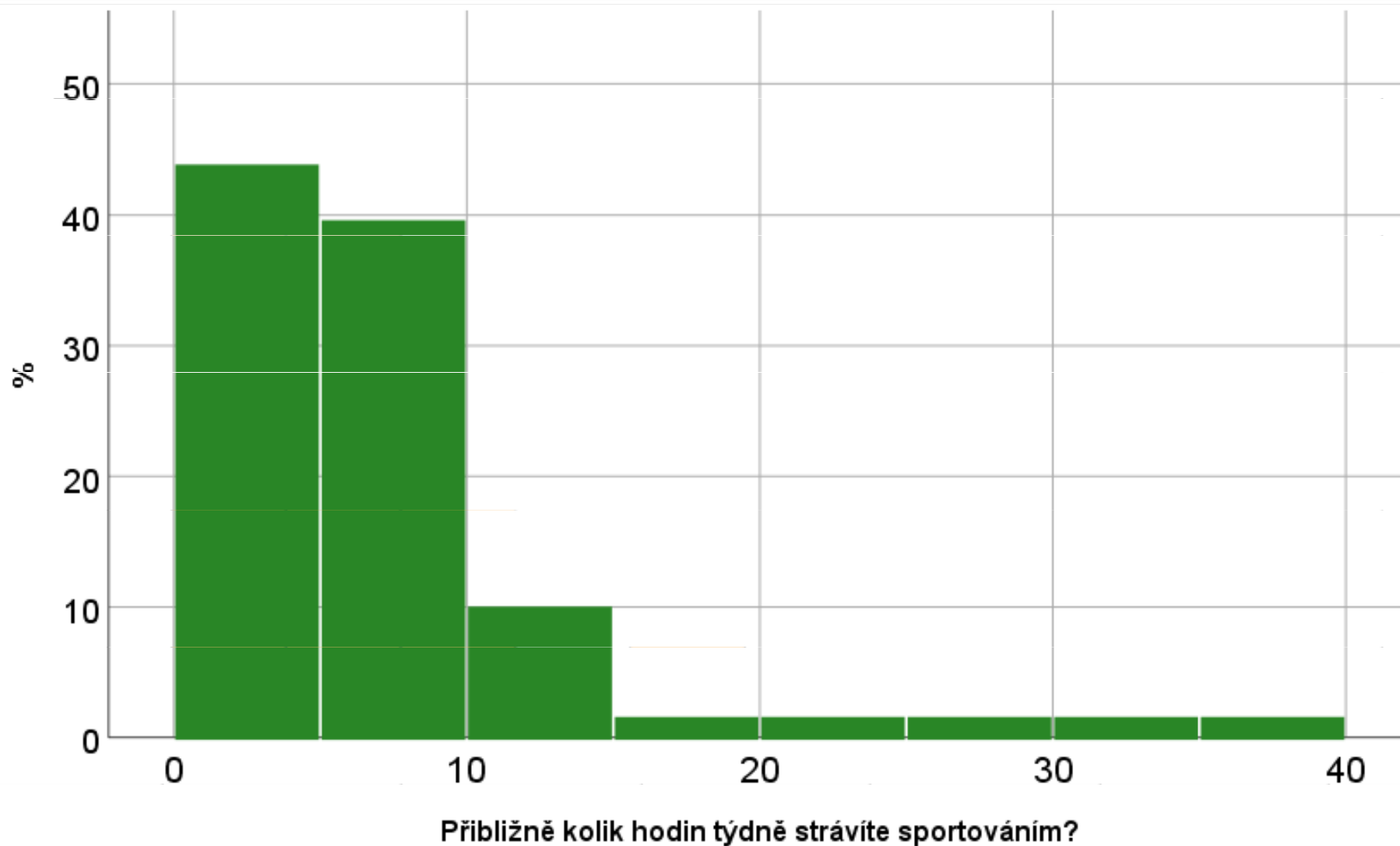


Tvary rozložení

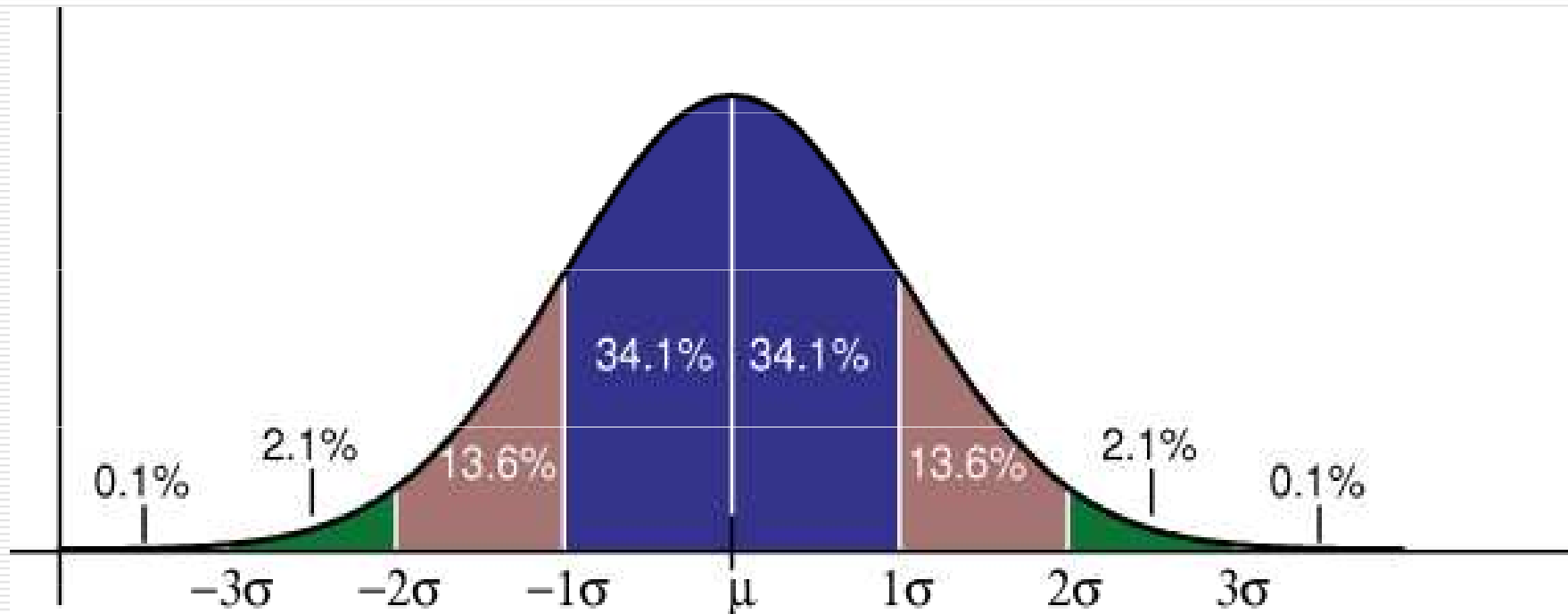
- Normální
- Uniformní
- Podle počtu vrcholů
 - Unimodální, bimodální, multimodální
- Zešikmení
 - Zešikmené zprava (pozitivně)
 - Zešikmené zleva (negativně)
- Strmost
 - Leptokurtické, platykurtické



Histogram s relativními četnostmi (%)



Normální (Gaussovo) rozložení

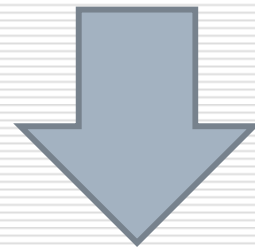


http://en.wikipedia.org/wiki/Image:Standard_deviation_diagram.png

- „Normální“ ve smyslu „velmi běžné“
- Tam, kde se setkává mnoho nezávislých vlivů.
- Ne vždy, nesouvisí s „kvalitou“ dat.

Změna „první otázky“

Jaké hodnoty proměnné se vyskytují a jak často?



Jaké je rozložení proměnné?

Cíl: popsat rozložení

Parametrický popis rozložení

- **Rozložení** je úplně popsáno (určeno) četnostmi jednotlivých hodnot, popř. intervalů.
- Je tedy popsáno množstvím statistik (četností), přesněji $k-1$ četnostmi, pokud proměnná nabývá k hodnot (či k intervalů).

Lze rozložení popsat efektivněji, méně statistikami (**parametry**)?

- Všechny hodnoty jsou stejně četné (1 parametr)
 - $f_k = k/N$ kde k je konstanta **UNIFORMNÍ** rozložení
- Četnosti jsou výsledkem procesu, který se dá připodobnit k opakovanému házení korunou, kdy nás zajímá počet „hlav“
 - $p_k = p^k(1-p)^{n-k} \frac{n!}{(n-k)!}$ kde n = počet hodů, k = počet hlav
 p = pravděpodobnost „hlavy“
 - **BINOMICKÉ** rozložení pro diskrétní proměnné (2 parametry)

- Normální rozložení

$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Nedalo by se rozložení hodnot proměnné popsat úsporněji než pomocí tabulky četností, histogramu?

