

Intervalová data

Příklad

Výše pojistných plnění 227 klientů jedné pojistovny:

výše plnění	počet klientů
0 - 7500	99
7500 - 17500	42
17500 - 32500	29
32500 - 67500	28
67500 - 125000	17
125000 - 300000	9
přes 300000	3

jak reprezentovat taková data?

a) interval reprezentoval jedním bodem

b) reprezentoval si data z rovnoměrných rozdělení na přísl. intervaly

c) neparametrické odhady pro intervalová data

d) parametrické modely -||-

empirická distribuční funkce: $\hat{F}_m(x) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{X_i \leq x\} = \frac{1}{m} \{\text{počet pozorování} \leq x\}$.

omezené hranice intervalů $c_0 < c_1 < c_2 < \dots < c_k$ a počty pozorování v intervalu $(c_{j-1}, c_j]$ jako m_j .

(c_0 a c_k může být $\pm \infty$)

$(c_0, c_1] \dots m_1$
 $(c_1, c_2] \dots m_2$
 \vdots
 $(c_{k-1}, c_k] \dots m_k$
 $\frac{\quad}{m}$

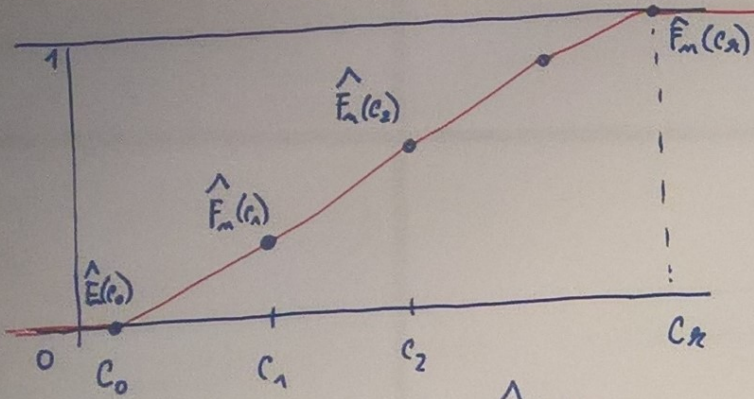
Podle definice můžeme určit hodnoty emp. dist. fun. v bodech c_0, c_1, \dots, c_k :

$$\hat{F}_m(c_0) = 0$$

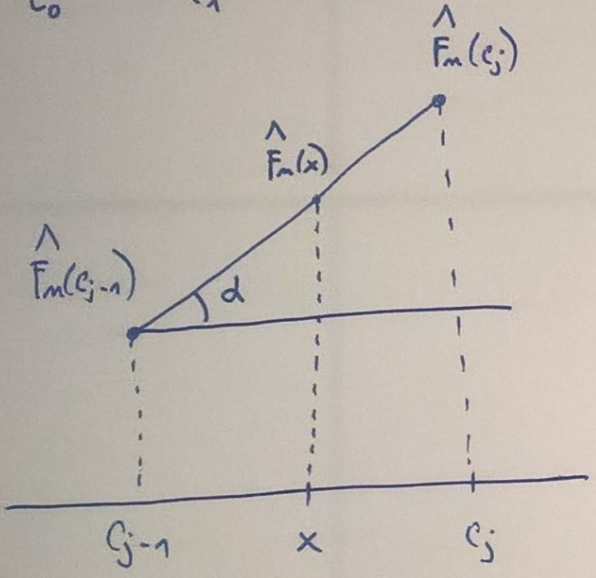
$$\hat{F}_m(c_j) = \frac{m_1 + m_2 + \dots + m_j}{m}$$

a mezi nimi spojitě dodefinovat (ogive)

$$\hat{F}_m(c_k) = 1$$



$$\hat{F}_m(x) = \begin{cases} 0 & x \leq c_0 \\ \frac{c_j - x}{c_j - c_{j-1}} \cdot \hat{F}_m(c_{j-1}) + \frac{x - c_{j-1}}{c_j - c_{j-1}} \cdot \hat{F}_m(c_j) & c_{j-1} \leq x \leq c_j \\ 1 & x \geq c_k \end{cases}$$



$$d = \frac{\hat{F}_m(c_j) - \hat{F}_m(c_{j-1})}{c_j - c_{j-1}}$$

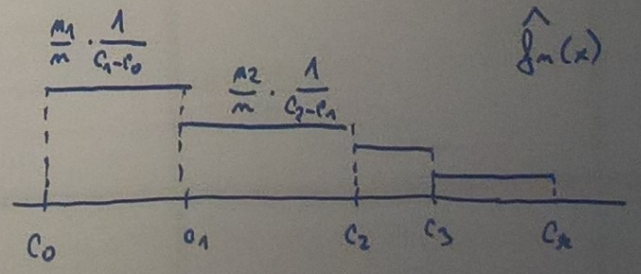
$$\hat{F}_m(x) - \hat{F}_m(c_{j-1}) = (x - c_{j-1}) \cdot d$$

