

Regresja na kwantylach

Jerzy Mycielski

2008

- W przypadku regresji dla wartości oczekiwanej modelujemy zależność między wartością oczekiwaną a zmiennymi objaśniającymi
- Hiperpłaszczyzna regresji y na \mathbf{x} dla wartości oczekiwanej jest warunkową wartością oczekiwaną

$$E(y|\mathbf{x}) = \mu(\mathbf{x})$$

- Intuicyjnie, hiperpłaszczyzna regresji podaje zależność między oczekiwaną wielkością y dla znanego \mathbf{x}

- Kwantyl τ jest taką wielkością która jest większa lub równa od y z prawdopodobieństwem τ .
- Bezwarunkowy kwantyl τ będziemy więc definiować jako takie ξ , dla którego

$$F(\xi_\tau) = \tau.$$

- Dla znanej postaci dystrybuanty F kwantyl τ można policzyć jako

$$\xi_\tau = F_y^{-1}(\tau).$$

- W przypadku regresji na kwantylach modelujemy zależność między wielkością kwantyla τ a wielkością zmiennych niezależnych.
- Warunkowy kwantyl τ można więc zdefiniować jako

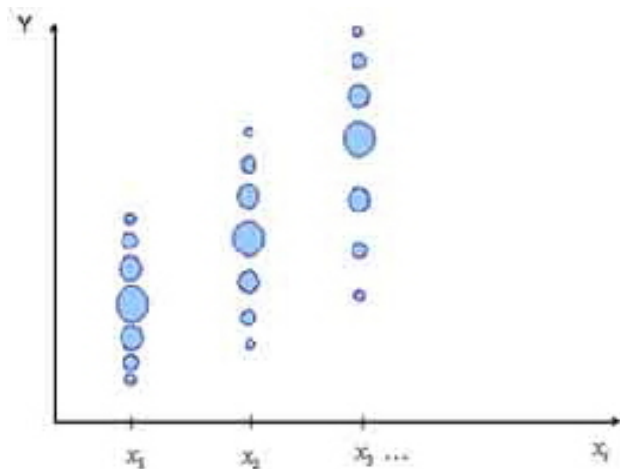
$$F_y(\xi_\tau | \mathbf{x}) = \tau.$$

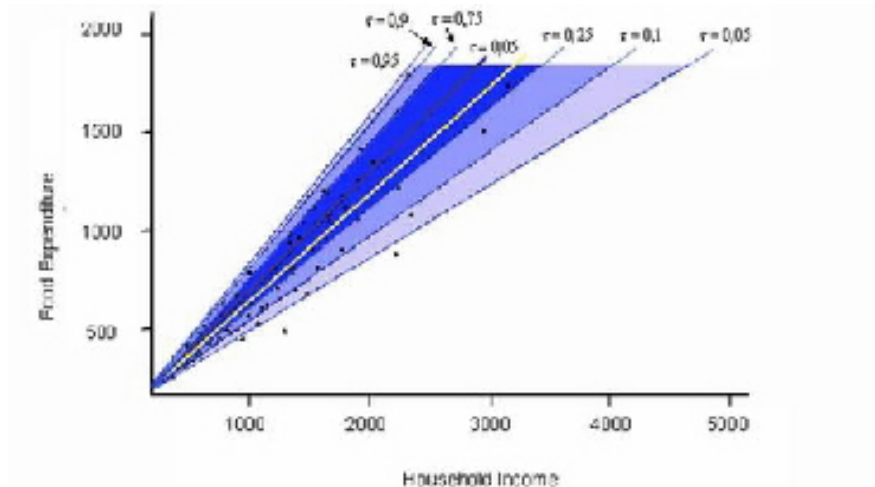
- Wynika z tego, że

$$\xi_\tau(\mathbf{x}) = F_y^{-1}(\tau | \mathbf{x})$$

- Zauważmy, że dla każdego kwantyla, postać funkcji regresji może wyglądać nieco inaczej

Kwantyle





- W przypadku regresji dla wartości oczekiwanej standardowo stosuje się metodę najmniejszych kwadratów.
- Zakładamy, że $\mu(\mathbf{x})$ jest dana funkcją parametryczną $\mu(\mathbf{x}, \boldsymbol{\beta})$ o nieznanym wektorze parametrów $\boldsymbol{\beta}$.
- Wektor $\boldsymbol{\beta}$ znajdujemy minimalizując sumę kwadratów reszt:

$$\min_{\boldsymbol{\beta} \in \mathbf{R}^k} \sum_{i=1}^N [y_i - \mu(\mathbf{x}_i, \boldsymbol{\beta})]^2$$