Kde na měřené škále se rozložení nalézá?

UKAZATEL CENTRÁLNÍ TENDENCE

Jak moc jsou hodnoty proměnné rozptýlené?

UKAZATEL VARIABILITY

+ tvar rozložení (často implicitně)

Centrální tendence (=střední hodnoty, umístění)

- CT je jeden údaj, jímž se snažíme popsat rozložení četností jedné proměnné
 - Kouzlo i zrádnost je právě v tom, že je to 1 údaj.
- CT udává průměrnou, typickou, reprezentativní, *očekávanou* hodnotu
 - Co přesně tím míníme, záleží na tom, jakou míru CT se rozhodneme použít
- CT udává, kde na číselné ose si představujeme rozložení proměnné – odtud ukazatel *lokace*, umístění

Modus, medián a průměr

Modus - **kategorická** typická hodnota

\hat{X}, Mo

- nejčastější hodnota, h. s nejvyšší četností
- jediná možnost u nominálních dat, u vyšších úrovní často užitečnou volbou

Medián – **pořadová** střední hodnota

- hodnota prvku uprostřed uspořádaného souboru, 50. percentil (P_{50})
- při sudém počtu prvků je mediánem kterékoli číslo z intervalu mezi nejbližší vyšší a nejbližší nižší hodnotou (konsensuálně střed intervalu)
- hodnota minimalizující sumu absolutních odchylek
- pořadová data a výše

\tilde{X}, Md

Aritmetický průměr – deviační, odchylková, **momentová** střední h.

- jak ho znáte ze školy
- pouze intervalová a poměrová data
- velmi citlivý na extrémní hodnoty
- hodnota minimalizující sumu kvadratických odchylek

\bar{X}, M, m

Jak spočítat M_o , M_d , M

M_o

- vyčteme z tabulky četností – hodnota/interval s nejvyšší četností
- Excel: =MODE(rozsah_s_daty_proměnné)

M_d

- kategorické p.: vyčteme z tabulky četností, nejsnáze kum.
- spojité p. --> intervalové četnosti --> interpolujeme...
- formálně,
 - je-li N liché, je to X_k (k -tý prvek seříděné řady hodnot proměnné), kde $k=(N+1)/2$,
 - je-li N sudé, je to průměr X_k a X_{k+1} , kde $k=N/2$
- Excel: =MEDIAN(rozsah_s_daty_proměnné)
- =PERCENTIL(rozsah_s_daty_proměnné;0,5)

M

- Excel: =PRŮMĚR(rozsah_s_daty_proměnné)
-

Medián u intervalových četností a spojitéch proměnných s celými hodnotami - *interpolací*

	f	%	cum %
0 – 1>	3	11,5	11,5
1 – 2>	4	19,2	30,8
2 – 3>	9	30,8	61,5
3 – 4>	2	7,7	69,2
...
7 – 8>	1	3,8	100,0
Celkem	26	100,0	

1. Identifikujeme interval, v němž kumulativní četnost přesáhne 50% $(2;3>$
2. Četnost tohoto intervalu = $f_m = 9$
3. Kumulativní četnost pro předchozí interval = $f_p = 7$
4. Horní mez předchozího intervalu = $L_p = 2$
5. Šířka intervalu = $W = 1$
6. Vypočítáme medián

