

## Průzkumová analýza jednorozměrných dat, diagnostické grafy

### Motivace

Průzkumová analýza dat je odvětví statistiky, které pomocí různých postupů odhaluje zvláštnosti v datech. Při zpracování dat se často používají metody, které jsou založeny na předpokladu, že data pocházejí z nějakého konkrétního rozložení, nejčastěji normálního. Tento předpoklad nemusí být vždy splněn, protože data

- mohou pocházet z jiného rozložení
- mohou být zatížena hrubými chybami
- mohou pocházet ze směsi několika rozložení.

Proto je důležité provést průzkumovou analýzu dat, abychom se vyvarovali neadekvátního použití statistických metod.

Data zkoumáme pomocí **funkcionálních** a **číselných charakteristik** a pomocí **diagnostických grafů**.

### Osnova:

- datový soubor
- bodové a intervalové rozložení četností
- typy znaků, číselné charakteristiky znaků
- krabicový diagram, N-P plot, P-P plot, Q-Q plot, histogram

## Funkcionální charakteristiky datového souboru

### Označení

Na množině objektů  $\{\varepsilon_1, \dots, \varepsilon_n\}$  zjišťujeme hodnoty znaku  $X$  (např. u 6 domácností zjišťujeme počet členů). Hodnotu znaku  $X$  na objektu  $\varepsilon_i$  označíme  $x_i$ ,  $i = 1, \dots, n$ .

Tyto hodnoty zaznamené do **jednorozměrného datového souboru**  $\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$  (např.  $\begin{pmatrix} 2 \\ 1 \\ 2 \\ 3 \\ 1 \\ 2 \end{pmatrix}$ ).

Uspořádané hodnoty  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$  tvoří **uspořádaný datový soubor**  $\begin{pmatrix} x_{(1)} \\ \vdots \\ x_{(n)} \end{pmatrix}$ , v našem případě  $\begin{pmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 2 \\ 3 \end{pmatrix}$ .

Vektor  $\begin{pmatrix} x_{[1]} \\ \vdots \\ x_{[r]} \end{pmatrix}$ , kde  $x_{[1]} < \dots < x_{[r]}$  jsou navzájem různé hodnoty znaku  $X$ , se nazývá **vektor variant**, v našem případě  $\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$ .

## Bodové rozložení četností

Je-li počet variant znaku X malý, přiřazujeme četnosti jednotlivým variantám a hovoříme o bodovém rozložení četností.

$n_j$  – absolutní četnost varianty  $x_{[j]}$

$$p_j = \frac{n_j}{n} \text{ – relativní četnost varianty } x_{[j]}$$

$N_j = n_1 + \dots + n_j$  – absolutní kumulativní četnost prvních  $j$  variant

$$F_j = \frac{N_j}{n} = p_1 + \dots + p_j \text{ – relativní kumulativní četnost prvních } j \text{ variant}$$

Absolutní a relativní četnosti zapisujeme do tabulky rozložení četností nebo je znázorňujeme graficky např. pomocí sloupkového diagramu či polygonu četností.

$$\text{Četnostní funkce: } p(x) = \begin{cases} p_j \text{ pro } x = x_{[j]}, j=1, \dots, r \\ 0 \text{ jinak} \end{cases}$$

$$\text{Empirická distribuční funkce: } F(x) = \begin{cases} 0 \text{ pro } x < x_{[1]} \\ F_j \text{ pro } x_{[j]} \leq x < x_{[j+1]}, j=1, \dots, r-1 \\ 1 \text{ pro } x \geq x_{[r]} \end{cases}$$

**Příklad 1.:** U 30 domácností byl zjišťován počet členů.

Počet členů	1	2	3	4	5	6
Počet domácností	2	6	4	10	5	3

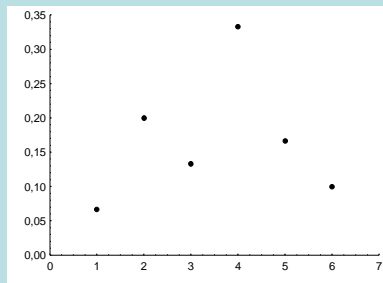
Vytvořte tabulku rozložení četností. Nakreslete grafy četnostní funkce a empirické distribuční funkce. Dále nakreslete sloupkový diagram a polygon četností.

**Řešení:**

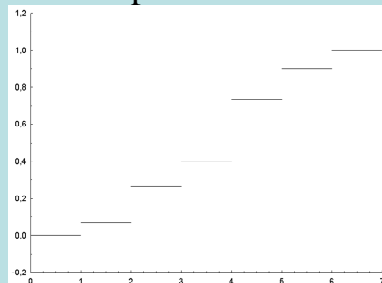
Tabulka rozložení četností

$x_{[j]}$	$n_j$	$p_i$	$N_i$	$F_i$
1	2	2/30	2	2/30
2	6	6/30	8	8/30
3	4	4/30	12	12/30
4	10	10/30	22	22/30
5	5	5/30	27	27/30
6	3	3/30	30	1

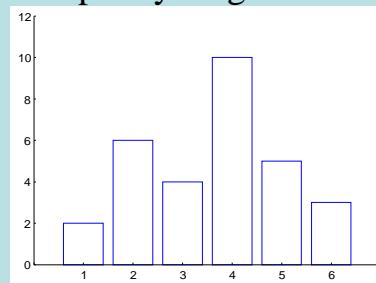
Graf četnostní funkce



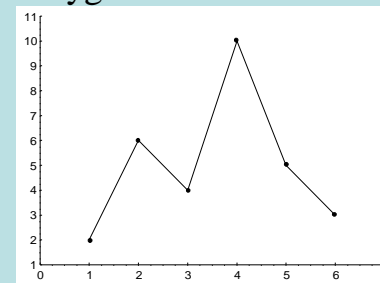
Graf empirické distribuční funkce



Sloupkový diagram



Polygon četností



## Intervalové rozložení četností

Je-li počet variant znaku  $X$  velký, přiřazujeme četnosti nikoli jednotlivým variantám, ale třídícím intervalům  $(u_1, u_2)$ , ...,  $(u_r, u_{r+1})$  a hovoříme o intervalovém rozložení četností. Názvy četností jsou podobné jako u bodového rozložení četností, navíc zavádíme **četnostní hustotu**  $j$ -tého třídícího intervalu  $f_j = \frac{p_j}{d_j}$ , kde  $d_j = u_{j+1} - u_j$ . Stanovení počtu třídících intervalů je dosti subjektivní záležitost. Často se doporučuje volit  $r$  blízké  $\sqrt{n}$ .

