



Politechnika  
Wrocławska

# Regresja jądrowa Nadaraya-Watsona

Jakub Koral

Wrocław, 8 grudnia 2021



# Estymacja gęstości

Rozpiszmy gęstość jako pochodną dystrybuanty

$$\begin{aligned} f(x) = F'(x) &= \lim_{h \rightarrow 0} \frac{F(x+h) - F(x-h)}{2h} \\ &= \lim_{h \rightarrow 0} \frac{\mathbb{P}[x-h \leq X \leq x+h]}{2h}. \end{aligned} \quad (1)$$

Teraz wystarczy oszacować prawdopodobieństwo, by otrzymać jądrowy estymator gęstości

$$\begin{aligned} \hat{f}(x, h) &= \frac{1}{2nh} \sum_{i=1}^n \mathbb{1}_{\{x-h \leq x_i \leq x+h\}} \\ &= \frac{1}{nh} \sum_{i=1}^n \frac{1}{2} \cdot \mathbb{1}_{\left\{\left|\frac{x-x_i}{h}\right| \leq 1\right\}} = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) \end{aligned} \quad (2)$$

Widać, że  $K(z) = \frac{1}{2} \cdot \mathbb{1}_{\{|z| \leq 1\}} = f_U(z)$ , gdzie  $U \sim \mathcal{U}(-1, 1)$ .

# Formalne definicje

## Definicja (Jądrowy estymator gęstości)

*Jądrowy estymator gęstości dla próbki  $\{x_i\}_{i=1}^n$  zadany jest wzorem*

$$\hat{f}(x, h) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right),$$

*gdzie  $K(\cdot)$  to jądro, a  $h$  to parametr wygładzenia.*

## Definicja (Jądro)

*Funkcję  $K : \mathbb{R} \rightarrow [0, \infty)$ , która spełnia następujące warunki*

- 1.  $\int_{\mathbb{R}} K(z) dz = 1$  (normalizacja),*
- 2.  $K(-z) = K(z)$  (symetria),*

*nazywamy jądrem. Często stosujemy zapis  $K_h(z) := \frac{1}{h} K\left(\frac{z}{h}\right)$ .*

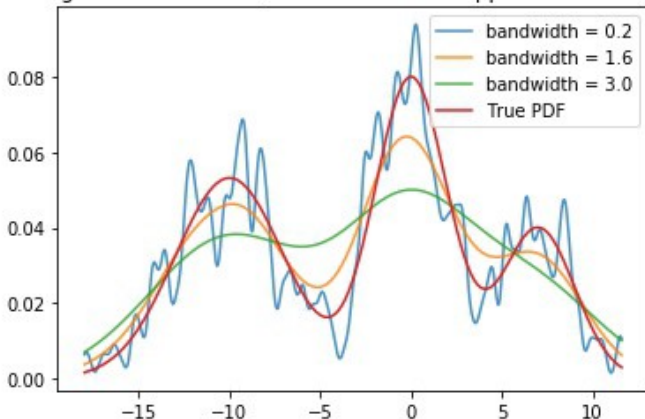
## Przykłady jąder

| Nazwa                           | Wzór   |
|---------------------------------|--|
| jądro jednostajne (prostokątne) | $\frac{1}{2} \cdot \mathbb{1}_{\{ z  \leq 1\}}$                        |
| jądro trójkątne                 | $(1 -  z ) \cdot \mathbb{1}_{\{ z  \leq 1\}}$                          |
| jądro normalne                  | $\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$                   |
| jądro Epanecznikowa             | $\frac{3}{4} (1 - z^2) \cdot \mathbb{1}_{\{ z  \leq 1\}}$              |
| jądro cosinusowe                | $\frac{\pi}{4} \cos \frac{\pi z}{2} \cdot \mathbb{1}_{\{ z  \leq 1\}}$ |

Wybór jądra ma niewielkie znaczenie (w porównaniu do wyboru parametru  $h$ ). Najlepsze pod względem scałkowanego błędu średniokwadratowego jest jądro Epanecznikowa.