$$Md = L_p + W((N+1)/2 - f_p) / f_m =$$

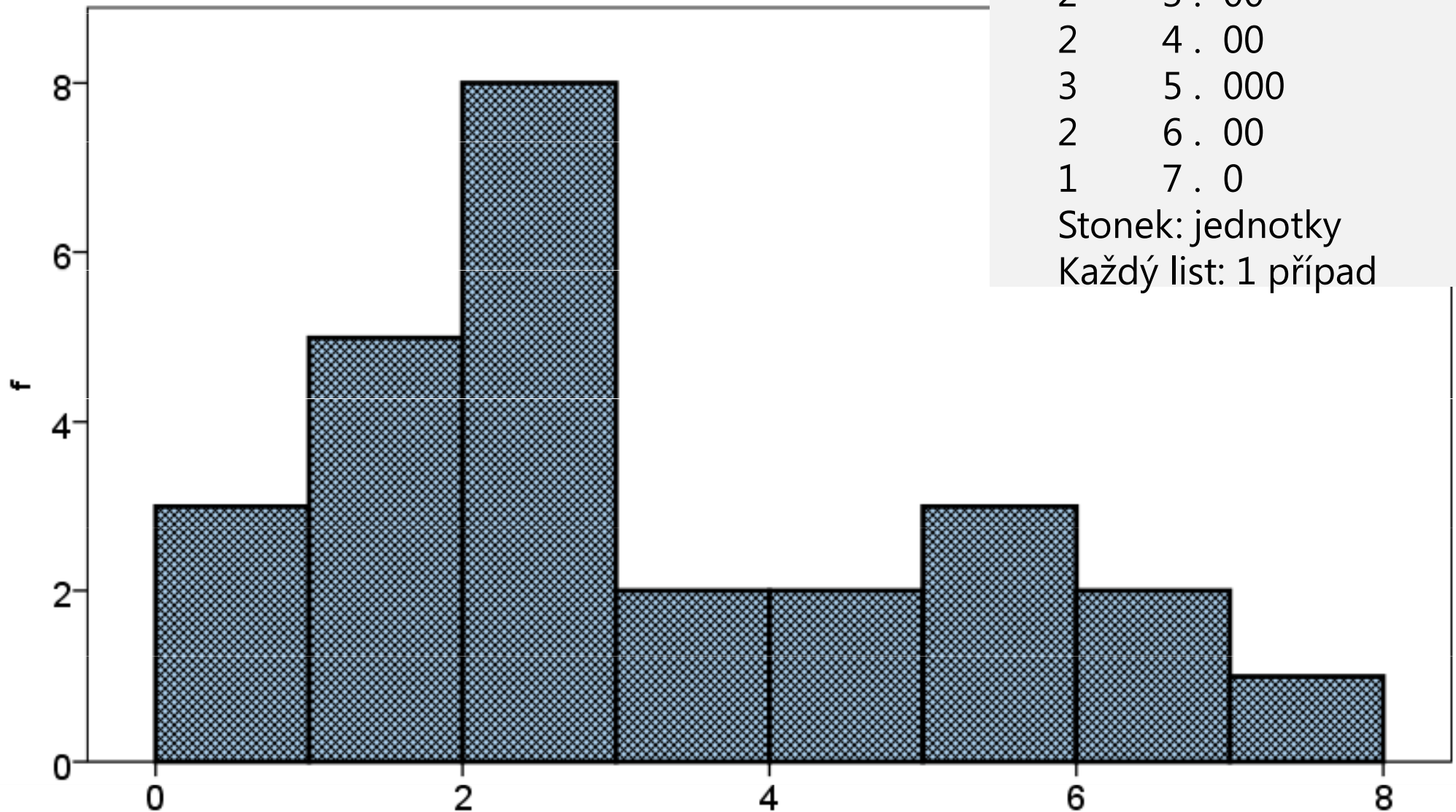
$$= 2 + 1(27/2 - 7) / 9 = 2,7$$

*Takto odhadnutý medián závisí na tom, jak jsou stanoveny hranice intervalů.

$M_o=2$ $M_d=2$ $M=2,68$
 $M_d=2,7$

f	Stonek a list
3	0 . 002
5	1 . 00005
8	2 . 00000000
2	3 . 00
2	4 . 00
3	5 . 000
2	6 . 00
1	7 . 0

Stonek: jednotky
Každý list: 1 případ



Přibližně kolik hodin týdně strávíte sportováním?

Míry variability (rozptýlenosti)

- Druhé číslo, jímž popisujeme rozložení hodnot proměnné
 - Udává, jak moc či málo jsou data na škále rozptýlená.
 - Malá variabilita = většina hodnot v souboru je stejných nebo velmi blízkých
 - Vysoká variabilita = hodnoty jsou velmi rozmanité (n. rozložení je bimodální)
-

Rozpětí, rozptyl, směrodatná odchylka

Nominální statistika – entropie – nepoužívá se

Pořadové statistiky

- (variační) rozpětí = $X_{max} - X_{min}$ (extrémně roste s velikostí vzorku)
- (inter)**kvartilové rozpětí** = $Q_3 - Q_1$, IQR

Odchylkové (deviační, momentové) statistiky

- založené na odchylkách od průměru: $x = X - m$
- průměrná absolutní odchylka MAD = $(\sum|x| / n)$ – řídká, ale ...
- průměrná odchylka na druhou – **rozptyl** – s^2 , $VAR(X)$
 - populační $(\sum x^2 / n)$ vs. výběrový $(\sum x^2 / (n - 1))$
 - součet odchylek na druhou = **suma čtverců**
- **směrodatná odchylka** (standardní odchylka) – s , SD
 - odmocnina rozptylu - návrat k původní jednotce

Směrodatná odchylka

- lze interpretovat přibližně jako
 - Průměrná odchylka od průměru
 - Očekávaná odchylka od průměru
-

Jak spočítat ukazatele variability

□ $IQR = Q_3 - Q_1$

- $Q_1 = X_k^*)$, kde $k = (N+1) * 0,25$ zaokrouhleno dolů
- $Q_3 = X_k'$, kde $k = (N+1) * 0,75$ zaokrouhleno dolů
- =PERCENTIL(rozsah_s_daty_proměnné; 0,25) resp. 0,75

U spojitých proměnných lze využít intervalového výpočtu jako u mediánu.