**Hustota četnosti:**  $f(x) = \begin{cases} f_j & \text{pro } u_j < x \leq u_{j+1}, j=1, \dots, r \\ 0 & \text{jinak} \end{cases}$  (grafem hustoty četnosti je histogram)

**Intervalová empirická distribuční funkce:**  $F(x) = \int_{-\infty}^x f(t) dt$ .

**Příklad 2.:** U 70 domácností byly zjišťovány týdenní výdaje na nealkoholické nápoje (v Kč).

Výdaje	(35,65)	(65,95)	(95,125)	(125,155)	(155,185)	(185,215)
Počet dom.	7	16	27	14	4	2

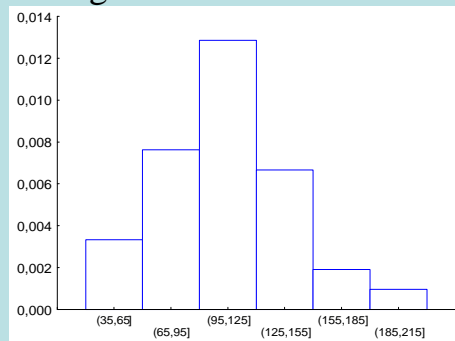
Sestavte tabulku rozložení četností, nakreslete histogram a graf intervalové empirické distribuční funkce.

**Řešení:**

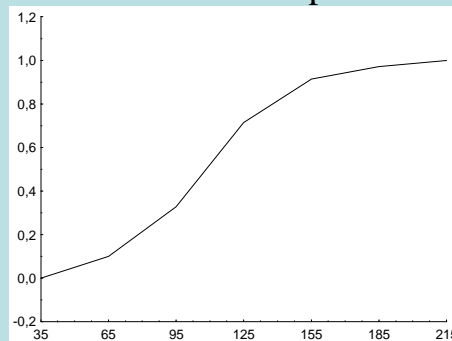
Tabulka rozložení četností

$(u_j, u_{j+1}]$	$n_j$	$p_j$	$f_j$	$N_j$	$F_j$
(35,65)	7	7/70	7/2100	7	7/70
(65,95)	16	16/70	16/2100	23	23/70
(95,125)	27	27/70	27/2100	50	50/70
(125,155)	14	14/70	14/2100	64	64/70
(155,185)	4	4/70	4/2100	68	68/70
(185,215)	2	2/70	2/2100	70	1

Histogram



Graf intervalové empirické distribuční funkce



## Číselné charakteristiky datového souboru

### Znaky nominálního typu

Tyto znaky umožňují obsahovou interpretaci pouze u relace rovnosti.

Příklady nominálních znaků: lékařská diagnóza, typ profese, barva očí, rodinný stav, národnost, ...

Charakteristikou polohy je **modus**, tj. nejčetnější varianta či střed nejčetnějšího intervalu.

### Znaky ordinálního typu

Lze u nich navíc obsahově interpretovat relaci uspořádání.

Příklad ordinálního znaku: školní klasifikace vyjadřuje menší nebo větší znalosti zkoušených žáků – jedničkář je lepší než dvojkař, ale intervaly mezi známkami nemají obsahovou interpretaci. Nelze tvrdit, že rozdíl ve znalostech mezi jedničkářem a dvojkařem je stejný jako mezi trojkařem a čtyřkařem.

Další příklady: Různá bodování ve sportovních a uměleckých soutěžích, posuzování různých rysů sociálního chování, posuzování stavu pacientů, hodnocení postojů respondentů k různým otázkám, ...

Charakteristikou polohy je  **$\alpha$ -kvantil**. Je-li  $\alpha \in (0;1)$ , pak  $\alpha$ -kvantil  $x_\alpha$  je číslo, které rozděluje uspořádaný datový soubor na dolní úsek, obsahující aspoň podíl  $\alpha$  všech dat a na horní úsek obsahující aspoň podíl  $1 - \alpha$  všech dat. Pro výpočet  $\alpha$ -kvantilu slouží algoritmus:

$$n\alpha = \begin{cases} \text{celé číslo } c \Rightarrow x_\alpha = \frac{x_{(c)} + x_{(c+1)}}{2} \\ \text{necelé číslo} \Rightarrow \text{zaokrouhlíme nahoru na nejbližší celé číslo } c \Rightarrow x_\alpha = x_{(c)} \end{cases}$$

Pro speciálně zvolená  $\alpha$  užíváme názvů:

$x_{0,50}$  – **medián**,  $x_{0,25}$  – **dolní kvartil**,  $x_{0,75}$  – **horní kvartil**,  $x_{0,1}, \dots, x_{0,9}$  – **decily**,  $x_{0,01}, \dots, x_{0,99}$  – **percentily**.

Jako charakteristika variability slouží **kvartilová odchylka**:  $q = x_{0,75} - x_{0,25}$ .

**Příklad 3.:** Během semestru se studenti podrobili písemnému testu z matematiky, v němž bylo možno získat 0 až 10 bodů. Výsledky jsou uvedeny v tabulce:

Počet bodů	0	1	2	3	4	5	6	7	8	9	10
Počet studentů	1	4	6	7	11	15	19	17	12	6	3

Zjistěte modus, medián, 1. decil, 9. decil a kvartilovou odchylku počtu bodů.

### Řešení:

Modus je nejčetnější varianta znaku, v tomto případě tedy 6.

Pro výpočet kvantilů musíme znát rozsah datového souboru:  $n = 1 + 4 + \dots + 3 = 101$ . Výpočty uspořádáme do tabulky.

$\alpha$	$n\alpha$	c	$x_{\alpha} = X_{(c)}$
0,50	50,5	51	6
0,10	10,1	11	2
0,90	90,9	91	8
0,25	25,25	26	4
0,75	75,75	76	7

$$q = 7 - 4 = 3$$

### Výpočet pomocí systému STATISTICA:

Otevřeme nový datový soubor o 2 proměnných a 11 případech. První proměnnou nazveme X, druhou četnost a zapíšeme do nich počet bodů a odpovídající absolutní četnosti.

Statistiky – Základní statistiky/tabulky – Popisné statistiky – zapneme proměnnou vah četnost – OK – OK – Proměnné X – OK – Detailní výsledky – vybereme Medián, Dolní a horní kvartily, Kvantilové hranice – Výpočet – ve výstupní tabulce upravíme počet desetinných míst.

Proměnná	Popisné statistiky (pocet bodu.sta)					
	N platných	Medián	Spodní kvartil	Horní kvartil	Kvantil 10,00000	Kvantil 90,00000
X	101	6	4	7	2	8



## Znaky intervalového a poměrového typu

U těchto znaků lze navíc obsahově interpretovat operaci rozdílu resp. podílu.

Příklad intervalového znaku: teplota měřená ve stupních Celsia. Např. naměříme-li ve čtyřech po sobě jdoucích dnech polední teploty 0, 2, 4, 6 °C, znamená to, že každým dnem stouply teploty o 2 °C. Nelze však říci, že z druhého na třetí den vzrostla teplota dvojnásobně, kdežto ze třetího na čtvrtý den pouze jeden a půl krát.

Další příklady: kalendářní systémy, směr větru, inteligenční kvocient, ...

Společný znak intervalových znaků: nula byla stanovena uměle, pouhou konvencí.

Příklad poměrového znaku: délka předmětu měřená v cm. Má-li jeden předmět délku 8 cm a druhý 16 cm, má smysl prohlásit, že druhý předmět je dvakrát delší než první předmět.