# Dlaczego parametr wygładzenia jest istotny?

Effect of various bandwidth values  
The larger the bandwidth, the smoother the approximation becomes



**Rysunek:** Porównanie estymatorów gęstości jądrowej dla różnych parametrów  $h$  (ang. *bandwidth*).

## Dobór parametru wygładzenia

Nie da się wyznaczyć optymalnego (w sensie średniokwadratowym) estymatora parametru  $h$  bez znajomości postaci funkcji gęstości  $f$ . Można jednak znaleźć dobre jego przybliżenie metodą cross-validation. Zdefiniujmy

$$\text{LSCV}(h) := \int \hat{f}(x, h)^2 dx - \frac{2}{n(n-1)} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n K_h(x_i - x_j), \quad (3)$$

który nazywamy selektorem cross-validation o najmniejszym błędzie średniokwadratowym. Wtedy możemy oszacować  $h$  przy pomocy wzoru

$$\hat{h}_{\text{LSCV}} := \underset{h>0}{\operatorname{argmin}} \text{LSCV}(h). \quad (4)$$

## Regresja nieparametryczna

W przypadku regresji nieparametrycznej będziemy estymować funkcję

$$m(x) = \mathbb{E}[Y|X = x] = \int y f(y|x) dy = \int y \frac{f(x, y)}{f(x)} dy. \quad (5)$$

Skorzystamy z jądrowego estymatora gęstości

$$\begin{aligned} \hat{\mathbb{E}}[Y|X = x] &= \int y \frac{\sum_{i=1}^n K_h(x - x_i) K_h(y - y_i)}{\sum_{j=1}^n K_h(x - x_j)} dy \\ &= \frac{\sum_{i=1}^n K_h(x - x_i) \int y K_h(y - y_i) dy}{\sum_{j=1}^n K_h(x - x_j)} \\ &= \frac{\sum_{i=1}^n K_h(x - x_i) y_i}{\sum_{j=1}^n K_h(x - x_j)}. \end{aligned} \quad (6)$$

Estymator jądrowy Nadaraya-Watsona dany jest wzorem

$$\hat{m}_{NW}(x, h) = \frac{\sum_{i=1}^n K_h(x - x_i) y_i}{\sum_{j=1}^n K_h(x - x_j)}. \quad (7)$$

# Kod

Implementacją regresji jądrowej Nadaraya-Watsona w R jest chociażby funkcja `ksmooth`:

```
require(graphics)
library(MASS)

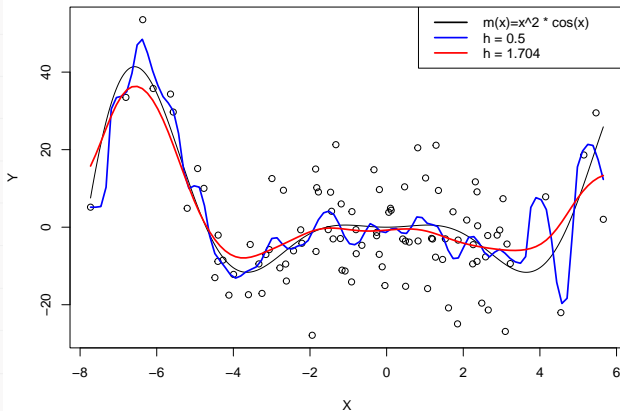
n <- 100
X <- rnorm(n, sd = 3)
m <- function(x) x^2 * cos(x)
eps <- rnorm(n, sd = 10)
Y <- m(X) + eps
h_LSCV <- ucv(X, nb = n)
xGrid <- seq(min(X), max(X), l = 500)

plot(X, Y)
lines(xGrid, m(xGrid))
lines(ksmooth(X, Y, "normal", bandwidth = 0.5),
col = "blue", lwd = 2)
lines(ksmooth(X, Y, "normal", bandwidth = h_LSCV),
col = "red", lwd = 2)
title("Przykład wykorzystania regresji jądrowej\nNadraya-Watsona")
legend(x="topright", c("m(x)=x^2 * cos(x)", "h = 0.5",
paste("h =", toString(round(h_LSCV, 3))))), lty = 1,
col=c("black", "blue", "red"), lwd = 2)
```



# Przykład

Przykład wykorzystania regresji jądrowej  
Nadaraya-Watsona



**Rysunek:** Zależność dla próbki  $\{x_i\}_{i=1}^{100}$  z rozkładu  $\mathcal{N}(0, 3)$  i  $y_i = x_i^2 \cos(x_i) + \varepsilon_i$ , gdzie  $\varepsilon_i \sim \mathcal{N}(0, 10)$ . Dorysowane krzywe regresji jądrowej Nadaraya-Watsona dla parametrów  $h = 0.5$  i  $h = \hat{h}_{\text{LSCV}}$ .

## Literatura:

Portugués E.G., *Notes for Predictive Modeling*:

- ▶ **Jądrowy estymator gęstości:**  
<https://bookdown.org/egarpor/PM-UC3M/npreg-npdens.html>
- ▶ **Regresja Nadaraya-Watsona:**  
<https://bookdown.org/egarpor/PM-UC3M/npreg-kre.html>

## Rysunek:

- ▶ <https://deeptai.org/machine-learning-glossary-and-terms/kernel-density-estimation>