□ SD/VAR

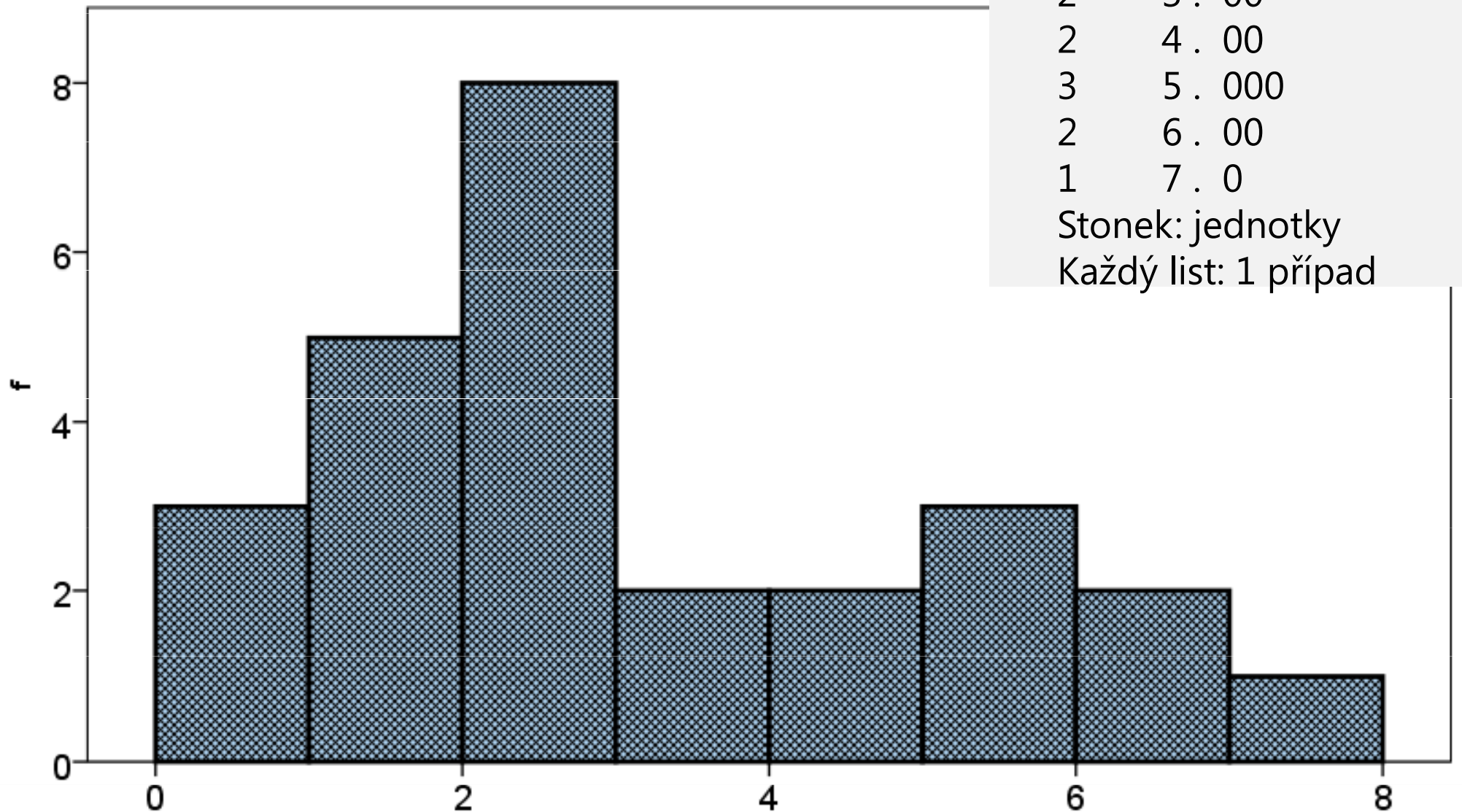
1. pro každý skór spočítáme deviační skór $x_i = X_i - M$
 2. deviační skóry umocníme na druhou
 3. druhé mocniny deviačních skórů sečteme a podělíme $(N-1)$
 4. pro SD výsledek ještě odmocníme
- =VAR.VÝBĚR(rozsah_s_daty_proměnné)
 - =SMODCH.VÝBĚR(rozsah_s_daty_proměnné)

