



**AGH**

Akademia Górniczo-Hutnicza  
im. Stanisława Staszica  
w Krakowie

---

Praca licencjacka

# Testowanie regresji liniowej przeciwko wypukłej

Grzegorz Mika

Kierunek: Matematyka

Nr albumu: 267543

Promotor  
dr Konrad Nosek



Wydział Matematyki Stosowanej

---

Kraków 2016

## Oświadczenie autora

*Ja, niżej podpisany Grzegorz Mika oświadczam, że praca ta została napisana samodzielnie i wykorzystywała (poza zdobytą na studiach wiedzę) jedynie wyniki prac zamieszczonych w spisie literatury.*

.....

(Podpis autora)

## Oświadczenie promotora

*Oświadczam, że praca spełnia wymogi stawiane pracom licencjackim.*

.....

(Podpis promotora)

# Spis treści

Wstęp	1
1 Stożki wypukłe	3
2 Regresja wypukła	5
3 Test statystyczny i jego rozkład	8
Bibliografia	12

# Wstęp

Rozważmy pewien zestaw danych  $\{(x_i, y_i)\}_{i=1,2,\dots,n}$  i spróbujmy dopasować pewną funkcję  $f$  do danych według modelu

$$y_i = f(x_i) + \varepsilon_i$$

gdzie zakładamy, że błędy  $\varepsilon_i$  są niezależnymi zmiennymi losowymi o tym samym rozkładzie normalnym.

Najprostszym związkiem między obserwacjami  $x_i$  a odpowiedziami  $y_i$  jest zależność liniowa, możliwy jest jednak również inny związek między obserwacjami a odpowiedziami, co prowadzi do sformułowania hipotezy

$$H_0: f(x) = ax + b \quad vs. \quad H_1: f \in \mathcal{F},$$

gdzie  $\mathcal{F}$  jest klasą funkcji wypukłych.

W niniejszej pracy postaramy się skonstruować odpowiedni do postawionego problemu test statystyczny. Zaproponowane zostanie rozwiązanie oparte o iloraz wiarygodności w przypadku modelu regresji z ograniczeniami w postaci nierówności.

W pierwszym rozdziale zostaną omówione podstawowe własności stożków wypukłych traktowanych jako podzbiór przestrzeni liniowej. Drugi rozdział będzie traktował o konstrukcji estymatora regresji wypukłej jako rzutu wektora danych na wypukły stożek wielościenne. W trzecim rozdziale zostanie wyznaczony rozkład szukanego testu w przypadku ze znaną wariancją błędu obserwacji.

Praca została napisana na podstawie [2], natomiast rozdział o algorytmie baz prymalno-dualnych został napisany w dużym stopniu na podstawie [1].

# 1 Stożki wypukłe

Poszukiwany test zostanie wyznaczony metodą rzutowania wektora danych na wielościan powstały w wyniku narzuconych ograniczeń liniowych. W tym rozdziale zostaną przedstawione podstawowe definicje i własności wypukłych stożków wielościenne użyteczne w dalszych rozważaniach.

**Definicja 1 (Ortant).** *Ortantem w  $n$ -wymiarowej przestrzeni  $\mathbb{R}^n$  nazywamy podzbiór powstały przez ograniczenie każdej ze współrzędnych do bycia nieujemną lub niedodatnią, czyli*

$$O = \{(x_1, x_2, \dots, x_n) \in \mathbb{R}^n : \epsilon_i x_i \geq 0, |\epsilon_i| = 1, i = 1, 2, \dots, n\} \quad (1.1)$$

**Definicja 2 (Stożek wypukły).** *Niech  $V$  będzie przestrzenią wektorową. Stożkiem wypukłym nazywamy przecięcie skończonej ilości półprzestrzeni przestrzeni  $V$ .*

Rozważmy  $n$ -wymiarową przestrzeń wektorową  $V$ . Dowolną półprzestrzeń  $H$  przestrzeni  $V$  można wyrazić jako

$$H = \{(x_1, x_2, \dots, x_n) \in \mathbb{R}^n : a_1 x_1 + a_2 x_2 + \dots + a_n x_n \geq b\} \quad (1.2)$$

gdzie  $a_1, a_2, \dots, a_n, b$  są pewnymi, ustalonymi liczbami rzeczywistymi.

