

Neparametrické modelování hustot a distribucí.

K danému náhodnému výběru $\{x_1, x_2, \dots, x_n\}$ pevného rozsahu n hledáme neznámou hustotu

ve tvaru: $f_a(x) = \frac{1}{nh} \sum_{i=1}^n \kappa\left(\frac{x-x_i}{h}\right)$ kde:

$h; h > 0$ je parametr měřítka – pokud na něčem závislý, pak pouze nanejvýš na rozsahu výběru n (někdy též i na pozorovaných hodnotách), vyhlazovací parametr, scaling parameter, ... a

$\kappa(x)$ je jádro aproximace – nějaká hustota, splňující následující podmínky:

1. $\kappa(x) \geq 0; \forall x \in R_1; \int_{-\infty}^{+\infty} \kappa(x) dx = 1$, triviální podmínky toho, že za jádro budeme považovat jen hustoty (v obecné teorii tomu tak nemusí být, často je opouštěna podmínka nezápornosti).
2. $\int_{-\infty}^{+\infty} x \kappa(x) dx = 0$; jádro je centrované, jaderná hustota má nulovou střední hodnotu
3. $\int_{-\infty}^{+\infty} x^2 \kappa(x) dx = 1$; jádro je normované, jaderná hustota má jednotkový rozptyl (samozřejmě za platnosti podmínky 2.).

Základní vlastnosti odhadu $f_a(x)$

$$\int_{-\infty}^{+\infty} f_a(x) dx = \frac{1}{nh} \sum_{i=1}^n \int_{-\infty}^{+\infty} \kappa\left(\frac{x-x_i}{h}\right) dx = (A) \text{ po substituci } \frac{x-x_i}{h} = y \Rightarrow x = hy + x_i \Rightarrow dx = hdy$$

dostaneme $(A) = \frac{1}{nh} \sum_{i=1}^n h \int_{-\infty}^{+\infty} \kappa(y) dy = \frac{1}{nh} \sum_{i=1}^n h = 1$. Tedy zvolený typ aproximace je opět hustotou.

$$\mu_a = E_{f_a}\{x\} = \int_{-\infty}^{+\infty} x f_a(x) dx = \frac{1}{nh} \sum_{i=1}^n \int_{-\infty}^{+\infty} x \kappa\left(\frac{x-x_i}{h}\right) dx = (B) \text{ po substituci}$$

$$\frac{x-x_i}{h} = y \Rightarrow x = hy + x_i \Rightarrow dx = hdy; \quad (B) = \frac{1}{nh} \sum_{i=1}^n h \int_{-\infty}^{+\infty} (hy + x_i) \kappa(y) dy =$$

$$= \frac{1}{n} \sum_{i=1}^n h \int_{-\infty}^{+\infty} y \kappa(y) dy + \frac{1}{n} \sum_{i=1}^n x_i \int_{-\infty}^{+\infty} \kappa(y) dy = \frac{1}{n} \sum_{i=1}^n x_i$$

=0 podle
podmínky 2.

= 1 podle
podmínky 1.

Tedy střední hodnota počítaná podle aproximační hustoty je rovna výběrovému průměru.

$$\sigma_a^2\{x\} = E \int_{-\infty}^{+\infty} (x - \mu_a)^2 f_a(x) dx = \frac{1}{nh} \sum_{i=1}^n \int_{-\infty}^{+\infty} (x - \mu_a)^2 \kappa\left(\frac{x-x_i}{h}\right) dx = (C) \text{ po substituci}$$

$$\frac{x-x_i}{h} = y \Rightarrow x = hy + x_i \Rightarrow dx = hdy; \quad (C) = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{+\infty} (hy + x_i - \mu_a)^2 \kappa(y) dy =$$

$$= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{+\infty} \left[(hy)^2 + (x_i - \mu_a)^2 + 2hy(x_i - \mu_a) \right] \kappa(y) dy = \frac{1}{n} \sum_{i=1}^n \left[\int_{-\infty}^{+\infty} (hy)^2 \kappa(y) dy + \int_{-\infty}^{+\infty} (x_i - \mu_a)^2 \kappa(y) dy \right] +$$

$$+ \frac{1}{n} \sum_{i=1}^n \left[\int_{-\infty}^{+\infty} 2hy(x_i - \mu_a) \kappa(y) dy \right] = h^2 + \frac{1}{n} \sum_{i=1}^n (x_i - \mu_a)^2$$

= h² podle podmínky 3.

= 0 podle podmínky 2.

= (x_i - μ_a)² podle podmínky 1.

Pokud označíme $s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu_a)^2$, výběrový rozptyl, dostaneme:

$$\sigma_a^2\{x\} = h^2 + \frac{1}{n} \sum_{i=1}^n (x_i - \mu_a)^2 = h^2 + \left(1 - \frac{1}{n}\right) s_n^2. \text{ Pokud bychom požadovali } \sigma_a^2\{x\} = s_n^2,$$

$$\text{dostaneme volbu } h = \sqrt{\frac{s_n^2}{n}}.$$

(1)

„Nestrannost“ odhadu $f_a(x)$

Podmínky náhodného výběru nám dávají shodnou hustotu $f(x)$, nám nedostupnou, pro všechna pozorování tohoto náhodného výběru. Dále podle podmínek kladených na náhodný výběr (iid) platí pro sdruženou hustotu $f(x_1, \dots, x_n) = \prod_{i=1}^n f(x_i)$. Pak:

$$\begin{aligned} E_{x_1, x_2, \dots, x_n} \{f_a(x)\} &= \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} f_a(x) \left[\prod_{i=1}^n f(x_i) \right] dx_1 \dots dx_n = \\ &= \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} \frac{1}{nh} \sum_{j=1}^n \kappa\left(\frac{x-x_j}{h}\right) \left[\prod_{i=1}^n f(x_i) \right] dx_1 \dots dx_n = \frac{1}{nh} \sum_{j=1}^n \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} \kappa\left(\frac{x-x_j}{h}\right) \left[\prod_{i=1}^n f(x_i) \right] dx_1 \dots dx_n = \\ &= \frac{1}{nh} \sum_{j=1}^n \int_{-\infty}^{+\infty} \kappa\left(\frac{x-x_j}{h}\right) f(x_j) dx_j = \frac{1}{h} \int_{-\infty}^{+\infty} \kappa\left(\frac{x-z}{h}\right) f(z) dz = (D) \end{aligned}$$

$$\text{po substituci: } \frac{x-z}{h} = y; z = x - hy; dz = -hdy \text{ dostaneme } (D) = \int_{-\infty}^{+\infty} \kappa(y) f(x - hy) dy.$$

$$\text{Shrnuto: } E_{x_1, x_2, \dots, x_n} \{f_a(x)\} = \int_{-\infty}^{+\infty} \kappa(y) f(x - hy) dy.$$

Potom $\lim_{h \rightarrow 0_+} E_{x_1, x_2, \dots, x_n} \{f_a(x)\} = \lim_{h \rightarrow 0_+} \int_{-\infty}^{+\infty} \kappa(y) f(x - hy) dy = (E)$ a pokud jsou splněny předpoklady záměny limity a integrálu (zde v nevlastních mezích) a funkce $f(x)$ je spojitá, platí:

$$(E) = \int_{-\infty}^{+\infty} \kappa(y) \left(\lim_{h \rightarrow 0_+} f(x - hy) \right) dy = \int_{-\infty}^{+\infty} \kappa(y) \left(\lim_{h \rightarrow 0_+} f(x - \lim_{h \rightarrow 0_+}(hy)) \right) dy = f(x).$$

Důkaz tohoto tvrzení za daleko obecnějších předpokladů je v: L. Devroye and L. Györfi: Nonparametric Density Estimation: The L1 View, str. 6-11. Pokud tedy je $h = h(n) \xrightarrow{n \rightarrow +\infty} 0$ je odhad hustoty $f_a(x)$ asymptoticky nestranný.

V L. Devroye and L. Györfi: Nonparametric Density Estimation: The L1 View, str. 12-19 je

dokázáno pro $J_n = \int_{-\infty}^{+\infty} |f_a(x) - f(x)| dx$: Následující podmínky jsou ekvivalentní:

1. $J_n \xrightarrow{n \rightarrow +\infty} 0$ v pravděpodobnosti pro některou hustotu $f(x)$.
2. $J_n \xrightarrow{n \rightarrow +\infty} 0$ v pravděpodobnosti pro každou hustotu $f(x)$.
3. $J_n \xrightarrow{n \rightarrow +\infty} 0$ skoro všude pro každou hustotu $f(x)$.
4. $J_n \xrightarrow{n \rightarrow +\infty} 0$ „exponenciálně“¹ pro každou hustotu $f(x)$. $J_n \xrightarrow{n \rightarrow +\infty} 0$ „exponenciálně“ znamená: $\forall \varepsilon > 0; \exists r > 0, n_0$ takové, že $n \geq n_0 \Rightarrow P(J_n \geq \varepsilon) \leq e^{-r^n}$. To ale také znamená konvergenci řady $\sum_{n=1}^{+\infty} P(J_n \geq \varepsilon)$ neboť má konvergentní majorantu e^{-r^n} .

Tam uvedený důkaz potřebuje jen předpoklad 1. na jádro $\kappa(x)$.

Diskuse „nestrannosti“ jádrového odhadu hustoty

Platí: $E_{x_1, x_2, \dots, x_n} \{f_a(x)\} = \int_{-\infty}^{+\infty} \kappa(y) f(x - hy) dy$, odtud:

$E_{x_1, x_2, \dots, x_n} \{f_a(x)\}$ je hustotou

- Nezápornost je zřejmá
- $\int_{-\infty}^{+\infty} E_{x_1, x_2, \dots, x_n} \{f_a(x)\} dx = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \kappa(y) f(x - hy) dy dx = 1$

Dále, použijeme-li ve vztahu $E_{x_1, x_2, \dots, x_n} \{f_a(x)\} = \int_{-\infty}^{+\infty} \kappa(y) f(x - hy) dy$ substituci $z = hy \Rightarrow dy = \frac{dz}{h}$,

dostaneme: $E_{x_1, x_2, \dots, x_n} \{f_a(x)\} = \int_{-\infty}^{+\infty} \kappa(y) f(x - hy) dy = \int_{-\infty}^{+\infty} \frac{1}{h} \kappa\left(\frac{z}{h}\right) f(x - z) dz$.