*) hodnota k-tého prvku seřazené řady hodnot proměnné X

$M_o=2$ $M_d=2$ $M=2,68$
 $IQR=3$ $SD=1,97$

f	Stonek a list
3	0 . 002
5	1 . 00005
8	2 . 00000000
2	3 . 00
2	4 . 00
3	5 . 000
2	6 . 00
1	7 . 0

Stonek: jednotky
 Každý list: 1 případ



Přibližně kolik hodin týdně strávíte sportováním?

Ukazatele centrální tendence a variability - poznámky

- je třeba je umět spočítat ručně (a zopakovat si práci se sumačním symbolem Σ)
- i vážený průměr
- jak je ovlivní datové transformace přičtení konstanty a násobení konstantou
- vhodnost použití ukazatelů centrální tendence (Hendl s.95)

Očekávaná hodnota a její chyba

Střední hodnoty a ukazatele variability lze také interpretovat z perspektivy pravděpodobnostních očekávání plynoucích ze statistik.

Když přijde někdo náhodný, jakou hodnotu budeme očekávat - hádat?

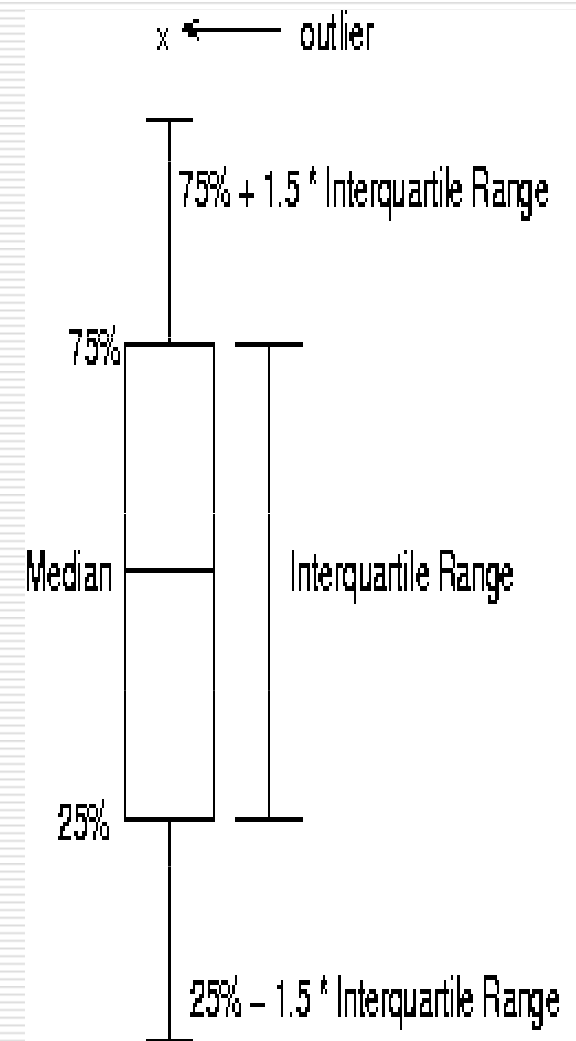
To záleží na tom, na čem nám nejvíc záleží.