Korzystając z tego przedstawienia możemy dowolny stożek wypukły  $P$  zapisać jako

$$K = \bigcap_{i=1}^m H_i, \quad (1.3)$$

gdzie

$$H_j = \{\mathbf{x} \in \mathbb{R}^n : \sum_{i=1}^n a_i^j x_i \geq b^j\}. \quad (1.4)$$

Stąd możemy zapisać, że

$$K = \left\{ (x_1, x_2, \dots, x_n) \in \mathbb{R}^n : \begin{cases} \sum_{i=1}^n a_i^1 x_i \geq b^1 \\ \sum_{i=1}^n a_i^2 x_i \geq b^2 \\ \vdots \\ \sum_{i=1}^n a_i^m x_i \geq b^m \end{cases} \right\} \quad (1.5)$$

co będziemy zapisywać skrótowo jako

$$K = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} \geq \mathbf{b}\} \quad (1.6)$$

gdzie

$$\mathbf{A} = \begin{bmatrix} a_1^1 & a_2^1 & \dots & a_n^1 \\ a_1^2 & a_2^2 & \dots & a_m^2 \\ \vdots & \vdots & \dots & \vdots \\ a_1^n & a_2^n & \dots & a_m^n \end{bmatrix}, \mathbf{b}^T = (b^1, b^2, \dots, b^m), \mathbf{x} = (x_1, x_2, \dots, x_n)^T. \quad (1.7)$$

Symbolem  $\langle \cdot, \cdot \rangle$  będziemy oznaczać iloczyn skalarny w przestrzeni wektorowej  $V$ . Oznaczmy przez  $\gamma_i$  kolejne wiersze macierzy  $-\mathbf{A}$ . Bez straty ogólności możemy założyć ponadto, że tworzą one układ wektorów liniowo niezależnych, gdyż w przeciwnym wypadku któreś ograniczenie stanowiłoby kombinację pozostałych oraz że  $m \leq n$ . Wtedy stożek  $K$  możemy też zapisać w sposób

$$K = \{\mathbf{x} \in \mathbb{R}^n : \langle \mathbf{x}, \gamma_i \rangle \leq 0, i = 1, 2, \dots, m\} \quad (1.8)$$

Uzupełniając zbiór wektorów  $\{\gamma_i\}$  do bazy przestrzeni  $\mathbb{R}^n$  o wektory ortogonalne i definiując bazę dualną złożoną z wektorów  $\beta_i$  w następujący sposób

$$\beta_i^T \gamma_j = \begin{cases} -1, & i = j \\ 0, & i \neq j \end{cases} \quad (1.9)$$

możemy zapisać równoważne przedstawienie stożka  $K$

$$K = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} = \sum_{i=1}^m b_i \beta_i + \sum_{i=m+1}^n c_i \beta_i, b_i \geq 0, c_i \in \mathbb{R}\} \quad (1.10)$$

**Twierdzenie 1. Przedstawienia**

$$K = \{\mathbf{x} \in \mathbb{R}^n : \langle \mathbf{x}, \gamma_i \rangle \leq 0, i = 1, 2, \dots, m\} \quad (1.11)$$

$$K = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} = \sum_{i=1}^m b_i \beta_i + \sum_{i=m+1}^n c_i \beta_i, b_i \geq 0, c_i \in \mathbb{R}\} \quad (1.12)$$

są równoważne.

*Dowód.* Wektory  $\beta_i, \gamma_i, i = 1, 2, \dots, n$  spełniają zależność  $\beta_i^T \gamma_i = -1$  oraz  $\beta_i^T \gamma_j = 0, i \neq j$ . Oznaczając przez  $\mathbf{B}, \mathbf{C}$  macierze, których kolumnami są odpowiednio wektory  $\beta_i, \gamma_i$ , związek ten możemy przedstawić jako  $\mathbf{B}^T \mathbf{C} = -\mathbf{I}$ . Ze związku  $\mathbf{B}^T \mathbf{C} = -\mathbf{I}$  dostajemy, że  $\mathbf{C}^T \mathbf{x} = -\mathbf{B}^{-1} \mathbf{x}$ . Wyrażenia  $\langle \mathbf{x}, \gamma_i \rangle, i = 1, 2, \dots, m$  są pierwszymi  $m$  współrzędnymi  $\mathbf{C}^T \mathbf{x}$ . Zatem wektor  $\mathbf{x}$  wyrażony w bazie złożonej z wektorów  $\beta_i$  ma pierwsze  $m$  współrzędnych nieujemnych wtedy i tylko wtedy, gdy  $\langle \mathbf{x}, \gamma_i \rangle \leq 0, i = 1, 2, \dots, m$ , co dowodzi równoważności przedstawień.  $\square$