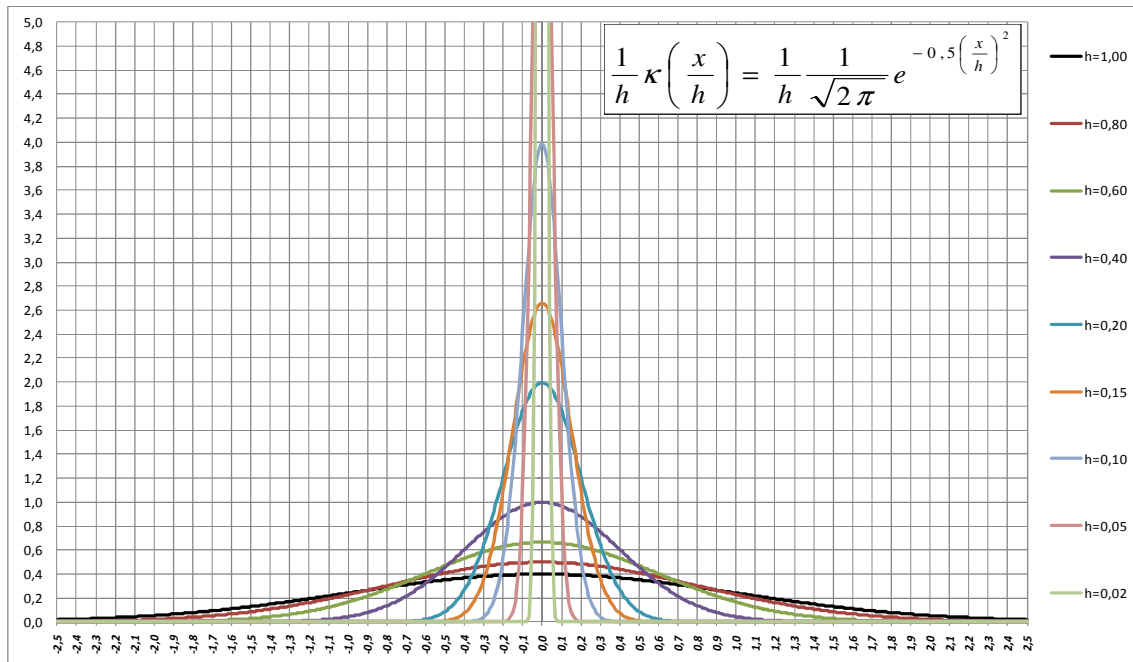
¹ $J_n \xrightarrow{n \rightarrow +\infty} 0$ „exponenciálně“ znamená: $\forall \varepsilon > 0; \exists r > 0, n_0$ takové, že $n \geq n_0 \Rightarrow P(J_n \geq \varepsilon) \leq e^{-r^n}$. To ale

také znamená, že řada $\sum_{n=1}^{+\infty} P(J_n \geq \varepsilon)$ je konvergentní neboť má konvergentní majorantu e^{-r^n} .

To je ale konvoluce dvou hustot $\frac{1}{h} \kappa\left(\frac{x}{h}\right)$ a $f(x)$.

Proto: hustota $E_{x_1, x_2, \dots, x_n} \{f_a(x)\}$ je hustotou náhodné proměnné, která je tvořena součtem původní (pozorované) náhodné proměnné a „na ní“ nezávislého šumu ε_h s hustotou $\frac{1}{h} \kappa\left(\frac{x}{h}\right)$.

Máme-li „šum“ ε_1 s hustotou $\kappa(x)$, šum s hustotou $\frac{1}{h} \kappa\left(\frac{x}{h}\right)$ bude $\varepsilon_h = h\varepsilon_1$. Příklad hustot takového „šumu“ je na následujícím obrázku:



„Vydatnost“ odhadu $f_a(x) = \frac{1}{nh} \sum_{i=1}^n \kappa\left(\frac{x-x_i}{h}\right)$

$$E_{x_1, x_2, \dots, x_n} \{f_a^2(x)\} = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} f_a^2(x) \left[\prod_{i=1}^n f(x_i) \right] dx_1 \dots dx_n =$$

$$= \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} \left[\frac{1}{nh} \sum_{j=1}^n \kappa\left(\frac{x-x_j}{h}\right) \right]^2 \left[\prod_{i=1}^n f(x_i) \right] dx_1 \dots dx_n =$$

$$= \frac{1}{n^2 h^2} \sum_{j=1}^n \sum_{k=1}^n \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} \kappa\left(\frac{x-x_j}{h}\right) \kappa\left(\frac{x-x_k}{h}\right) \left[\prod_{i=1}^n f(x_i) \right] dx_1 \dots dx_n =$$

$$\begin{aligned}
&= \frac{1}{n^2 h^2} \left[\sum_{i=1}^n \int_{-\infty}^{+\infty} \kappa^2 \left(\frac{x-x_k}{h} \right) f(x_k) dx_k + 2 \sum_{i=1}^n \sum_{j=i+1}^n \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \kappa \left(\frac{x-x_k}{h} \right) \kappa \left(\frac{x-x_j}{h} \right) f(x_k) f(x_j) dx_k dx_j \right] = \\
&= \frac{1}{nh^2} \int_{-\infty}^{+\infty} \kappa^2 \left(\frac{x-x_k}{h} \right) f(x_k) dx_k + \frac{n-1}{nh^2} \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \kappa \left(\frac{x-x_k}{h} \right) \kappa \left(\frac{x-x_j}{h} \right) f(x_k) f(x_j) dx_k dx_j = \\
&= \frac{1}{nh} \int_{-\infty}^{+\infty} \kappa^2(y) f(x-hy) dy + \frac{n-1}{n} \left(\int_{-\infty}^{+\infty} \kappa(y) f(x-hy) dy \right)^2.
\end{aligned}$$

2

Protože: $\sigma_{x_1, x_2, \dots, x_n}^2 \{f_a(x)\} = E_{x_1, x_2, \dots, x_n} \{f_a^2(x)\} - (E_{x_1, x_2, \dots, x_n} \{f_a(x)\})^2$ a
 $E_{x_1, x_2, \dots, x_n} \{f_a(x)\} = \int_{-\infty}^{+\infty} \kappa(y) f(x-hy) dy$, dostaneme pro rozptyl odhadu:

$\sigma_{x_1, x_2, \dots, x_n}^2 \{f_a(x)\} = \frac{1}{nh} \int_{-\infty}^{+\infty} \kappa^2(y) f(x-hy) dy - \frac{1}{n} \left(\int_{-\infty}^{+\infty} \kappa(y) f(x-hy) dy \right)^2$. Odtud bezprostředně (a
z asymptotické nestrannosti) plyne potřeba $n \rightarrow +\infty \Rightarrow [(h \rightarrow 0) \wedge (nh \rightarrow +\infty)]$ pro existenci
rozptylu a asymptotickou vydatnost (za předpokladu $h = h(n)$).

Shrnuto: $\sigma_{x_1, x_2, \dots, x_n}^2 \{f_a(x)\} = \frac{1}{nh} \int_{-\infty}^{+\infty} \kappa^2(y) f(x-hy) dy - \frac{1}{n} \left(\int_{-\infty}^{+\infty} \kappa(y) f(x-hy) dy \right)^2$.