Další příklady: počet dětí v rodině, výška kapesného v Kč, hmotnost osoby, ...

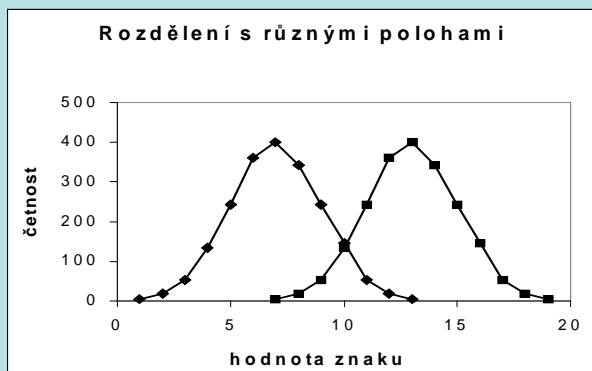
Společný znak poměrových znaků: poměrový znak má přirozený počátek, ke kterému jsou vztahovány všechny další hodnoty znaku.

Charakteristika polohy: **aritmetický průměr**  $m = \frac{1}{n} \sum_{i=1}^n x_i$ .

U poměrových znaků, které nabývají pouze kladných hodnot, lze použít **geometrický průměr**  $\sqrt[n]{x_1 \cdot \dots \cdot x_n}$ .

Pomocí průměru zavedeme **i-tou centrovanou hodnotu**  $x_i - m$  (podle znaménka poznáme, zda i-tá hodnota je podprůměrná či nadprůměrná).

Znázornění rozložení četností dvou datových souborů, které se liší aritmetickým průměrem



## Vlastnosti aritmetického průměru

- Aritmetický průměr si lze představit jako těžiště dat – součet podprůměrných hodnot je stejný jako součet nadprůměrných hodnot – oba součty jsou v rovnováze.

- Průměr centrovaných hodnot je nulový, protože  $\frac{1}{n} \sum_{i=1}^n (x_i - m) = \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{n} \sum_{i=1}^n m = m - \frac{1}{n} \cdot n \cdot m = 0 = 0$ .

- Výraz  $\sum_{i=1}^n (x_i - a)^2$  (tzv. kvadratická odchylka) nabývá svého minima pro  $a = m$ . Uvedený výraz charakterizuje celkovou chybu, které se dopustíme, když datový soubor nahradíme jedinou hodnotou  $a$ . Tato chyba je tedy nejmenší, když datový soubor nahradíme aritmetickým průměrem, přičemž za míru chyby považujeme kvadratickou odchylku.

- Pokud každou hodnotu  $x_i$  podrobíme lineární transformaci  $y_i = a + bx_i$ , pak průměr transformovaných hodnot je roven lineární transformaci původního průměru, tj.  $m_2 = a + bm_1$ .

- Mají-li znaky  $X$ ,  $Y$  průměry  $m_1$ ,  $m_2$ , pak znak  $Z = X + Y$  má průměr  $m_1 + m_2$ .

- Aritmetický průměr je silně ovlivněn extrémními hodnotami.

- Aritmetický průměr je vhodné použít, pokud je rozložení dat přibližně symetrické.

### Příklad na vlastnosti aritmetického průměru:

U skupiny 20 pracovníků v určité dílně byly zjišťovány měsíční mzdy. Průměr mezd činil 15 500 Kč. Určete průměr mezd, jestliže mzdy všech pracovníků se zvýší

a) o 300 Kč, b) 1,1 krát, c) o 20%.

### Řešení:

Označme  $m_1$  průměr hodnot  $x_1, \dots, x_n$  a  $m_2$  průměr hodnot  $y_1, \dots, y_n$ , přičemž  $y_i = a + bx_i$ ,  $i = 1, \dots, n$ . Pak  $m_2 = a + bm_1$ .

ad a)  $m_2 = 300 + m_1 = 15\ 800$

Průměr se zvýšil o 300 Kč na 15 800 Kč.

ad b)  $m_2 = 1,1 \cdot m_1 = 17\ 050$

Průměr se zvýšil na 17 050 Kč.

ad c)  $m_2 = 1,2 \cdot m_1 = 18\ 600$

Průměr se zvýšil na 18 600 Kč.

## Charakteristiky variability intervalových a poměrových znaků

**Variační rozpětí**  $R = x_{(n)} - x_{(1)}$  (nevýhoda – bere v úvahu pouze nejmenší a největší hodnotu datového souboru),

**rozptyl**  $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - m)^2$  (nevýhoda – vychází ve druhých mocninách jednotek, v nichž byl měřen znak X)

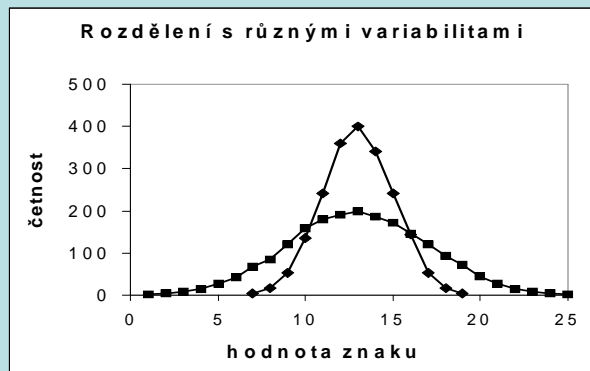
**směrodatná odchylka**  $s = \sqrt{s^2}$ .

Pomocí směrodatné odchylky zavedeme **i-tou standardizovanou hodnotu**  $\frac{x_i - m}{s}$  (vyjadřuje, o kolik směrodatných odchylek se i-tá hodnota odchýlila od průměru).

U poměrových znaků se jako charakteristika variability používá též:

**koefficient variace**  $\frac{s}{m}$  (často se udává v procentech a udává, kolika procent průměru dosahuje směrodatná odchylka),

Znázornění rozložení četností dvou datových souborů, které se liší rozptylem:



### Vlastnosti rozptylu:

- Rozptyl je nulový pouze tehdy, když jsou všechny hodnoty stejné, jinak je kladný.

- Rozptyl centrovaných hodnot je roven původnímu rozptylu, neboť  $\frac{1}{n} \sum_{i=1}^n [(x_i - m) - 0]^2 = \frac{1}{n} \sum_{i=1}^n (x_i - m)^2 = s^2$ .

- Rozptyl standardizovaných hodnot je 1, protože  $\frac{1}{n} \sum_{i=1}^n \left( \frac{x_i - m}{s} - 0 \right)^2 = \frac{1}{s^2} \cdot \frac{1}{n} \sum_{i=1}^n (x_i - m)^2 = \frac{s^2}{s^2} = 1$ .

- Rozptyl se zpravidla počítá podle vzorce  $s^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - m^2$ .