Estymacja parametrów - regresja na kwantylach

Mediana

- W przypadku regresji na kwantylach estymację można także sprowadzić do rozwiązywania problemu maksymalizacyjnego.
- Rozpatrzmy najprostszy przypadek szukania mediany y ($\tau = 0.5$) poprzez rozwiązywanie następującej funkcji celu

$$\min_{\zeta_{0.5} \in \mathbf{R}} \sum_{i=1}^N |y_i - \zeta_{0.5}|$$

- Warunki pierwszego rzędu dla tego problemu (dla N nieparzystego i $y_i \neq y_j$ dla $i \neq j$) będą miały następującą postać:

$$\sum_{i=1}^N [\mathbb{I}(y_i > \zeta_{0.5}) - \mathbb{I}(y_i < \zeta_{0.5})] = 0$$

lub

$$\frac{\sum_{i=1}^N \mathbb{I}(y_i > \zeta_{0.5})}{\sum_{i=1}^N \mathbb{I}(y_i < \zeta_{0.5})} = 1$$

- Warunek pierwszego rzędu będzie więc spełniony dla $\tilde{\zeta}_{0.5}$ równemu y_i , które jest medianą z próby.
- W ogólniejszym przypadku estymacji parametrów regresji na medianie rozwiązujemy następujący problem maksymalizacyjny:

$$\min_{\beta \in \mathbf{R}^k} \sum_{i=1}^N |y_i - \tilde{\zeta}_{0.5}(\mathbf{x}, \beta)|$$

- Zauważmy, że estymator dla regresji dla mediany jest równoważny estymatorowi *LAD* (**L**east **A**bsolute **D**eviation)

- Dla dowolnego τ oszacowanie kwantyla z próby można znaleźć rozwiązując zadanie:

$$\min_{\xi_\tau \in \mathbb{R}} \sum_{i=1}^N \rho_\tau(y_i - \xi_\tau)$$

gdzie

$$\rho_\tau(x) = \begin{cases} \tau x & \text{dla } x > 0 \\ -(1 - \tau)x & \text{dla } x < 0 \end{cases}$$

- Warunki pierwszego rzędu w tym przypadku będą miały postać

$$\frac{\sum_{i=1}^N \mathbb{I}(y_i > \xi_\tau)}{\sum_{i=1}^N \mathbb{I}(y_i < \xi_\tau)} = \frac{1 - \tau}{\tau}$$

- Estymatorem jest odpowiedni kwantyl empiryczny.

- W ogólniejszym przypadku estymacji parametrów regresji na medianie rozwiązujemy następujący problem maksymalizacyjny:

$$\min_{\beta \in \mathbf{R}^K} \sum_{i=1}^N \rho_{\tau}(y_i - \xi_{\tau}(\mathbf{x}, \beta))$$

- Rozwiązanie tego problemu maksymalizacji znajduje się metodami programowania liniowego (simplex)

Regresja na kwantylach - zalety

- Przeprowadzając regresję na szeregu kwantyli, opisujemy zależność całego rozkładu zmiennej zależnej od zmiennych niezależnych a nie tylko wartości oczekiwanej zmiennej zależnej
- Heteroskedastyczność można łatwo wykryć analizując wyniki regresji na kwantylach
- W przypadku występowania heteroskedastyczności estymacja regresji na medianie może okazać się bardziej efektywnym sposobem szukania wartości parametrów niż regresja na wartości oczekiwanej
- Kwantyl zmiennej losowej przekształconej za pomocą przekształcenia monotonicznie rosnącego $g(y)$ jest równy przekształconemu w ten sam sposób kwantylowi oryginalnej zmiennej

$$\tilde{\zeta}_{\tau}(\mathbf{x}) = F_{g(y)}^{-1}(\tau | \mathbf{x}) = g(F_y^{-1}(\tau | \mathbf{x}))$$

- Regresja na kwantylach jest odporna na problem outlier'ów (obserwacji błędnych)

- Największy problem związany z regresją na kwantylach związany jest z brakiem wzorów analitycznych na macierz wariancji i kowariancji uzyskanych oszacowań.
- W związku z tym trudniejsze jest znajdowanie przedziałów ufności oraz weryfikacja hipotez statystycznych
- Oszacowania macierzy wariancji i kowariancji można znaleźć metodą bootstrap