Proto, aby byl takový odhad asymptoticky nestranný a asymptoticky vydatný, stačí:

$$\begin{aligned}
n &\rightarrow +\infty \Rightarrow (h \rightarrow 0) \\
n &\rightarrow +\infty \Rightarrow (nh \rightarrow +\infty)
\end{aligned}$$

Příklady jádrových funkcí

„Obdélníkové jádro“:

$$\kappa(x) = \frac{1}{2a} \Leftrightarrow -a \leq x \leq +a; \quad a > 0; \text{ jinak } \kappa(x) = 0 \quad (\mathbf{K1})$$

1. Podmínka je zřejmá
2. Podmínka je též zřejmá, hustota je symetrická.

$$^2 \sum_{j=n+1}^n \dots = 0$$

$$3. \text{ Podmínka } \int_{-a}^{+a} x^2 \frac{1}{2a} dx = \frac{1}{2a} \int_{-a}^{+a} x^2 dx = \frac{1}{2a} \left[\frac{x^3}{3} \right]_{-a}^{+a} = \frac{1}{2a} \left[\frac{a^3}{3} + \frac{a^3}{3} \right] = \frac{a^2}{3} = 1 \Rightarrow a = \sqrt{3}$$

„Epanechnikovo jádro“:

$$\kappa(x) = b \left(1 - \left(\frac{x}{a} \right)^2 \right) \Leftrightarrow -a \leq x \leq +a; \quad a, b > 0; \text{ jinak } \kappa(x) = 0 \quad (\mathbf{K2})$$

$$1. \text{ Podmínka: nezápornost je zřejmá, } \int_{-a}^{+a} b \left(1 - \left(\frac{x}{a} \right)^2 \right) dx = b \int_{-a}^{+a} \left(1 - \left(\frac{x}{a} \right)^2 \right) dx =$$

$$= b \left[\left(1 - \frac{x^3}{3a^2} \right) \right]_{-a}^{+a} = \frac{4}{3} ab = 1 \Rightarrow b = \frac{3}{4a}.$$

2. Podmínka: je zřejmá, hustota je symetrická.

3. Podmínka:

$$\int_{-a}^{+a} bx^2 \left(1 - \left(\frac{x}{a} \right)^2 \right) dx = b \int_{-a}^{+a} \left(x^2 - \frac{x^4}{a^2} \right) dx = b \left[\frac{x^3}{3} - \frac{x^5}{5a^2} \right]_{-a}^{+a} = b \left(\frac{2a^3}{3} - \frac{2a^5}{5a^2} \right) = \frac{4a^3b}{15} = 1$$

$$\text{Řešení obou podmínek } b = \frac{3}{4a} \text{ a } \frac{4a^3b}{15} = 1 \text{ dá } a = \sqrt{5}; \quad b = \frac{3}{4\sqrt{5}}.$$

„Trojúhelníkové jádro“:

$$\kappa(x) = b \left(1 - \left| \frac{x}{a} \right| \right) \Leftrightarrow -a \leq x \leq +a; \quad a, b > 0; \text{ jinak } \kappa(x) = 0 \quad (\mathbf{K3})$$

1. Podmínka: nezápornost je zřejmá,

$$2 \int_0^{+a} b \left(1 - \left(\frac{x}{a} \right) \right) dx = 2b \left[x - \frac{x^2}{2a} \right]_0^{+a} = 2b \left(a - \frac{a^2}{2a} \right) = ab = 1$$

2. Podmínka: je zřejmá, hustota je symetrická

$$3. \text{ Podmínka: } 2 \int_0^{+a} bx^2 \left(1 - \left(\frac{x}{a} \right) \right) dx = 2b \left[\frac{x^3}{3} - \frac{x^4}{4a} \right]_0^{+a} = 2b \left(\frac{a^3}{3} - \frac{a^4}{4a} \right) = \frac{a^3b}{6} = 1$$

$$\text{Řešení obou podmínek } ab = 1 \text{ a } \frac{a^3b}{6} = 1 \text{ dá } a = \sqrt[3]{6}; \quad b = \frac{1}{\sqrt[3]{6}}.$$

„Mocninné jádro“:

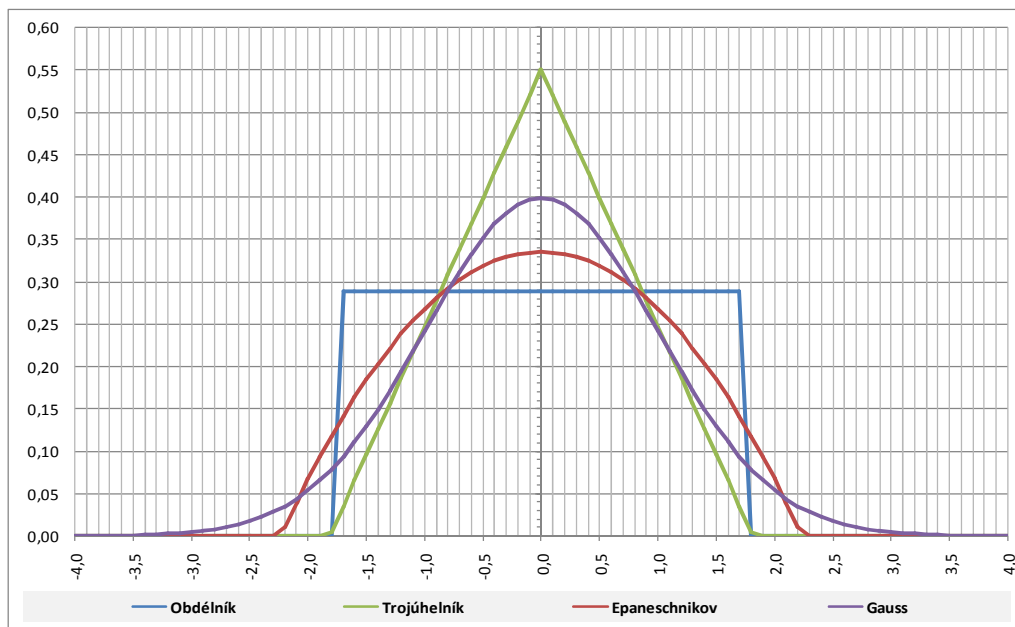
$$\kappa(x) = b \left(1 - \left| \frac{x}{a} \right|^\alpha \right) \Leftrightarrow -a \leq x \leq +a; \quad \alpha > 0; \quad a, b > 0; \text{ jinak } \kappa(x) = 0 \quad (\mathbf{K4})$$