- Pokud každou hodnotu  $x_i$  podrobíme lineární transformaci  $y_i = a + bx_i$ , pak rozptyl transformovaných hodnot je roven původnímu rozptylu vynásobenému  $b^2$ , tj.  $s_2^2 = b^2 s_1^2$ .

- Rozptyl je stejně jako průměr silně ovlivněn extrémními hodnotami.

- Rozptyl se nehodí jako charakteristika variability, je-li rozložení dat nesymetrické.

**Příklad 4.:** Kurzy akcií společnosti AAA Auto Group v průběhu 23 dní v měsíci srpnu 2010 byly následující: 17,75; 17,74; 17,85; 17,59; 17,92; 17,98; 18,39; 18,25; 18,30; 18,00; 18,15; 18,15; 18,22; 18,40; 18,25; 17,95; 18,25; 18,23; 17,95; 17,90; 17,80; 17,87; 17,87. Vypočtěte charakteristiky variability.

**Řešení:**

Nejprve vypočítáme variační rozpětí:  $R = x_{(n)} - x_{(1)} = 18,4 - 17,59 = 0,81$ .

Před výpočtem dalších charakteristik variability musíme získat aritmetický průměr:  $m = \frac{1}{23}(17,75 + 17,74 + \dots + 17,87) = 18,033$ .

$$\text{Rozptyl: } s^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - m^2 = \frac{1}{23}(17,75^2 + 17,74^2 + \dots + 17,87^2) - 18,033^2 = 0,049$$

$$\text{Směrodatná odchylka: } s = \sqrt{s^2} = \sqrt{0,049} = 0,2213$$

$$\text{Koefficient variace: } \frac{s}{m} 100\% = \frac{0,2213}{18,033} 100\% = 1,23\%$$

## Výpočet pomocí systému STATISTICA:

Otevřeme nový datový soubor o jedné proměnné X a 23 případech. Do proměnné X zapíšeme zjištěné kurzy akcií. Statistika – Základní statistiky/tabulky – Popisné statistiky – OK – Proměnné X – OK – Detailní výsledky – vybereme Průměr, Rozptyl, Rozpětí – Výpočet. Ve výstupní tabulce přidáme za proměnnou Rozptyl tři nové proměnné nazvané rozptyl, směr. odch. a koef. variace. Do Dlouhého jména proměnné rozptyl napíšeme  $=v3*22/23$ , Dlouhého jména proměnné směr. odch. napíšeme  $=\text{sqrt}(v4)$  a do Dlouhého jména proměnné koef. variace napíšeme  $=100*v5/v1$ .

Proměnná	Průměr	Rozpětí	Rozptyl	rozptyl $=v3*22/23$	směr. odch. $=\text{sqrt}(v4)$	koef. variace $=100*v5/v1$
x	18,03304	0,810000	0,051231	0,049004	0,221367976	1,22756858

## Vážené číselné charakteristiky

Známe-li absolutní četnosti  $n_1, \dots, n_r$  či relativní četnosti  $p_1, \dots, p_r$  variant  $x_{[1]}, \dots, x_{[r]}$ , můžeme spočítat

**vážený průměr**  $m = \frac{1}{n} \sum_{j=1}^r n_j x_{[j]} = \sum_{j=1}^r p_j x_{[j]}$ ,

**vážený rozptyl**  $s^2 = \frac{1}{n} \sum_{j=1}^r n_j (x_{[j]} - m)^2 = \sum_{j=1}^r p_j (x_{[j]} - m)^2$  (výpočetní vzorec:  $s^2 = \frac{1}{n} \sum_{j=1}^r n_j x_{[j]}^2 - m^2 = \sum_{j=1}^r p_j x_{[j]}^2 - m^2$ ).



**Příklad 5.:** U 35 zaměstnanců byl zjištěn počet odpracovaných hodin za měsíc.

Počet odpracovaných hodin	184	185	186	187	188	189
Počet zaměstnanců	4	6	7	6	7	5

Vypočítejte průměr, směrodatnou odchylku a koeficient variace počtu odpracovaných hodin.

**Řešení:**

$$\text{Vážený průměr: } m = \frac{1}{n} \sum_{j=1}^r n_j x_{[j]} = \frac{1}{35} (4 \cdot 184 + 6 \cdot 185 + 7 \cdot 186 + 6 \cdot 187 + 7 \cdot 188 + 5 \cdot 189) = 186,6$$

$$\text{Vážený rozptyl: } s^2 = \frac{1}{n} \sum_{j=1}^r n_j x_{[j]}^2 - m^2 = \frac{1}{35} (4 \cdot 184^2 + 6 \cdot 185^2 + 7 \cdot 186^2 + 6 \cdot 187^2 + 7 \cdot 188^2 + 5 \cdot 189^2) - 186,6^2 = 2,5257$$

$$\text{Vážená směrodatná odchylka: } s = \sqrt{s^2} = \sqrt{2,5257} = 1,59 \text{ h} = 1 \text{ h } 35 \text{ min}$$

$$\text{Koeficient variace: } \frac{s}{m} 100\% = \frac{1,59}{186,6} 100\% = 0,85\%$$

Vidíme, že zaměstnanci odpracovali za měsíc v průměru 186,6 h, přičemž směrodatná odchylka dosahuje 0,85 % průměrné odpracované doby.

## Výpočet pomocí systému STATISTICA:

Otevřeme nový datový soubor o 2 proměnných a 6 případech. První proměnnou nazveme X, druhou četnost a zapíšeme do nich počet odpracovaných hodin a odpovídající počty zaměstnanců.

Statistiky – Základní statistiky/tabulky – Popisné statistiky – zapneme proměnnou vah četnost – OK – OK – Proměnné X – OK – Detailní výsledky – vybereme Průměr, Rozptyl – Výpočet. Ve výstupní tabulce přidáme za proměnnou Rozptyl dvě nové proměnné nazvané směr. odch. a koef. variace. Do Dlouhého jména proměnné směr. odch. napíšeme  $=\text{sqrt}(v2*34/35)$  a do Dlouhého jména proměnné koef. variace napíšeme  $=100*v3/v1$ .

Proměnná	Průměr	Rozptyl	směr.odch. $=\text{sqrt}(v2*34/35)$	koef. variace $=100*v3/v1$
X	186,6	2,6	1,5892496	0,851687888

Převod desetinných částí hodiny na minuty můžeme provést např. pomocí aplikace na adrese <http://www.prevody-jednotek.cz/>.

## Počáteční a centrální momenty

Aritmetický průměr a rozptyl jsou speciální případy momentů. Zavedeme

**k-tý počáteční moment**  $m'_k = \frac{1}{n} \sum_{i=1}^n x_i^k$ ,  $k = 1, 2, \dots$ ,

**k-tý centrální moment**

$$m_k = \frac{1}{n} \sum_{i=1}^n (x_i - m)^k, \quad k = 1, 2, \dots$$

Pomocí 3. a 4. počátečního momentu se definuje šikmost a špičatost.