## 2 Regresja wypukła

Podobnie jak w przypadku zwykłego estymatora regresji liniowej, który jest rzutem wektora danych na pewną mniej wymiarową podprzestrzeń, tak w przypadku estymatora regresji wypukłej jest on rzutem na pewen wielościan wypukły powstały w wyniku stosownych ograniczeń.

Zbiór nad którym będziemy minimalizować kwadrat błędu powstaje w sposób następujący. Przypuśmy, że wartości  $x$  są różne między sobą i uporządkowane rosnąco oraz niech  $\theta_i = f(x_i), i = 1, 2, \dots, n$ . Rozważając kawałkami liniowe przybliżenie funkcji regresji z węzłami w punktach  $x_i$ , wymóg wypukłości może zostać zapisany jako zbiór ograniczeń w postaci nierówności linowych następującej postaci:

$$\theta_i(x_{i+2} - x_{i+1}) - \theta_{i+1}(x_{i+2} - x_i) + \theta_{i+2}(x_{i+1} - x_i) \geq 0, i = 1, 2, \dots, n-2 \quad (2.1)$$

Zgodnie z definicją 1 możemy zbiór tych ograniczeń zapisać jako

$$K = \{\mathbf{A}\boldsymbol{\theta} \geq 0\} \quad (2.2)$$

gdzie  $\mathbf{A}$  jest rzeczywistą macierzą wymiaru  $(n-2) \times n$ .

W tym momencie problem znalezienia estymatora regresji wypukłej przyjmuje postać

$$\text{minimalizuj } \|\mathbf{y} - \boldsymbol{\theta}\|^2 \text{ po } \boldsymbol{\theta} \in K, \quad (2.3)$$

gdzie  $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}}$ .

Niech  $B = (\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n)$  oznacza bazę kanoniczną przestrzeni  $\mathbb{R}^n$ . Oznaczmy przez  $\boldsymbol{\gamma}_i = -\mathbf{e}_i^T \mathbf{A}^T, i = 1, 2, \dots, n-2$ . Wtedy zbiór  $K$  możemy zapisać jako  $K = \{\boldsymbol{\theta} \in \mathbb{R}^n: -\mathbf{e}_i^T \mathbf{A}^T \boldsymbol{\theta} \leq 0, i = 1, 2, \dots, n-2\} = \{\boldsymbol{\theta} \in \mathbb{R}^k: \langle \boldsymbol{\gamma}_i, \boldsymbol{\theta} \rangle \leq 0, i = 1, 2, \dots, n-2\}$ .

Z określenia macierzy  $\mathbf{A}$  oraz wektorów  $\boldsymbol{\gamma}_i, i = 1, 2, \dots, n-2$ , widać, że tworzą one układ wektorów liniowo niezależnych. Zatem zbiór  $B'_\gamma = \{\boldsymbol{\gamma}_i, i = 1, 2, \dots, n-2\}$  można uzupełnić do bazy  $B_\gamma$  przestrzeni  $\mathbb{R}^n$  o wektory  $\boldsymbol{\gamma}_{n-1}, \boldsymbol{\gamma}_n$  tak, żeby były one ortogonalne do wszystkich wektorów z bazy  $B'_\gamma$ . Łatwo sprawdzić, że warunek ten spełniają wektory  $\boldsymbol{\gamma}_{n-1} = \mathbf{1}$  oraz  $\boldsymbol{\gamma}_n = (x - \bar{x}\mathbf{1})$ , gdzie  $x = (x_1, x_2, \dots, x_n)$ ,  $\bar{x}$  oznacza wartość średnią,  $\mathbf{1} = (1, 1, \dots, 1)^T$ , a norma  $\|\cdot\|$  jest normą zadaną wcześniej.