Podmínky a hodnoty konstant a, b se odvodí obdobným způsobem jako u „trojúhelníkového jádra“

„Gaussovo jádro“:

$$\kappa(x) = \frac{1}{\sqrt{2\pi}} e^{-0,5x^2} \quad (\text{K5})$$

1. Podmínka je zřejmá, jde o normální, centrovanou a normovanou hustotu.
2. Podmínka je opět zřejmá, hustota je symetrická kolem nuly.
3. Podmínka je také zřejmá, hustota je centrovaná a normovaná.



Odpovídající modely distribučních funkcí

K danému tvaru jádrové funkce přiřadíme: $K(x) = \int_{-\infty}^x \kappa(z) dz$. Pak pro aproximaci distribuční

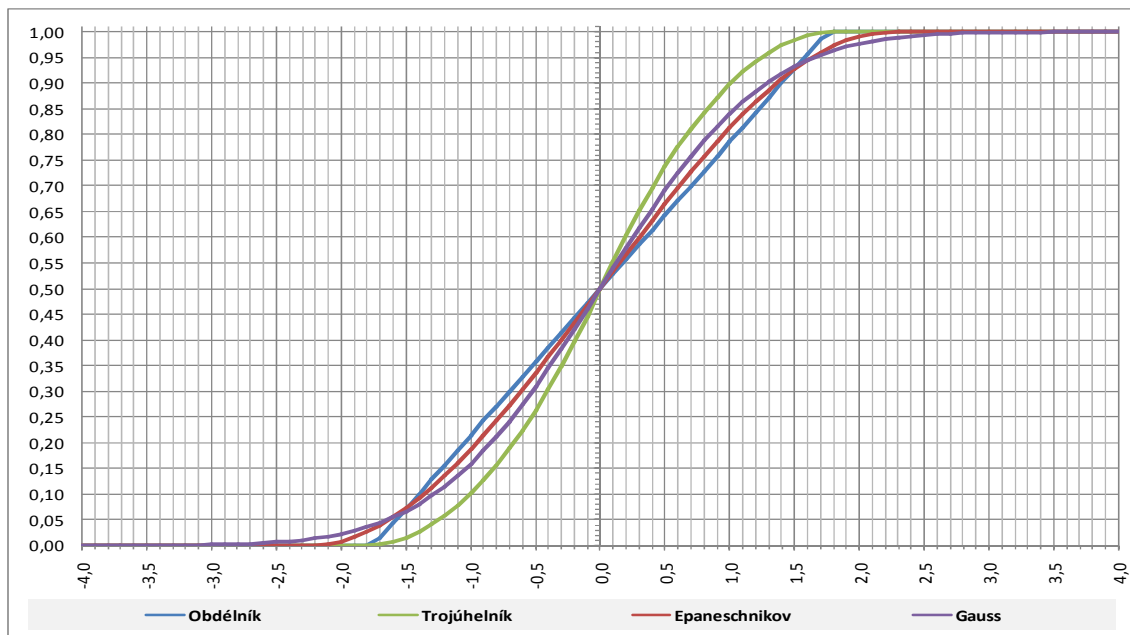
funkce dostaneme: $F_a(x) = \frac{1}{n} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$.

Např.: pro „obdélníkové jádro“: $\kappa(x) = \frac{1}{2a} \Leftrightarrow -a \leq x \leq +a$; $a > 0$; jinak $\kappa(x) = 0$ je odpovídající „distribuční jádro“:

$$K(x) = \frac{1}{a} x + 1 \Leftrightarrow -a \leq x \leq a; \quad K(x) = 0 \Leftrightarrow x < -a; \quad K(x) = 1 \Leftrightarrow x > a$$

Po dosazení $a = \sqrt{3}$ máme:

$$K(x) = \frac{1}{\sqrt{2}} x + 1 \Leftrightarrow -\sqrt{2} \leq x \leq \sqrt{2}; \quad K(x) = 0 \Leftrightarrow x < -\sqrt{2}; \quad K(x) = 1 \Leftrightarrow x > \sqrt{2}$$



Volba parametru h

Jedna možnost byla uvedena výše: $h = \sqrt{\frac{s_n^2}{n}}$ z podmínky aby neparametrická hustota měla stejný rozptyl jako výběrový. Další přístupy vycházejí z některých kritérií:

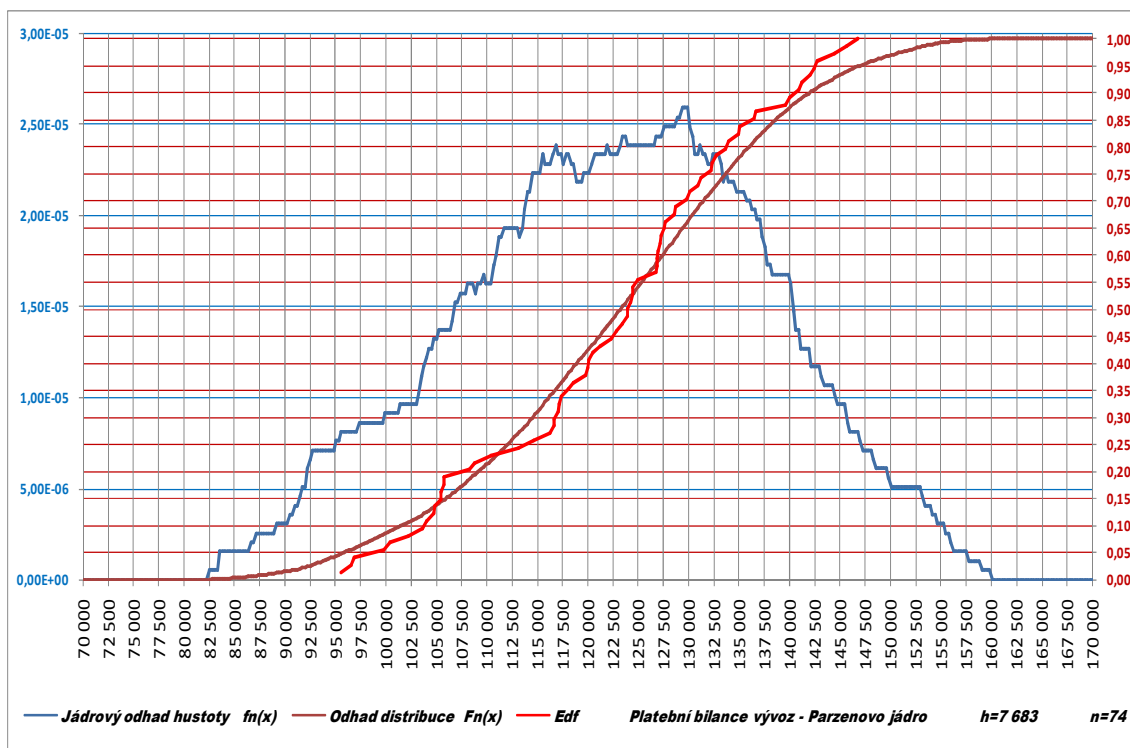
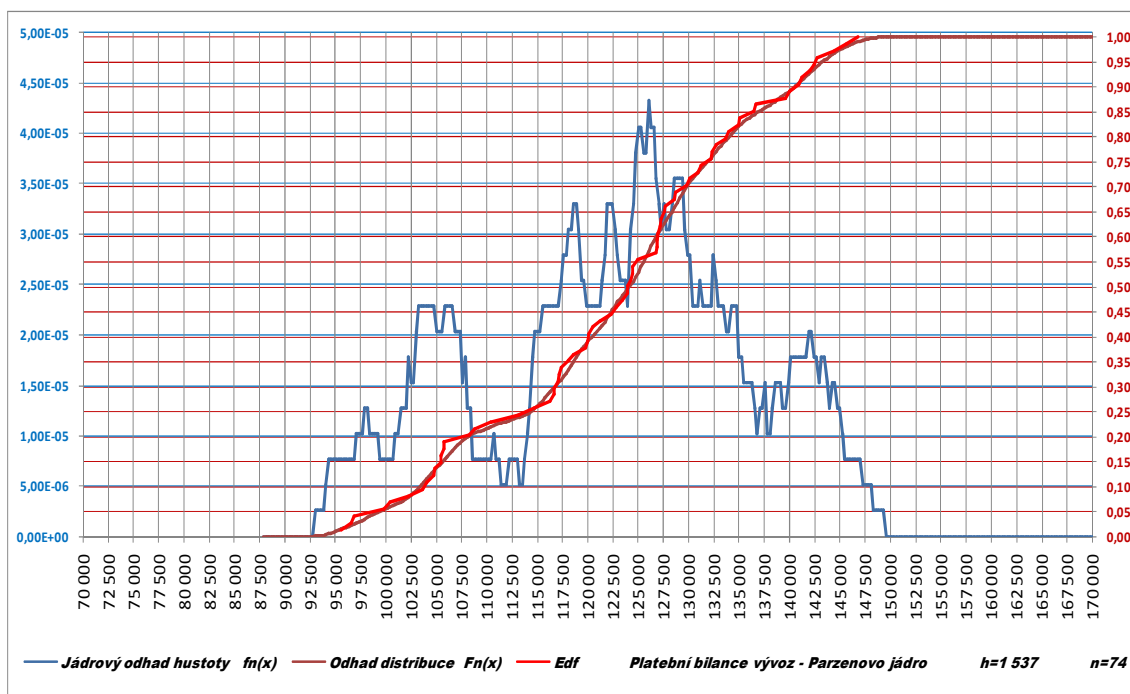
$$E\left\{f_{a,h,n}(x) - f(x)\right\} \xrightarrow{h} \min \quad \forall x; \quad E\left\{\sup_x |f_{a,h,n}(x) - f(x)|\right\} \xrightarrow{h} \min(\inf);$$

$$\sup_f E\left\{f_{a,h,n}(x) - f(x)\right\} \xrightarrow{h} \min(\inf) \quad \forall x; \quad \dots \text{ . Jejich přehled lze nalézt např. v:}$$

L. Devroye and L. Györfi Nonparametric Density Estimation: The L1 View. John Wiley, New York, 1985. Nebo v:

Luc Devroye and Gábor Lugosi: Variable Kernel Estimates: on the Impossibility of Tuning the Parameters. (School of Computer Science, McGill University, Montreal, Canada H3A2A 7, E-mail: luc@cs.mcgill.ca; Department of Economics, Pompeu Fabra University, Ramon Trias Fargas, 25-27, 08005 Barcelona, Spain, E-mail: gabor.lugosi@econ.upf.ed)

Příklady



Jádrové modely dvou proměnných

$$f_a(x, y) = \frac{1}{nh_x h_y} \sum_{i=1}^n \kappa_x \left(\frac{x-x_i}{h_x} \right) \kappa_y \left(\frac{y-y_i}{h_y} \right)$$

$$F_a(x, y) = \frac{1}{n} \sum_{i=1}^n K_x \left(\frac{x-x_i}{h_x} \right) K_y \left(\frac{y-y_i}{h_y} \right)$$

Je evidentní, že obě aproximace jsou „konzistentní“ vůči operaci přechodu k marginálním rozdělením.

$$\text{Tj.: } f_a(x) = \int_{-\infty}^{+\infty} f_a(x, y) dy = \frac{1}{nh_x h_y} \sum_{i=1}^n \kappa_x \left(\frac{x-x_i}{h_x} \right) \int_{-\infty}^{+\infty} \kappa_y \left(\frac{y-y_i}{h_y} \right) dy = \frac{1}{nh_x} \sum_{i=1}^n \kappa_x \left(\frac{x-x_i}{h_x} \right)$$

Pro další výpočty je užitečný následující triviální vztah: $\int_{-\infty}^{+\infty} x \frac{1}{h_x} \kappa_x \left(\frac{x-x_i}{h_x} \right) dx = x_i$.