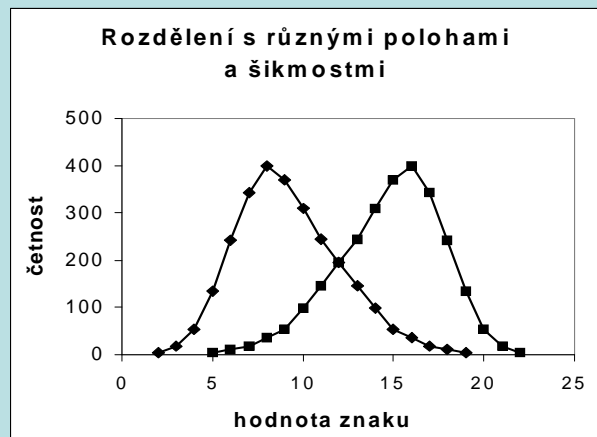
**Šikmost:**  $\alpha_3 = \frac{m_3}{s^3}$  - měří nesouměrnost rozložení četností kolem průměru.

Je-li rozložení dat symetrické kolem aritmetického průměru, pak  $\alpha_3 = 0$ .

Má-li rozložení dat prodloužený pravý konec, jde o **kladně zešikmené rozložení**,  $\alpha_3 > 0$ .

Má-li rozložení dat prodloužený levý konec, jde o **záporně zešikmené rozložení**,  $\alpha_3 < 0$ .

Znárodnění rozložení četností dvou datových souborů, které se liší aritmetickým průměrem a šikmostí



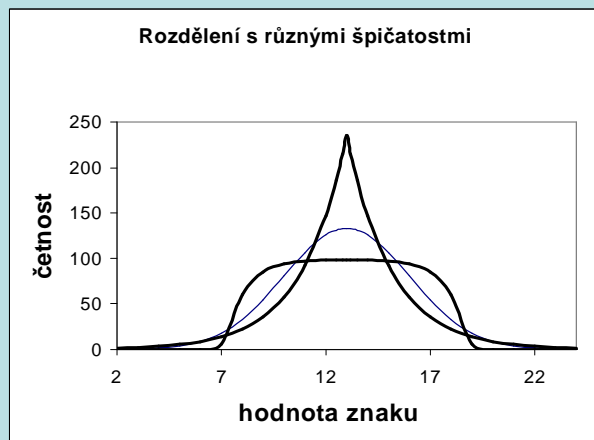
**Špičatost:**  $\alpha_4 = \frac{m_4}{s^4} - 3$  - měří koncentraci rozložení četností kolem průměru.

Je-li rozložení dat normální (Gaussovo), pak  $\alpha_4 = 0$ .

Je-li rozložení dat strmé, pak  $\alpha_4 > 0$ .

Je-li rozložení dat ploché, pak  $\alpha_4 < 0$ .

Znázornění rozložení četností dvou datových souborů, které se liší špičatostí

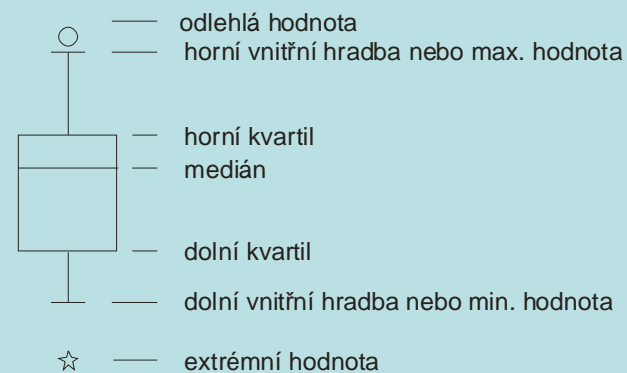


## Diagnostické grafy

### Krabicový diagram

Umožňuje posoudit symetrii a variabilitu datového souboru a existenci odlehlých či extrémních hodnot.

Způsob konstrukce



**Odlehlá hodnota** leží mezi **vnějšími a vnitřními hradbami**, tj. v intervalu

$(x_{0,75} + 1,5q, x_{0,75} + 3q)$  či v intervalu  $(x_{0,25} - 3q, x_{0,25} - 1,5q)$ .

**Extrémní hodnota** leží za vnějšími hradbami, tj. v intervalu  $(x_{0,75} + 3q, \infty)$  či v intervalu  $(-\infty, x_{0,25} - 3q)$ .

**Příklad 6.:** Pro údaje z příkladu 1 sestrojte krabicový diagram.

**Řešení:**

Počet členů	1	2	3	4	5	6
Počet domácností	2	6	4	10	5	3

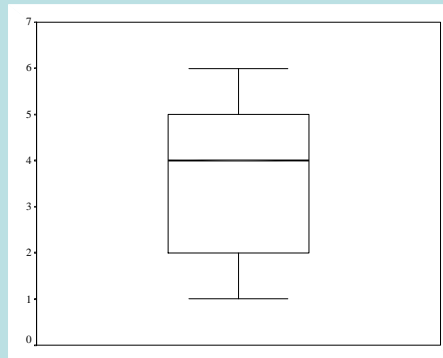
Rozsah souboru  $n = 30$ . Výpočty potřebných kvantilů uspořádáme do tabulky.

$\alpha$	$n\alpha$	$c$		$x_\alpha$
0,25	7,5	8	$x_{(c)}=x_{(8)}$	2
0,50	15	15	$\frac{x_{(15)} + x_{(16)}}{2}$	4
0,75	22,5	23	$x_{(c)}=x_{(23)}$	5

$$q = 5 - 2 = 3$$

Dolní vnitřní hradba:  $x_{0,25} - 1,5q = 2 - 1,5 \cdot 3 = -2,5$

Horní vnitřní hradba:  $x_{0,75} + 1,5q = 5 + 1,5 \cdot 3 = 9,5$

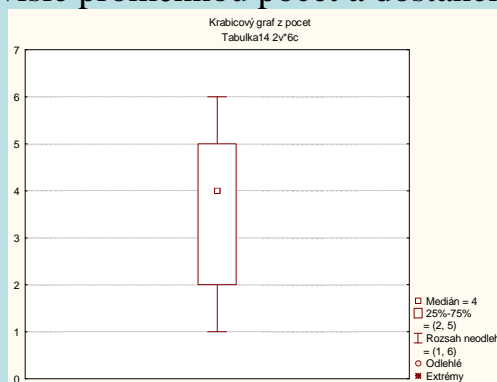


Vidíme, že datový soubor vykazuje určitou nesymetrii – medián je posunut směrem k hornímu kvartilu, soubor je tedy záporně sešikmen. V souboru se nevyskytují žádné odlehlé ani extrémní hodnoty.