Teraz możemy zdefiniować bazę  $B_\beta = \{\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_n\}$  dualną do bazy  $B_\gamma$  w następujący sposób:

$$\boldsymbol{\beta}_i^T \boldsymbol{\gamma}_j = \begin{cases} -1, & i = j \\ 0, & i \neq j \end{cases} \quad (2.4)$$

Oznaczając przez  $\mathbf{B}$  i  $\mathbf{C}$  macierze, których kolumnami są odpowiednio wektory  $\beta_i$  i  $\gamma_i$  związek między nimi możemy wyrazić jako

$$\mathbf{B}^T \mathbf{C} = -\mathbf{I}, \quad (2.5)$$

gdzie  $\mathbf{I}$  oznacza macierz jednostkową.

Niech  $E$  oznacza podprzestrzeń przestrzeni  $\mathbb{R}^n$  rozpiętą przez wektory  $\beta_{n-1}, \beta_n$ , natomiast  $\mathcal{L}(K)$  oznacza przestrzeń rozpiętą przez wektory  $\beta_i, i = 1, 2, \dots, n-2$ . Przestrzeń  $E$  oraz  $\mathcal{L}(K)$  są do siebie ortogonalne, zatem wektor obserwacji  $\mathbf{y}$  możemy zapisać jako sumę  $\mathbf{y}_E + \mathbf{z}$ , gdzie  $\mathbf{y}_E$  i  $\mathbf{z}$  są rzutami wektora  $\mathbf{y}$  odpowiednio na podprzestrzeń  $E$  oraz  $\mathcal{L}(K)$ .

**Przykład 1.** *Prześledźmy powyższe rozważania na przykładzie dla przypadku czterowymiarowego i równoodległych punktów  $x_i$  odległych o 1.*

*Macierz ograniczeń  $G$  przybiera wtedy postać*

$$G = \begin{bmatrix} 1 & -2 & 1 & 0 \\ 0 & 1 & -2 & 1 \end{bmatrix} \quad (2.6)$$

*Zatem stożek powstały z ograniczeń jest postaci*

$$K = \{\theta \in \mathbb{R}^4: \begin{bmatrix} 1 & -2 & 1 & 0 \\ 0 & 1 & -2 & 1 \end{bmatrix} \theta \geq 0\} \quad (2.7)$$

*Baza wektorów  $B_\gamma$  jest postaci*

$$B_\gamma = \left( (-1, 2, -1, 0), (0, -1, 2, -1), (1, 1, 1, 1), \left(-\frac{3}{2}, -\frac{1}{2}, \frac{1}{2}, \frac{3}{2}\right) \right) \quad (2.8)$$

*Wektory  $\beta_i$  spełniające warunek  $\langle \beta_i, \gamma_j \rangle = \begin{cases} -1, & i=j \\ 0, & i \neq j \end{cases}$  przybierają następującą postać*

$$B_\beta = ((3, -4, -1, 2), (2, -1, -4, 3), (-1, -1, -1, -1), (3, 1, -1, -3)) \quad (2.9)$$

*Przestrzeń na którą będziemy rzutować wektor danych przybierają postać*

$$E = \{t_1(-1, -1, -1, 1) + t_2(3, 1, -1, -3), t_1, t_2 \in \mathbb{R}\} \quad (2.10)$$

$$\mathcal{L}(K) = \{t_1(3, -4, -1, 2) + t_2(2, -1, -4, 3), t_1, t_2 \in \mathbb{R}\} \quad \diamond$$

Zadanie znalezienia rzutu wektora danych na stożek  $K$  sprowadza się w tym momencie do znalezienia rzutu jego składowych na stożek  $K$ . Wszystkie elementy podprzestrzeni  $E$  należą do stożka  $K$ , więc rzut wektora  $\mathbf{y}$  na stożek  $K$  jest tym samym co jego rzut na podprzestrzeń  $E$  i wyraża się wzorem

$$\mathbf{y}_E = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad (2.11)$$



gdzie macierz  $\mathbf{X}$  jest podmacierzą macierzy  $\mathbf{A}$  złożoną z pierwszych  $n - 2$  kolumn. Pozostaje zagadnienie znalezienia rzutu  $\mathbf{z}$  na stożek  $K$ . Sprowadza się ono do znalezienia rzutu  $\mathbf{z}$  na stożek

$$K' = K \cap \mathcal{L}(K) = \{\boldsymbol{\theta} \in \mathbb{R}^n : \boldsymbol{\theta} = \sum_{i=1}^{n-2} b_i \boldsymbol{\beta}_i, b_i \geq 0\}. \quad (2.12)$$