empirická hustota = histogram

$$\hat{f}_m(x) = \hat{f}_m^*(x) \quad \text{pro } x \neq c_0, c_1, \dots, c_k$$

$$\hat{f}_m(x) = \frac{-\hat{F}_m(c_{j-1}) + \hat{F}_m(c_j)}{c_j - c_{j-1}} = \frac{\frac{M_1 + \dots + M_{j-1}}{m} + \frac{M_1 + \dots + M_j}{m}}{c_j - c_{j-1}} = \frac{M_j}{m} \cdot \frac{1}{c_j - c_{j-1}}$$

pro $c_{j-1} \leq x < c_j$.



empirická 'kumulovaná' funkce:

učíme ji jako součet funkce k emp. data. f_i (ogive)

$$\hat{Q}_m(\alpha) = \hat{F}_m^{-1}(\alpha), \quad 0 < \alpha < 1$$

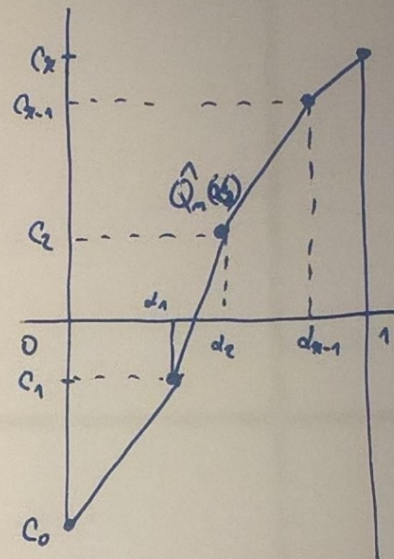
průběh: $\hat{Q}_m(0) = c_0$

$$\hat{Q}_m(\alpha) = c_j$$

pro $d = \frac{m_1 + m_2 + \dots + m_j}{n} =: d_j$

$$\hat{Q}_m(1) = c_k$$

} a mezi nimi spojitě dodefinujeme



$$\hat{Q}_m(\alpha) = \frac{d_j - d}{d_j - d_{j-1}} \cdot \hat{Q}_m(d_{j-1}) + \frac{d - d_{j-1}}{d_j - d_{j-1}} \cdot \hat{Q}_m(d_j), \quad d_{j-1} \leq d \leq d_j$$

Parametrické modely

původní data: X_1, \dots, X_m máh. ojetí X_i má distribuční funkci $F(x, \theta)$, resp. hustotu $f(x, \theta)$, kde θ_j neznámý parametr.

metoda maximální věrohodnosti

Počítáme pravd. že dané pozorování padne do intervalu $(c_{j-1}, c_j]$: $P(X_i \in (c_{j-1}, c_j]) = F(c_j, \theta) - F(c_{j-1}, \theta) = \int_{c_{j-1}}^{c_j} f(x, \theta) dx$

$$\begin{aligned} \text{věrohodnostní funkce } L(\theta) &= \prod_{i=1}^m P(X_i \in (c_{j-1}, c_j]) = [F(c_1, \theta) - F(c_0, \theta)]^{m_1} \cdot [F(c_2, \theta) - F(c_1, \theta)]^{m_2} \cdot \dots \cdot [F(c_k, \theta) - F(c_{k-1}, \theta)]^{m_k} = \\ &= \prod_{j=1}^k [F(c_j, \theta) - F(c_{j-1}, \theta)]^{m_j} \end{aligned}$$

logaritmická rešhodnotní funkce $l(\theta) = \log L(\theta) = \sum_{j=1}^k m_j \cdot \log [F(c_j, \theta) - F(c_{j-1}, \theta)]$

$\hat{\theta}$ je odhad parametru θ metodou max. věr. ještě $\hat{\theta} = \operatorname{argmax} L(\theta) = \operatorname{argmax} l(\theta)$.

Příklad

mezi se výše popsaného plnění X_i řidi exponenciálním rozdělením s parametrem $\frac{1}{\lambda}$, $\lambda > 0$

$f(x, \lambda) = \frac{1}{\lambda} \cdot e^{-\frac{x}{\lambda}}$, $x > 0$

$F(x, \lambda) = 1 - e^{-\frac{x}{\lambda}}$, $x > 0$

$F(c_j, \lambda) - F(c_{j-1}, \lambda) = e^{-\frac{c_{j-1}}{\lambda}} - e^{-\frac{c_j}{\lambda}}$

