

Regresja

- Regresją nazywamy związek między wartością oczekiwaną zmiennej zależnej i wielkościami zmiennych niezależnych.
- W sensie matematycznym regresja jest warunkową wartością oczekiwaną

$$E(y_i | x_{1i}, \dots, x_{ki}) = m(x_{1i}, \dots, x_{ki})$$

- Regresja liniowa:

$$m(x_{1i}, \dots, x_{ki}) = \beta_1 x_{1i} + \dots + \beta_k x_{ki}$$

- Regresji nieliniowa:

$$m(x_{1i}, \dots, x_{ki}) = m(x_{1i}, \dots, x_{ki}, \gamma_1, \dots, \gamma_k)$$

- W przypadku regresji nieparametrycznej nie zakładamy konkretnej formy $m(x_{1i}, \dots, x_{ki})$. Przymuje się jednak w tym przypadku pewne założenia dotyczące gładkości funkcji $m(x_{1i}, \dots, x_{ki})$.

Estymatory nieparametryczne funkcji gęstości

- Zdefiniujmy jądro (*kernel*) jako funkcję $K(u)$ o następujących własnościach:
 - ciągła, ograniczona, symetryczna
 - $\int K(u) du = 1$
- Najczęściej stosowane jądra:

| | |
|--------------|--|
| Jednostajne | $\frac{1}{2} \mathbb{I}(u \leq 1)$ |
| Trójkątne | $(1 - u) \mathbb{I}(u \leq 1)$ |
| Epanechnikov | $\frac{3}{4} (1 - u^2) \mathbb{I}(u \leq 1)$ |
| Gaussssian | $\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}u^2\right)$ |

- Estymator jądrowy Rosenblatta-Parzena funkcji gęstości:

$$\hat{f}_{h_n}(x) = n^{-1} \sum_{i=1}^n \frac{1}{h_n} K\left(\frac{x - X_i}{h_n}\right)$$

- h_n - jest oknem (*bandwidth*).
 - szerokość okna determinuje stopień wygładzenia funkcji gęstości
 - szersze okno zmniejsza wariancję estymatora ale zwiększa obciążenie
 - okno wybieramy często arbitralnie na podstawie rysunku
 - istnieją także wzory na optymalną szerokość okna
- Uogólnienie na kilka wymiarów

$$\hat{f}_{h_{1n}, \dots, h_{kn}}(x_1, \dots, x_K) = n^{-1} \sum_{i=1}^n \prod_{k=1}^K \frac{1}{h_{kn}} K\left(\frac{x_k - X_{ki}}{h_{kn}}\right)$$

Regresja nieparametryczna, estymator jądrowy

- Dla przypadku regresji dla jednej zmiennej

$$m(x) = \mathbb{E}(y|x) = \int y f(y|x) dy = \frac{1}{f(x)} \int y f(y, x) dy$$

- Oszacowanie $f(y, x)$ jest dane wzorem

$$\hat{f}_{h_n, g_n}(y, x) = n^{-1} \sum_{i=1}^n \frac{1}{h_n} K\left(\frac{x - X_i}{h_n}\right) \frac{1}{g_n} K\left(\frac{y - Y_i}{g_n}\right)$$

- Przybliżony wzór na

$$\int y f(y, x) dy \approx \sum_{i=1}^n \hat{f}_{h_n, g_n}(Y_i, x) Y_i = n^{-1} \sum_{i=1}^n \frac{1}{h_n} K\left(\frac{x - X_i}{h_n}\right) Y_i$$

ma tą samą wartość dla wszystkich elementów sumy i możemy przyjąć, że $\frac{1}{g_n} K(0) = 1$.

- Estymatorem jądrowy Nadaraya-Watsona funkcji regresji ma więc postać

$$\begin{aligned} m_h(x) &= \frac{n^{-1} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - X_i}{h}\right) Y_i}{n^{-1} \sum_{i=1}^n \frac{1}{h_n} K\left(\frac{x - X_i}{h_n}\right)} \\ &= \sum_{i=1}^n W_{ih_n} Y_i \end{aligned}$$

gdzie W_{ih_n} są wagami postaci

$$W_{ih} = \frac{\frac{1}{h} K\left(\frac{x-X_i}{h}\right) Y_i}{\sum_{i=1}^n \frac{1}{h_n} K\left(\frac{x-X_i}{h_n}\right)}$$

- Uogólnienia wielowymiarowe opiera się na wielowymiarowym opiera się na wielowymiarowym uogólnieniu estymatora funkcji gęstości

Lokalna regresja wielomianowa

- Wiemy, że wystarczająco gładką funkcję można przybliżyć wokół t za pomocą następującego wzoru Taylora

$$m(t) \approx m(x) + m'(x)(t-x) + \dots + m^{(p)}(x)(t-x)^p \frac{1}{p!}$$

- Lokalna regresja wielomianowa polega na przeprowadzeniu ważonej regresji następującej postaci:

$$\min_{\beta} \sum_{i=1}^n \left\{ Y_i - \beta_0 - \beta_1 (X_i - x) - \dots - \beta_p (X_i - x)^p \right\} \frac{1}{h} K \left(\frac{x - X_i}{h} \right)$$

- Wagi nadawane poszczególnym obserwacjom związane są z ich odległością od szacowanego punktu



Przekleństwo wymiaru

- Zarówno funkcja gęstości jak funkcja regresji przy estymacji metodami nieparametrycznymi jest szacowana na podstawie wielkości sąsiednich obserwacji
- Dla dużej liczby wymiarów precyzja oszacowań szybko spada, ponieważ liczba obserwacji, która sąsiaduje z dowolnie dobranym punktem się szybko zmniejsza wraz ze wzrostem liczby wymiarów
- Dodatkowy problem:
 - jak zdefiniować sąsiedztwo w wielu wymiarach w sytuacji, w której poszczególne zmienne mają różne jednostki

- Rozwiązaniem jest przyjęcie mniej ogólnej formy funkcji gęstości
- Jednym z proponowanych rozwiązań jest przyjęcie formy addytywnej

$$m(x_1, \dots, x_k) = a + f_1(x_1) + \dots + f_K(x_k)$$

- Można do takiego modelu dodać nieparametryczne interakcje między zmiennymi

$$m(x_1, \dots, x_k) = a + f_1(x_1) + f_2(x_k) + f_{12}(x_1, x_2)$$

- Możliwe jest także częściowe sparametryzowanie modelu

$$m(x_1, \dots, x_k) = a + \beta_1 x_1 + f_2(x_2) \dots + f_K(x_k)$$