W pracy [1] zostało pokazane, że przestrzeń  $\mathcal{L}(K)$  może zostać podzielona na  $2^{n-2}$  rozłącznych regionów w taki sposób, że każdy z nich może być opisany jako nieujemny ortant w bazie  $B_J = \{\boldsymbol{\beta}_i, i \in J, \boldsymbol{\gamma}_i, i \in L \setminus J\}$ , gdzie  $J$  jest pewnym podzbiorem zbioru  $L = \{1, 2, \dots, n - 2\}$ . Zatem każdy element  $z$  należący do  $\mathcal{L}(K)$  może być przedstawiony w następujący sposób

$$z = \sum_{i \in J} b_i \boldsymbol{\beta}_i + \sum_{i \in L \setminus J} c_i \boldsymbol{\gamma}_i, \quad b_i > 0, c_i \geq 0 \quad (2.13)$$

Dla dowolnego zbioru  $J \subset L$   $B_J$  jest bazą przestrzeni  $\mathcal{L}(K)$ , ponadto  $\boldsymbol{\beta}_i, i \in J$  oraz  $\boldsymbol{\gamma}_i, i \in L \setminus J$  są wzajemnie ortogonalne, zatem rzutem  $z$  na  $K'$  jest wektor postaci

$$z_{K'} = \sum_{i \in J} b_i \boldsymbol{\beta}_i, \quad b_i > 0 \quad (2.14)$$

Podsumowując, dowolny wektor  $\mathbf{y}$  z przestrzeni  $\mathbb{R}^n$  można przedstawić w następującej postaci

$$\mathbf{y} = \mathbf{z} + \mathbf{y}_E = \sum_{i \in J} b_i \boldsymbol{\beta}_i + \sum_{i \in L \setminus J} c_i \boldsymbol{\gamma}_i + d_1 \boldsymbol{\gamma}_{n-1} + d_2 \boldsymbol{\gamma}_n, \quad b_i > 0, c_i \geq 0, d_1, d_2 \in \mathbb{R} \quad (2.15)$$

Wtedy rzut tego wektora na stożek

$$K = \{\mathbf{A}\boldsymbol{\theta} \geq 0\} \quad (2.16)$$

jest postaci

$$\hat{\boldsymbol{\theta}} = \sum_{i \in J} b_i \boldsymbol{\beta}_i + d_1 \boldsymbol{\gamma}_{n-1} + d_2 \boldsymbol{\gamma}_n. \quad (2.17)$$

Natomiast wektor błędu  $\hat{\boldsymbol{\rho}} = \mathbf{y} - \hat{\boldsymbol{\theta}}$  jest postaci

$$\hat{\boldsymbol{\rho}} = \sum_{i \in L \setminus J} c_i \boldsymbol{\gamma}_i. \quad (2.18)$$

### 3 Test statystyczny i jego rozkład

Na początek wprowadzimy kilka oznaczeń i udowodnimy cztery lematy z których skorzystamy w dalszej części rozważań.

$$C_{L \setminus J} = \{\mathbf{y} \in \mathbb{R}^n : y = \sum_{i \in L \setminus J} b_i \boldsymbol{\gamma}_i + \sum_{i \in J} c_i \boldsymbol{\beta}_i + d_1 \boldsymbol{\gamma}_{n-1} + d_2 \boldsymbol{\gamma}_n, b_i > 0, c_i \geq 0, d_1, d_2 \in \mathbb{R}\} \quad (3.1)$$

$$S_{L \setminus J} = \text{span}\{\boldsymbol{\gamma}_i, i \in L \setminus J\} \quad (3.2)$$

$$d = |L \setminus J| = n - 2 - |J| \quad (3.3)$$

Niech

$$\mathbf{A}_{L \setminus J} \quad (3.4)$$

oznacza macierz wymiaru  $(n - 2) \times n$  taką, że pierwsze  $d$  wierszy to wektory  $-\boldsymbol{\gamma}_i, i \in L \setminus J$  natomiast pozostałe  $n - 2 - d$  wierszy to wektory  $-\boldsymbol{\beta}_i, i \in J$ .