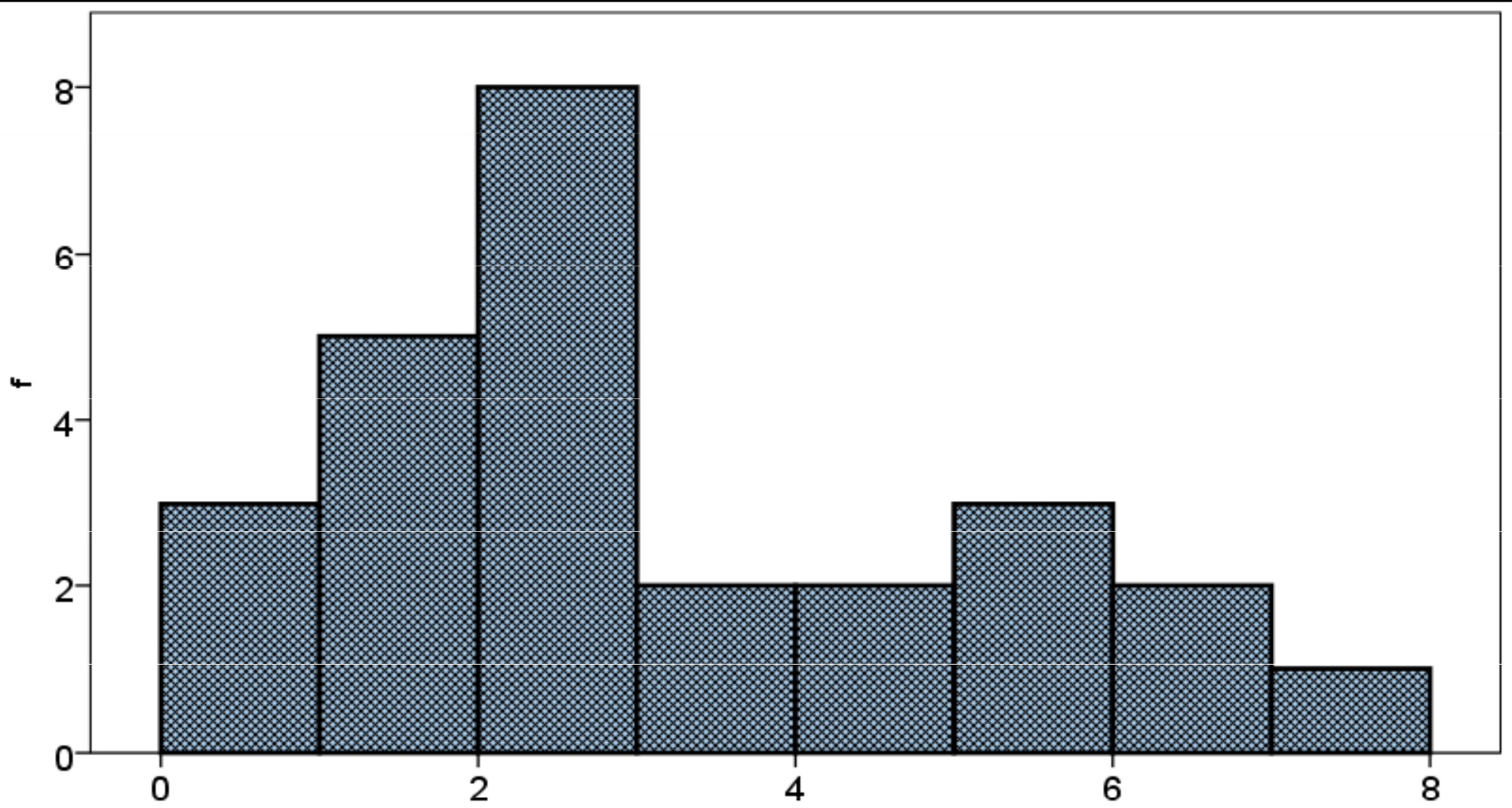
- ❑ Chceme-li se co nejčastěji přesně trefit (když za trefení dostaneme bod a ze netrefení ne), pak modus.
- ❑ Chceme-li minimalizovat součet velikostí chyb odhadu, pak medián
- ❑ Chceme-li minimalizovat součet kvadratických chyb odhadu, pak průměr

V tomto kontextu jsou ukazatele rozptýlenosti vlastně ukazateli velikosti chyb.