## Výpočet pomocí systému STATISTICA:

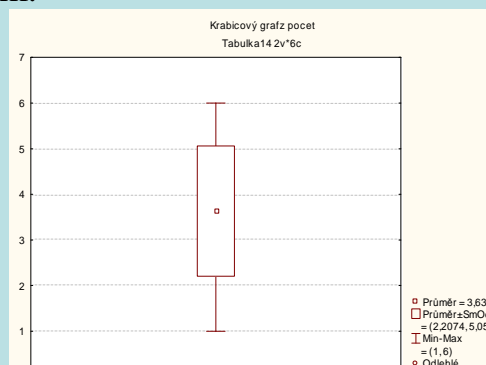
Otevřeme nový datový soubor o 2 proměnných a 6 případech. První proměnnou nazveme počet, druhou četnost a zapíšeme do nich počet členů domácnosti a odpovídající absolutní četnosti. Zvolíme Grafy – 2D Grafy – Krabicové grafy.

Zapneme proměnnou vah četnost, zadáme závisle proměnnou pocet a dostaneme krabicový diagram:



**Upozornění:** Máme-li data intervalového či poměrového charakteru, o nichž lze předpokládat, že pocházejí z nějakého symetrického rozložení (například normálního), je možné použít jinou variantu krabicového diagramu: bod či čára uvnitř krabice reprezentuje průměr, vodorovné hrany krabice jsou ve výšce průměr  $\pm$  směrodatná odchylka a svorky končí v minimu či maximu.

V našem případě dostaneme krabicový diagram:



Před uvedením dalších diagnostických grafů je nutné zavést pojem pořadí čísla v posloupnosti čísel.

### Pojem pořadí

Nechť  $x_1, \dots, x_n$  je posloupnost reálných čísel.

a) Jsou-li čísla navzájem různá, pak pořadím  $R_i$  čísla  $x_i$  rozumíme počet těch čísel  $x_1, \dots, x_n$ , která jsou menší nebo rovna číslu  $x_i$ .

b) Vyskytují-li se mezi danými čísly skupinky stejných čísel, pak každé takové skupince přiřadíme průměrné pořadí.

### Příklad na stanovení pořadí

a) Jsou dána čísla 9, 4, 5, 7, 3, 1. Stanovte pořadí těchto čísel.

b) Jsou dána čísla 6, 7, 7, 9, 6, 10, 8, 6, 6, 9.

### Řešení

ad a)

usp. čísla	1	3	4	5	7	9
pořadí	1	2	3	4	5	6

ad b)

usp. čísla	6	6	6	6	7	7	8	9	9	10
pořadí	1	2	3	4	5	6	7	8	9	10
prům. pořadí	2,5	2,5	2,5	2,5	5,5	5,5	7	8,5	8,5	10



## Normální pravděpodobnostní graf (N-P plot)

N- P plot umožňuje graficky posoudit, zda data pocházejí z normálního rozložení.

### Způsob konstrukce:

Na vodorovnou osu vynášíme uspořádané hodnoty  $x_{(1)} \leq \dots \leq x_{(n)}$ ,

na svislou osu kvantily  $u_{\alpha_j}$  standardizovaného normálního rozložení, kde  $\alpha_j = \frac{3j-1}{3n+1}$ , přičemž  $j$  je pořadí  $j$ -té uspořádané

hodnoty (jsou-li některé hodnoty stejné, pak za  $j$  bereme průměrné pořadí odpovídající takové skupince).

Pocházejí-li data z normálního rozložení, pak všechny dvojice  $(x_{(j)}, u_{\alpha_j})$  budou ležet na přímce.

Pro data z rozložení s kladnou šikmostí se dvojice  $(x_{(j)}, u_{\alpha_j})$  budou řadit do **konkávní křivky**,

pro data z rozložení se zápornou šikmostí se dvojice  $(x_{(j)}, u_{\alpha_j})$  budou řadit do **konvexní křivky**.

### Příklad na konstrukci N – P plotu:

Desetkrát nezávisle na sobě byla změřena jistá konstanta. Výsledky měření: 2 1,8 2,1 2,4 1,9 2,1 2 1,8 2,3 2,2. Pomocí normálního pravděpodobnostního grafu posuďte, zda se tato data řídí normálním rozložením.

### Řešení:

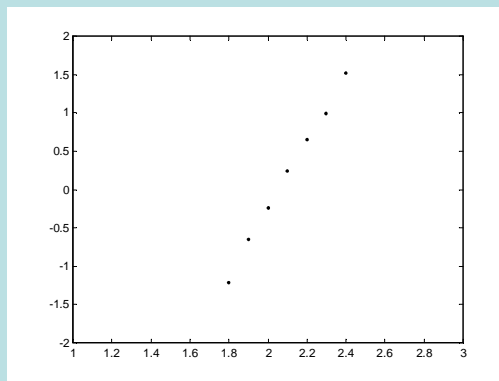
usp. hodnoty	1,8	1,8	1,9	2	2	2,1	2,1	2,2	2,3	2,4
pořadí	1	2	3	4	5	6	7	8	9	10
průměrné pořadí	1,5	1,5	3	4,5	4,5	6,5	6,5	8	9	10

Vektor hodnot průměrného pořadí:  $j = (1,5 \ 3 \ 4,5 \ 6,5 \ 8 \ 9 \ 10)$ ,

vektor hodnot  $\alpha_j = \frac{3j-1}{3n+1} = (0,1129; 0,2581; 0,4032; 0,5968; 0,7419; 0,8387; 0,9355)$ ,

vektor kvantilů  $u_{\alpha_j} = (-1,2112; -0,6493; -0,245; 0,245; 0,6493; 0,9892; 1,5179)$ .

Normální pravděpodobnostní graf

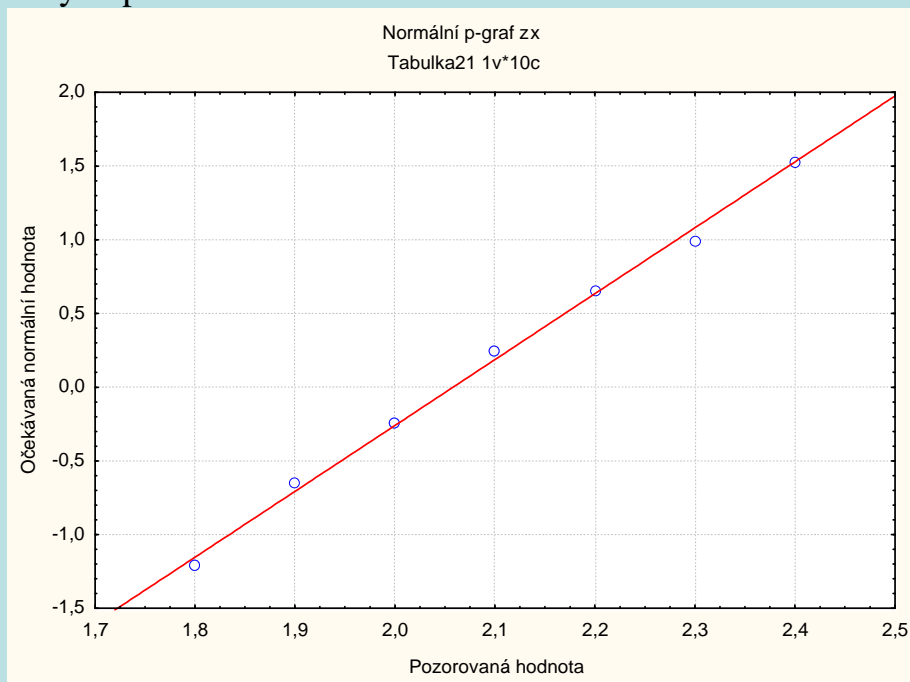


Protože dvojice  $(x_{(j)}, u_{\alpha_j})$  téměř leží na přímce, lze usoudit, že data pocházejí z normálního rozložení.