**Lemat 1.** Niech  $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n)^T \sim N_n(0, \mathbf{I})$  oraz niech  $\mathbf{A}$  będzie rzeczywistą macierzą wymiaru  $m \times n$ . Wtedy rozkładem warunkowym  $\|\mathbf{Z}\|^2$  pod warunkiem  $\mathbf{AZ} \geq 0$  jest  $\chi_n^2$ , o ile zbiór  $\{\mathbf{Z} : \mathbf{AZ} \geq 0\}$  jest niepusty.

*Dowód.* Chcemy pokazać, że  $P(\|\mathbf{Z}\|^2 \leq a | \mathbf{AZ} \geq 0) = \chi_n^2(a)$ . W tym celu zapiszmy wektor  $\mathbf{Z}$  we współrzędnych biegunowych

$$\begin{aligned} Z_1 &= r \cos \phi_1 \cos \phi_2 \cos \phi_3 \dots \cos \phi_{n-1} \\ Z_2 &= r \sin \phi_1 \cos \phi_2 \cos \phi_3 \dots \cos \phi_{n-1} \\ Z_3 &= r \sin \phi_2 \cos \phi_3 \dots \cos \phi_{n-1} \\ Z_4 &= r \sin \phi_3 \dots \cos \phi_{n-1} \\ &\vdots \\ Z_n &= r \sin \phi_{n-1}, \end{aligned} \quad (3.5)$$

gdzie  $r \in (0, \infty), \phi_i \in [0, 2\pi), i = 1, 2, \dots, n - 1$ .

Wtedy

$$\|\mathbf{Z}\|^2 = r^2 \quad (3.6)$$

oraz

$$\mathbf{AZ} \geq 0 \iff \mathbf{A} \begin{bmatrix} \cos \phi_1 \cos \phi_2 \cos \phi_3 \dots \cos \phi_{n-1} \\ \sin \phi_1 \cos \phi_2 \cos \phi_3 \dots \cos \phi_{n-1} \\ \sin \phi_2 \cos \phi_3 \dots \cos \phi_{n-1} \\ \sin \phi_3 \dots \cos \phi_{n-1} \\ \vdots \\ \sin \phi_{n-1} \end{bmatrix} \geq 0. \quad (3.7)$$

Widzimy zatem, że wartość  $\|\mathbf{Z}\|^2$  zależy jedynie od wartości  $r$  natomiast warunek  $\mathbf{AZ} \geq 0$  dotyczy jedynie kąta, który jest niezależny od promienia  $r$ , a zatem  $P(\|\mathbf{Z}\|^2 \leq a | \mathbf{AZ} \geq 0) = \chi_n^2(a)$ .  $\square$

**Lemat 2.** Niech  $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n)^T \sim N(0, \mathbf{I})$  oraz niech  $\hat{\mathbf{Z}}$  będzie rzutem  $\mathbf{Z}$  na przestrzeń liniową  $S$  wymiaru  $d < n$ . Ponadto niech  $\mathbf{A}$  będzie rzeczywistą macierzą wymiaru  $m \times n$  taką, że każdy jej wiersz jest ortogonalny do przestrzeni  $S$ . Wtedy rozkładem warunkowym  $\|\hat{\mathbf{Z}}\|^2$  pod warunkiem  $\mathbf{AZ} \geq 0$  jest  $\chi_d^2$ , o ile zbiór  $\{\mathbf{AZ} \geq 0\}$  jest niepusty.