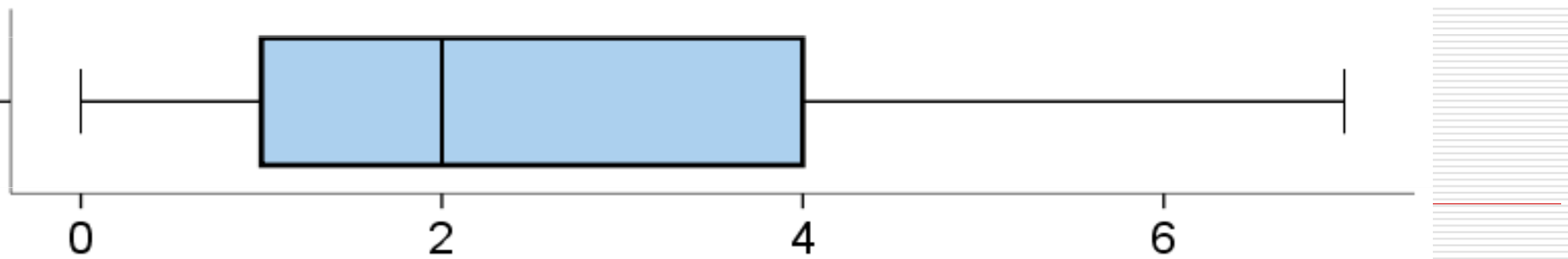
Boxplot – krabicový graf s anténami

- ❑ krabice je od Q_1 do Q_3
- ❑ v krabici se značí medián
- ❑ antény jsou X_{\min} do X_{\max} , **maximálně!** však 1,5x délka krabice (kvartilového rozpětí)
- ❑ hodnoty vzdálenější se značí jako body – odlehlé hodnoty
- ❑ hodnoty ještě vzdálenější (více než 3x délka krabice od Q_1 nebo Q_3) jsou někdy označovány jako extrémně odlehlé hodnoty





Přibližně kolik hodin týdně strávíte sportováním?



Popis rozložení pomocí percentilů

□ X -tý percentil

- hodnota, pro kterou platí, že X % lidí (jevů) ve vzorku má/získalo tuto nebo menší hodnotu
- lze odečíst z kumulativního histogramu či příčného sloupce tabulky četností

□ Rozložení popisujeme

- 10., 20., ..., 80., 90. percentilem – obecně
- min, 25., 50., 75., max – nejčastěji (...boxplot)
- min., 1., 5., 10., 25., 50., 75., 90., 95., 99. – v normách

□ Lze uvažovat v ještě menších částech rozložení než jsou procenta - obecně **kvantily**

„Deskriptivy“ – popisné statistiky – statistiky popisující rozložení

- Nominální deskriptivy
 - modus, (entropie)
 - Pořadové deskriptivy
 - medián, kvartily, percentily (a jiné *kvantily*)
 - kvartilové rozpětí
 - grafické znázornění rozložení pomocí pořadových deskriptiv - **BOXPLOT**
 - Odchylkové (deviační), momentové deskriptivy
 - aritmetický průměr
 - rozptyl, směrodatná odchylka ($k=2$)
 - zešikmení ($k=3$) $= (\sum x^k) / n$
 - špičatost (strmost) ($k=4$)
-

Volba popisných statistik

- Zvažujeme
 - úroveň měření
 - tvar rozložení – symetrie, normalita
 - cíl studie – pouze popis X usuzování, porovnávání
 - Podle komunikačních cílů...
 - Je-li cílem především deskripce dat(=rozložení), pak použijeme **POŘADOVÉ** ukazatele. Připojíme-li i odchylkové, nic nezkazíme.
 - ***N, min, Q₁, Md, Q₃, max***
 - **boxplot**
 - pro individuální skóry **percentily**
 - Je-li cílem další usuzování, porovnávání apod., používáme **ODCHYLKOVÉ** ukazatele ... pokud to úroveň měření dovoluje
 - ***N, m, s*** (*N, M, SD*)
 - popis rozložení
 - pro individuální skóry **z-skóry**
-

Prezentace deskriptiv ve studiích

- **Vždy!** Bez ohledu na to, jak složité statistiky následují.
 - Popis rozložení
 - Obvykle se neuvádějí tabulky četností a jejich grafické podoby, pokud ovšem není cílem studie právě statistická deskripce (např. manuál k testu inteligence).
 - Tvar rozložení obvykle podle potřeby zmiňujeme verbálně („přibližně normální, zleva zešikmené...“). Většinou se řeší pouze normalita a odchylky od ní.
 - Obvykle pouze pro proměnné, s nimiž pracujeme (interpretujeme...)
 - Minimální triáda: N, m, s (či jejich pořadové ekvivalenty Q_1, Md, Q_3, IQR)
 - Vhodná pětice: $N, X_{\min}, X_{\max}, m, s$
 - V případě potřeby: $N, X_{\min}, X_{\max}, m, s$, zešikmení, špičatost, zajímavé kvantily
 - Obvykle na 2-3 významné číslice (1-2 desetinná místa)
 - V českém textu česky, v anglickém anglicky!
 - Pozor na konvence spojené s jazykem: značky, desetinné tečky, chybějící nuly
 - Podoba tabulek je podchycena i normami, např. **publikační manuál APA**
-

Shrnutí

- Spíše než jednotlivé četnosti nás zajímá vzorec četností – **rozložení**
 - Neparametrický popis rozložení - tabulka četností, sloupcový dg./histogram
 - Parametrický popis četností
 - popisnými statistikami s ujištěním o tvaru rozložení
 - boxplotem
-