$l(\lambda) = 99 \cdot \log(1 - e^{-\frac{7500}{\lambda}}) + 42 \cdot \log(e^{-\frac{7500}{\lambda}} - e^{-\frac{17500}{\lambda}}) + \dots + 3 \cdot \log(e^{-\frac{30000}{\lambda}}) \dots$ numerická maximalizace $\rightarrow \hat{\lambda}$ = odhad průměrné výše plnění

III. Metoda minimálního χ^2

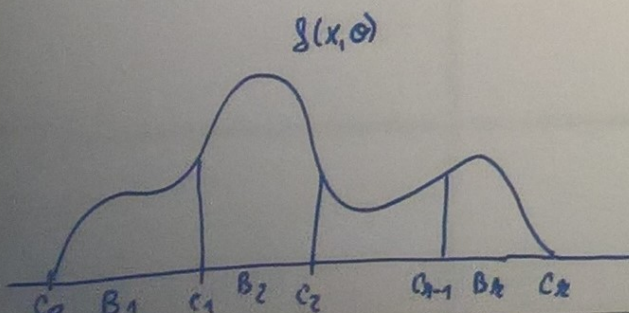
mějme opět intervalová data

$(c_0, c_1] = B_1 \dots m_1$ pozorování

$(c_1, c_2] = B_2 \dots m_2$ pozorování

$(c_{k-1}, c_k] = B_k \dots m_k$ pozorování

celkem $m = m_1 + m_2 + \dots + m_k$ pozorování



stejně jako v předchozím případě: $P(X_i \in B_j) = P(X_i \in (c_{j-1}, c_j]) = \int_{c_{j-1}}^{c_j} f(x, \theta) dx = F(c_j, \theta) - F(c_{j-1}, \theta) =: p_j(\theta)$

celkem m pozorování, tedy v intervalu $B_j = (c_{j-1}, c_j]$ by mělo být $m \cdot p_j(\theta)$ pozorování (očekávaná četnost), m_j pozorování \rightarrow pozorovanými m_j :

$$\chi^2(\theta) = \sum_{j=1}^k \frac{[m_j - m \cdot p_j(\theta)]^2}{m \cdot p_j(\theta)}$$

$\tilde{\theta}$ odhad parametru θ metodou minimálního χ^2 , jestliže $\tilde{\theta} = \underset{\theta}{\operatorname{argmin}} \chi^2(\theta)$.

Prozámky

$$\chi^2(\theta) = \sum_{j=1}^k \frac{m_j^2 - 2m m_j p_j(\theta) + m^2 p_j^2(\theta)}{m p_j(\theta)} = \sum_{j=1}^k \frac{m_j^2}{m p_j(\theta)} - 2 \cdot \underbrace{\sum_{j=1}^k m_j}_{=m} + m \underbrace{\sum_{j=1}^k p_j(\theta)}_{=1} = \sum_{j=1}^k \frac{m_j^2}{m p_j(\theta)} - m.$$

g-li θ p -rozměrný parametr a $p_j(\theta)$ jsou diferencovatelné podle θ , pak pro $\tilde{\theta}$ platí:

$$-\sum_{j=1}^k \frac{m_j^2}{m p_j^2(\theta)} \cdot \frac{\partial p_j(\theta)}{\partial \theta_h} = 0 \quad \text{pro } h=1, 2, \dots, p. \quad \tilde{\theta} \text{ není soustavou: } \sum_{j=1}^k \frac{m_j^2}{p_j^2(\theta)} \cdot \frac{\partial p_j(\theta)}{\partial \theta_h} = 0 \quad \text{pro } h=1, 2, \dots, p.$$

Prozámky
Metoda není použitelná na "klasická" data, když X_1, \dots, X_m tvoří náh. výběr z rozdělení $f(x, \theta)$. Intervaly B_1, \dots, B_k si vyhovíme sami a dopočítáme m_1, \dots, m_k .

Jak volit h ? heuristická pravidla: $h \approx 2n^{\frac{2}{5}}$, nebo $h \approx 15 \cdot \left(\frac{n}{100}\right)^{\frac{2}{5}}$.

Intervaly B_1, \dots, B_k se volí tak, aby všechny byly stejně pravděpodobné "lj". $\pi_1(\tilde{\Theta}) = \pi_2(\tilde{\Theta}) = \dots = \pi_k(\tilde{\Theta}) = \frac{1}{k}$.

$\tilde{\Theta}$ ale neznáme!

V praxi vezmeme nějaký jiný odhad Θ^* parametru Θ a požadujeme, aby $\pi_1(\Theta^*) = \pi_2(\Theta^*) = \dots = \pi_k(\Theta^*) = \frac{1}{k}$.

Jinými slovy c_1, c_2, \dots, c_{k-1} jsou $\frac{1}{k}, \frac{2}{k}, \dots, \frac{k-1}{k}$ - kvantily rozdělení s hustotou $f(x, \Theta^*)$.