*Dowód.* Niech  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$  będą wzajemnie ortonormalnymi wektorami w  $\mathbb{R}^n$  takimi, że wektory  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$  rozpinają przestrzeń  $S$ . Wektor  $\mathbf{Z}$  możemy zapisać jako  $\mathbf{Z} = \sum_{i=1}^n a_i \mathbf{v}_i$ , gdzie  $a_i = \langle \mathbf{v}_i, \mathbf{z} \rangle$ . Stąd  $a_i, i = 1, 2, \dots, n$  są niezależnymi zmiennymi losowymi o standardowym rozkładzie normalnym oraz  $\hat{\mathbf{Z}} = \sum_{i=1}^d a_i \mathbf{v}_i$ . Wtedy  $\|\hat{\mathbf{Z}}\|^2 = a_1^2 + a_2^2 + \dots + a_d^2$  co ma gęstość  $\chi_d^2$ . Niech teraz  $\mathbf{V}$  oznacza macierz taką, której poszczególne kolumny są kolejno wektorami  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$  oraz niech  $\mathbf{a} = (a_1, a_2, \dots, a_n)^T$ . Macierz  $\mathbf{V}$  możemy zapisać jako  $\mathbf{V} = [\mathbf{V}_1 | \mathbf{V}_2]$ , gdzie  $\mathbf{V}_1$  jest macierzą wymiaru  $n \times d$ , oznaczmy też przez  $\mathbf{a}^1$  wektor  $(a_1, a_2, \dots, a_d)^T$  a przez  $\mathbf{a}^2$  wektor  $(a_{d+1}, \dots, a_n)^T$ . Wtedy  $\mathbf{Z} = \mathbf{V}\mathbf{a} = \mathbf{V}_1\mathbf{a}^1 + \mathbf{V}_2\mathbf{a}^2$  a warunek  $\mathbf{AZ} \geq 0$  możemy zapisać jako  $\mathbf{AV}_1\mathbf{a}^1 + \mathbf{AV}_2\mathbf{a}^2 \geq 0$ . Zauważmy, że z założeń oraz konstrukcji macierzy  $\mathbf{V}$  dostajemy, że  $\mathbf{AV}_1 = 0$  oraz  $\mathbf{a}^1, \mathbf{a}^2$  są niezależne. Zatem warunek  $\mathbf{AZ} \geq 0$  nie wpływa na gęstość  $\|\hat{\mathbf{Z}}\|^2$ .  $\square$

**Lemat 3.** Niech  $\mathbf{y} \in C_J$  dla pewnego zbioru  $J \subset L = \{1, 2, \dots, n-2\}$  oraz niech  $a, b \in \mathbb{R}$ . Wtedy  $\mathbf{y}' = \mathbf{y} + a\boldsymbol{\gamma}_{n-1} + b\boldsymbol{\gamma}_n \in C_J$  oraz wektory błędów  $\boldsymbol{\rho} = \mathbf{y} - \hat{\boldsymbol{\theta}}$  i  $\boldsymbol{\rho}' = \mathbf{y}' - \hat{\boldsymbol{\theta}}'$  są sobie równe.

*Dowód.* Jeśli  $\mathbf{y} \in C_J$  to  $\mathbf{y}$  możemy zapisać jako  $\mathbf{y} = \sum_{i \in J} b_i \boldsymbol{\beta}_i + \sum_{i \in L \setminus J} c_i \boldsymbol{\gamma}_i + d_1 \boldsymbol{\gamma}_{n-1} + d_2 \boldsymbol{\gamma}_n$ ,  $b_i > 0, c_i \geq 0, d_1, d_2 \in \mathbb{R}$ . Wtedy  $\mathbf{y}' = \sum_{i \in J} b_i \boldsymbol{\beta}_i + \sum_{i \in L \setminus J} c_i \boldsymbol{\gamma}_i + (d_1 + a) \boldsymbol{\gamma}_{n-1} + (d_2 + b) \boldsymbol{\gamma}_n$ ,  $b_i > 0, c_i \geq 0, d_1, d_2 \in \mathbb{R}$ . Oczywiście  $d_1 + a, d_2 + b \in \mathbb{R}$  zatem  $\mathbf{y}' \in C_J$ .

Wektor  $\boldsymbol{\rho}$  jest postaci  $\boldsymbol{\rho} = \sum_{i \in L \setminus J} c_i \boldsymbol{\gamma}_i$ . Z postaci wektora  $\mathbf{y}'$  widzimy jednak, że  $\boldsymbol{\rho}' = \sum_{i \in L \setminus J} c_i \boldsymbol{\gamma}_i = \boldsymbol{\rho}$ .  $\square$

**Lemat 4.** Wektory losowe  $\mathbf{y} - \hat{\boldsymbol{\theta}}$  i  $\hat{\boldsymbol{\theta}} - \hat{\mathbf{y}}$  są niezależne.