## Výpočet pomocí systému STATISTICA:

Otevřeme nový datový soubor o jedné proměnné a 10 případech. Zjištěné hodnoty zapíšeme do proměnné X.

Grafy – 2D Grafy – Normální pravděpodobnostní grafy – Proměnná X – OK - odškrtneme Neurčovat průměrnou pozici svázaných pozorování - OK.



## Quantile - quantile plot (Q-Q plot)

Umožňuje graficky posoudit, zda data pocházejí z nějakého známého rozložení (např. STATISTICA nabízí 8 typů rozložení: beta, exponenciální, Gumbelovo, gamma, log-normální, normální, Rayleighovo a Weibulovo).

### Způsob konstrukce:

na svislou osu vynášíme uspořádané hodnoty  $x_{(1)} \leq \dots \leq x_{(n)}$ ,

na vodorovnou osu kvantily  $K_{\alpha_j}(X)$  vybraného rozložení, kde  $\alpha_j = \frac{j - r_{adj}}{n + n_{adj}}$ , přičemž  $r_{adj}$  a  $n_{adj}$  jsou korigující faktory  $\leq 0,5$ ,

implicitně  $r_{adj} = 0,375$  a  $n_{adj} = 0,25$ . (Jsou-li některé hodnoty  $x_{(1)} \leq \dots \leq x_{(n)}$  stejné, pak za  $j$  bereme průměrné pořadí odpovídající takové skupince.)

Pokud vybrané rozložení závisí na nějakých parametrech, pak se tyto parametry odhadnou z dat nebo je může zadat uživatel.

Body  $(K_{\alpha_j}(X), x_{(j)})$  se metodou nejmenších čtverců proloží přímkou. Čím méně se body odchylují od této přímky, tím je lepší soulad mezi empirickým a teoretickým rozložením.

**Příklad na konstrukci Q-Q plotu:** Desetkrát nezávisle na sobě byla změřena jistá konstanta. Výsledky měření: 2 1,8 2,1 2,4 1,9 2,1 2 1,8 2,3 2,2. Pomocí Q-Q plotu ověřte, zda se tato data řídí normálním rozložením.

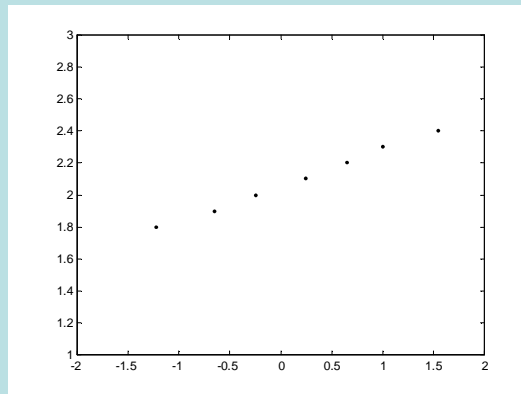
**Řešení:**

usp.hodnoty	1,8	1,8	1,9	2	2	2,1	2,1	2,2	2,3	2,4
pořadí	1	2	3	4	5	6	7	8	9	10
průměrné pořadí	1,5	1,5	3	4,5	4,5	6,5	6,5	8	9	10

Vektor hodnot průměrného pořadí:  $j = (1,5 \ 3 \ 4,5 \ 6,5 \ 8 \ 9 \ 10)$

vektor hodnot  $\alpha_j = \frac{j - 0,375}{n + 0,25} = (0,1098; 0,2561; 0,4024; 0,5976; 0,7439; 0,8415; 0,939)$

vektor kvantilů  $u_{\alpha_j} = (-1,2278; -0,6554; -0,247; 0,247; 0,6554; 1,0005; 1,566)$

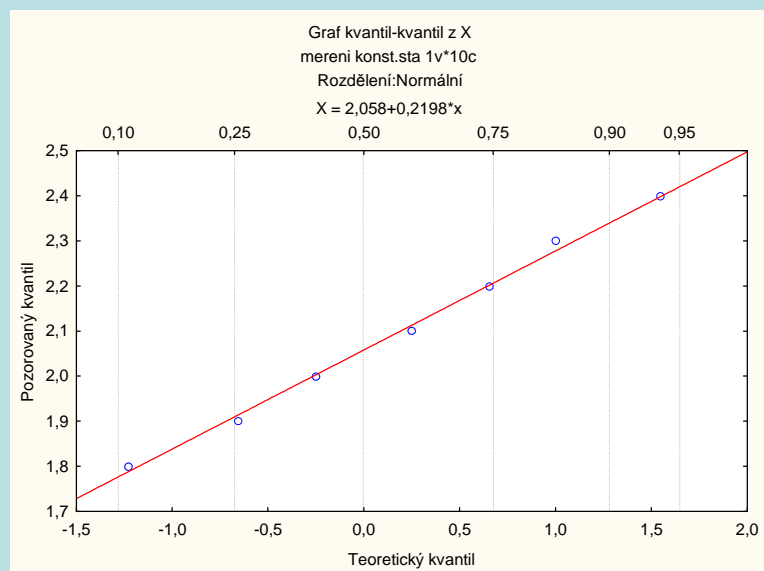


Vzhled grafu nasvědčuje tomu, že data pocházejí z normálního rozložení.

### Výpočet pomocí systému STATISTICA:

Otevřeme nový datový soubor o jedné proměnné a 10 případech. Zjištěné hodnoty zapíšeme do proměnné X.

Grafy – 2D Grafy – Grafy typu Q-Q– Proměnná X – OK - odškrtneme Neurčovat průměrnou pozici svázaných pozorování - OK.



## Probability - probability plot (P-P plot)

Používá se ke stejným účelům jako Q-Q plot, ale jinak se konstruuje.