K tomu aby platil, stačí předpoklady: $\int_{-\infty}^{+\infty} \kappa(x) dx = 1$, $\int_{-\infty}^{+\infty} x \kappa(x) dx = 0$.

Míry vztahu obou proměnných:

Korelace

$$E\{xy\} = \frac{1}{nh_x h_y} \sum_{i=1}^n \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} xy \kappa_x \left(\frac{x-x_i}{h_x} \right) \kappa_y \left(\frac{y-y_i}{h_y} \right) dx dy = \frac{1}{n} \sum_{i=1}^n x_i y_i \Rightarrow$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} \Rightarrow$$

$$\text{Corr}(x, y) = \frac{\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}}{\sqrt{h_x^2 + \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{h_y^2 + \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}}, \text{ a při volbě } h_x = \sqrt{\frac{s_x^2}{n}}, h_y = \sqrt{\frac{s_y^2}{n}}$$

$$\text{dostaneme } \text{Corr}(x, y) = \frac{\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}}{\sqrt{s_x^2} \sqrt{s_y^2}}.$$

Regrese

Podmíněné rozdělení:

$$f_a(x, y) = \frac{1}{nh_x h_y} \sum_{i=1}^n \kappa_x \left(\frac{x-x_i}{h_x} \right) \kappa_y \left(\frac{y-y_i}{h_y} \right) \Rightarrow$$

$$f_a(y/x) = \frac{\frac{1}{nh_x h_y} \sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right) \kappa_y\left(\frac{y-y_i}{h_y}\right)}{\frac{1}{nh_x} \sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right)} = \frac{\frac{1}{h_y} \sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right) \kappa_y\left(\frac{y-y_i}{h_y}\right)}{\sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right)} \quad \text{a odtud}$$

vlastní regrese:

$$r(x) = E\{y/x\} = \int_{-\infty}^{+\infty} y f_a(y/x) dy = \frac{\frac{1}{h_y} \sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right) \int_{-\infty}^{+\infty} y \kappa_y\left(\frac{y-y_i}{h_y}\right) dy}{\sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right)} = \frac{\sum_{i=1}^n y_i \kappa_x\left(\frac{x-x_i}{h_x}\right)}{\sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right)}$$

Skedastická funkce

$sked(x) = E\{(y - E\{y/x\})^2/x\} = E\{y^2/x\} - (E\{y/x\})^2$. Pak

$$E\{y^2/x\} = \frac{\sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right) \int_{-\infty}^{+\infty} y^2 \frac{1}{h_y} \kappa_y\left(\frac{y-y_i}{h_y}\right) dy}{\sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right)} = \frac{\sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right) (h_y^2 + y_i^2)}{\sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right)} = (F), \text{ protože}$$

$$\int_{-\infty}^{+\infty} y^2 \frac{1}{h_y} \kappa_y\left(\frac{y-y_i}{h_y}\right) dy = \int_{-\infty}^{+\infty} (h_y z + y_i)^2 \kappa_y(z) dz = \int_{-\infty}^{+\infty} (h_y^2 z^2 + 2h_y z y_i + y_i^2) \kappa_y(z) dz = h_y^2 + y_i^2,$$

$$\text{potom: } (F) = h_y^2 + \frac{\sum_{i=1}^n y_i^2 \kappa_x\left(\frac{x-x_i}{h_x}\right)}{\sum_{i=1}^n \kappa_x\left(\frac{x-x_i}{h_x}\right)} = h_y^2 + \sum_{i=1}^n y_i^2 \frac{\kappa_x\left(\frac{x-x_i}{h_x}\right)}{\sum_{j=1}^n \kappa_x\left(\frac{x-x_j}{h_x}\right)}.$$

Shrnuto:

$$sked(x) = h_y^2 + \sum_{i=1}^n y_i^2 \frac{\kappa_x\left(\frac{x-x_i}{h_x}\right)}{\sum_{j=1}^n \kappa_x\left(\frac{x-x_j}{h_x}\right)} - \left(\frac{\sum_{i=1}^n y_i \kappa_x\left(\frac{x-x_i}{h_x}\right)}{\sum_{j=1}^n \kappa_x\left(\frac{x-x_j}{h_x}\right)} \right)^2 = h_y^2 + \sum_{i=1}^n y_i^2 \frac{\kappa_x\left(\frac{x-x_i}{h_x}\right)}{\sum_{j=1}^n \kappa_x\left(\frac{x-x_j}{h_x}\right)} - r^2(x)$$

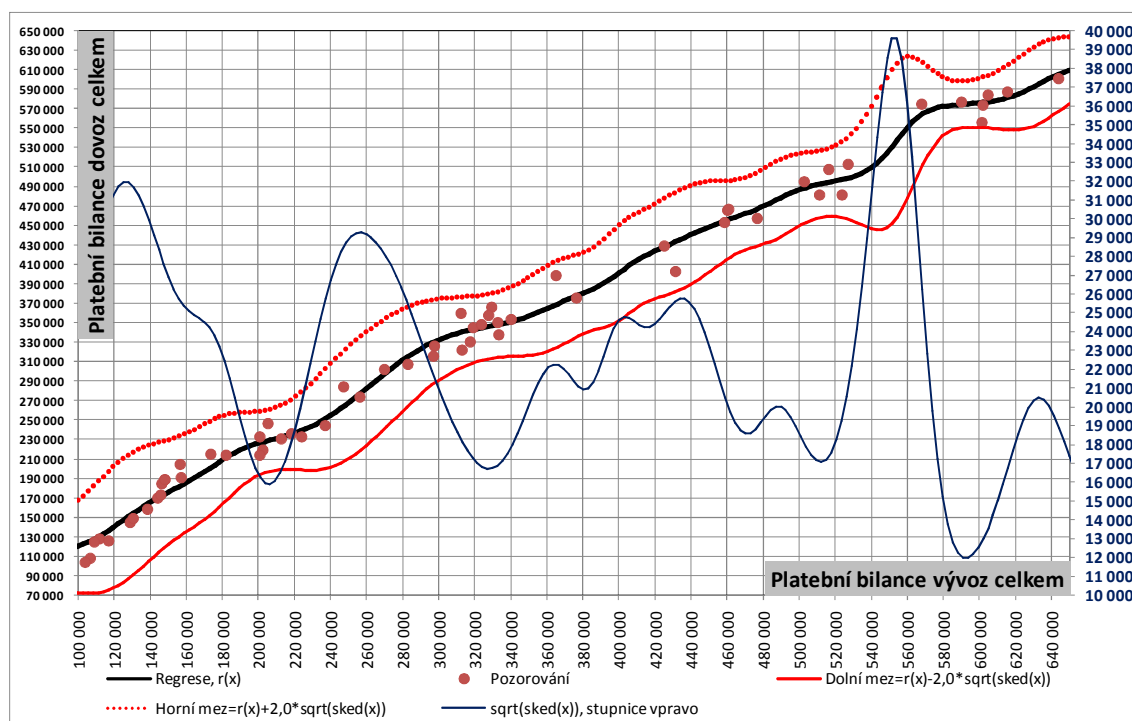
$$\text{V obou, jak regresní, tak skedastická funkci je důležitá řada funkcí } w_i(x) = \frac{\kappa_x\left(\frac{x-x_i}{h_x}\right)}{\sum_{j=1}^n \kappa_x\left(\frac{x-x_j}{h_x}\right)} =$$

vah. S tímto označením dostaneme:

$$r(x) = \sum_{i=1}^n y_i w_i(x), \text{ kde } \forall x \in \text{def} X; \quad w_i(x) \geq 0; \quad \sum_{i=1}^n w_i(x) = 1$$

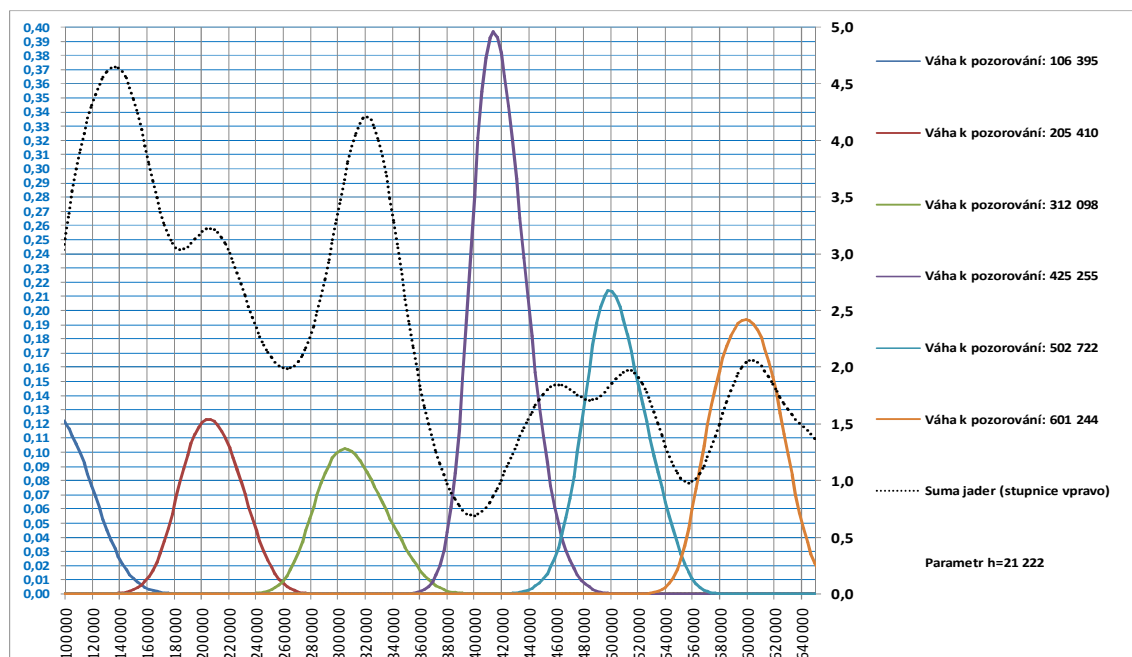
$$sked(x) = h_y^2 + \sum_{i=1}^n y_i^2 w_i(x) - r^2(x).$$

Příklad regrese a skedastické funkce



Gaussovo jádro $h_x=21\,222$, $h_y=19\,166$, $n=65$, $h = \sqrt{\frac{s_n^2}{n}}$. Data z časových řad nebyla stacionarizována.

Příklad váhových funkcí, odpovídajících výše zobrazené regresi



Doporučená a zdrojová literatura:

L. Devroye and L. Györfi	Nonparametric Density Estimation: The L1 View	John Wiley, New York, 1985.
A. Berlinet and L. Devroye	A Comparison of Kernel Density Estimates	Working paper, School of Computer Science, McGill University, Montreal, Canada, H3A 2A7.
Alfréd Rényi	Teorie pravděpodobnosti	ACADEMIA, Praha 1972.