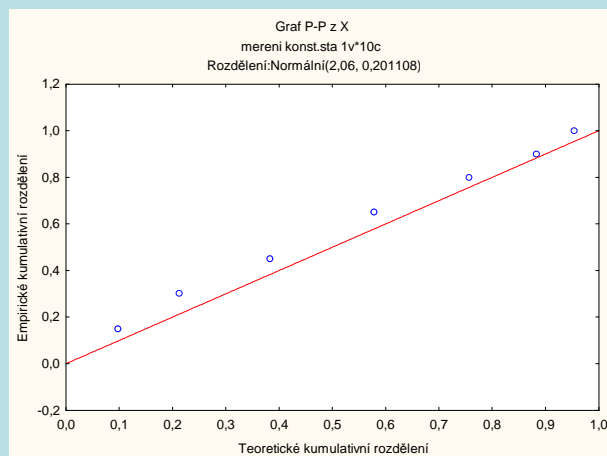
**Způsob konstrukce:** spočtou se standardizované hodnoty  $z_{(j)} = \frac{x_{(j)} - m}{s}$ ,  $j = 1, \dots, n$ . Na vodorovnou osu se vynesou hodnoty teoretické distribuční funkce  $\Phi(z_{(j)})$  a na svislou osu hodnoty empirické distribuční funkce  $F(z_{(j)}) = j/n$ . (Jsou-li některé hodnoty  $x_{(1)} \leq \dots \leq x_{(n)}$  stejné, pak za  $j$  bereme průměrné pořadí odpovídající takové skupince.) Pokud se body  $(\Phi(z_{(j)}), F(z_{(j)}))$  řadí kolem hlavní diagonály čtverce  $[0,1] \times [0,1]$ , lze usuzovat na dobrou shodu empirického a teoretického rozložení.

**Příklad na konstrukci P-P plotu pomocí systému STATISTICA:** Desetkrát nezávisle na sobě byla změřena jistá konstanta. Výsledky měření: 2 1,8 2,1 2,4 1,9 2,1 2 1,8 2,3 2,2. Pomocí P-P plotu ověřte, zda se tato data řídí normálním rozložením.

### Výpočet pomocí systému STATISTICA:

Otevřeme nový datový soubor o jedné proměnné a 10 případech. Zjištěné hodnoty zapíšeme do proměnné X.

Grafy – 2D Grafy – Grafy typu P-P – Proměnná X – OK - odškrtneme Neurčovat průměrnou pozici svázaných pozorování - OK.



## Histogram

Umožňuje porovnat tvar hustoty četnosti s tvarem hustoty pravděpodobnosti vybraného teoretického rozložení. (Ve STATISTICE je pojem histogramu širší, skrývá se za ním i sloupkový diagram.)

**Způsob konstrukce:** na vodorovnou osu vynášíme meze třídících intervalů. Nad každým třídícím intervalem sestrojíme obdélník o ploše odpovídající relativní četnosti příslušného třídícího intervalu, tj. výška obdélníku je rovna četnostní hustotě třídícího intervalu (četnostní hustota je relativní četnost třídícího intervalu dělená délkou tohoto intervalu).

**Způsob konstrukce ve STATISTICE:** na vodorovnou osu se vynášejí třídící intervaly (implicitně 10, jejich počet lze změnit, stejně tak i meze třídících intervalů) či varianty znaku a na svislou osu absolutní nebo relativní četnosti třídících intervalů či variant. Do histogramu se zakreslí tvar hustoty (či pravděpodobnostní funkce) vybraného teoretického rozložení.



### Příklad na konstrukci histogramu:

U 70 domácností byly zjišťovány týdenní výdaje na nealkoholické nápoje (v Kč).

Výdaje	(35,65)	(65,95)	(95,125)	(125,155)	(155,185)	(185,215)
Počet dom.	7	16	27	14	4	2

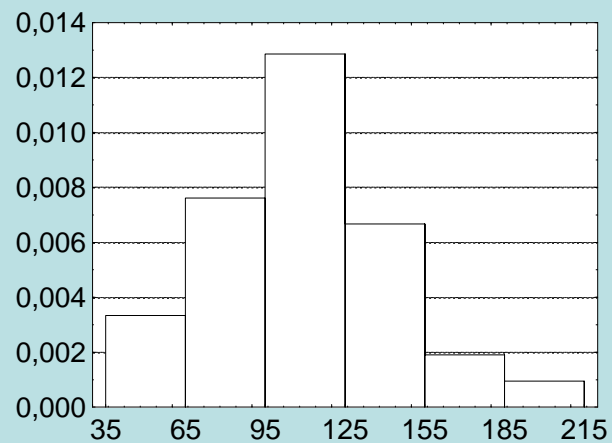
Nakreslete histogram.

### Řešení:

Nejprve sestavíme tabulku rozložení četností:

$(u_j, u_{j+1})$	$x_{[j]}$	$d_j$	$n_j$	$p_j$	$N_j$	$F_j$	$f_j$
(35,65)	50	30	7	$7/70=0,1$	7	$7/70=0,1$	$7/2100=0,0033$
(65,95)	80	30	16	$16/70=0,23$	23	$23/70=0,33$	$16/2100=0,0076$
(95,125)	110	30	27	$27/70=0,38$	50	$50/70=0,71$	$27/2100=0,0109$
(125,155)	140	30	14	$14/70=0,2$	64	$64/70=0,91$	$14/2100=0,0067$
(155,185)	170	30	4	$4/70=0,06$	68	$68/70=0,97$	$4/2100=0,0019$
(185,215)	200	30	2	$2/70=0,03$	70	$70/70=1$	$2/2100=0,0010$

S pomocí této tabulky sestojíme histogram:

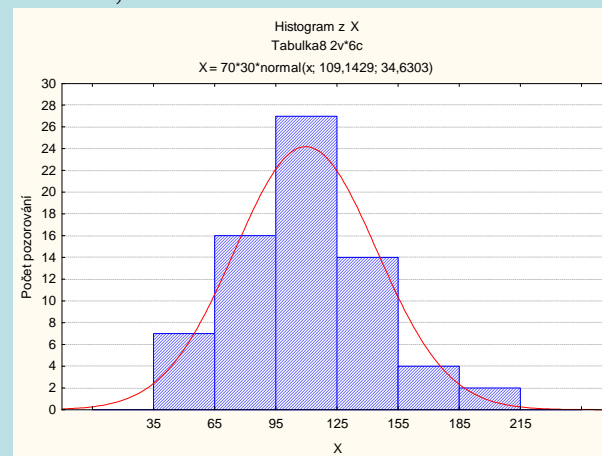


### Výpočet pomocí systému STATISTICA:

Otevřeme nový datový soubor o dvou proměnných a 6 případech. První proměnnou nazveme X, druhou četnost. Do proměnné X napíšeme středy třídících intervalů, do proměnné četnost odpovídající absolutní četnosti:

	1	2
	X	četnost
1	50	7
2	80	16
3	110	27
4	140	14
5	170	4
6	200	2

Grafy – Histogramy – zadáme proměnnou vah četnost – Proměnná X - zaškrtneme Hranice – Určit hranice – zaškrtneme Zadejte hraniční rozmezí: Minimum 35, Krok 30, Maximum 215 – OK – OK. Dostaneme graf:



Na rozdíl od histogramu konstruovaného ručně jsou na svislé ose absolutní četnosti, nikoliv četnostní hustoty. V porovnání s grafem hustoty normálního rozložení je vidět, že naše rozložení četností je lehce kladně zešikmené. Naše data tedy nepocházejí z normálního rozložení.

## Vzhled diagnostických grafů pro rozložení s různou šikmostí

Pro ilustraci se podívejme, jak se různá šikmost rozložení projeví na histogramu, N-P plotu a na krabicovém diagramu.