*Dowód.* Zauważmy, że  $\langle \mathbf{y} - \hat{\boldsymbol{\theta}}, \hat{\mathbf{y}} - \hat{\boldsymbol{\theta}} \rangle = \langle \mathbf{y} - \hat{\boldsymbol{\theta}}, \hat{\mathbf{y}} \rangle + \langle \mathbf{y} - \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\theta}} \rangle = 0 + 0 = 0$ .  $\square$

Dla ułatwienia rozważań założymy, że wariancja w zaproponowanym modelu  $\sigma^2$  jest znana.

Teraz możemy przystąpić do wyliczania rozkładu testu opartego o iloraz wiarygodności hipotezy

$$H_0: f(x) = ax + b \text{ vs. } H_1: f \in \mathcal{F} \quad (3.8)$$

gdzie  $\mathcal{F}$  jest klasą funkcji wypukłych.

Niech  $\hat{\mathbf{y}}$  oznacza estymator regresji liniowej, czyli rzut wektora danych  $\mathbf{y}$  na przestrzeń  $\text{span}\{\boldsymbol{\gamma}_{n-1}, \boldsymbol{\gamma}_n\}$ . Ponadto oznaczmy przez  $R_0 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$  oraz  $R_1 = \sum_{i=1}^n (y_i - \hat{\theta}_i)^2$ . Wtedy poszukiwany test przyjmuje postać

$$M = \frac{R_0 - R_1}{\sigma^2}. \quad (3.9)$$

W celu znalezienia rozkładu testu  $M$ , gdy prawdziwa jest hipoteza zerowa potrzebne będzie znalezienie wymiaru modelu i liczby stopni swobody błędu dla modelu regresji wypukłej. Z postaci rzutu  $\hat{\boldsymbol{\theta}}$  wektora  $\mathbf{y}$  można przypuszczać, że wymiar modelu wynosi  $n - d$  oraz liczba stopni swobody błędu wynosi  $d$ . Jednak zbiór  $J$  jest losowy, różne wartości wektora błędu  $\boldsymbol{\varepsilon}$  mogą umieścić wektor danych  $\mathbf{y}$  w różnych zbiorach  $C_J$ . Wyliczenie rozkładu testu zaczniemy w następujący sposób

$$P(M \leq a) = \sum_{J \in \mathcal{P}(L)} P(M \leq a, \mathbf{y} \in C_J) = \sum_{J \in \mathcal{P}(L)} P(M \leq a | \mathbf{y} \in C_J) P(\mathbf{y} \in C_J) \quad (3.10)$$

gdzie  $\mathcal{P}(L)$  oznacza zbiór potęgowy zbioru  $L$ .

Z lematu 3. możemy bez straty ogólności założyć, że  $f(x) = 0$  i  $\mathbf{y} = \boldsymbol{\varepsilon}$ . Dla dowolnej realizacji wektora  $\boldsymbol{\varepsilon} \in \mathbb{R}^n$  oznaczmy przez  $L \setminus J$  taki zbiór indeksów, że  $\boldsymbol{\varepsilon} \in C_{L \setminus J}$ . Niech  $\hat{\boldsymbol{\varepsilon}}$  będzie rzutem wektora  $\boldsymbol{\varepsilon}$  na przestrzeń  $S_{L \setminus J}$ . Zauważmy, że macierz  $\mathbf{A}_{L \setminus J}$  można zapisać jako  $[\mathbf{A}^1 | \mathbf{A}^2]$ , gdzie macierz  $\mathbf{A}^1$  jest wymiaru  $d \times n$ . Zatem kolumny macierzy  $\mathbf{A}^1$  rozpinają  $S_{L \setminus J}$ , natomiast kolumny macierzy  $\mathbf{A}^2$  są ortogonalne do przestrzeni  $S_{L \setminus J}$ . Dodatkowo, gdy  $\boldsymbol{\varepsilon} \in C_J$ , zachodzi  $\mathbf{A}^1 \boldsymbol{\varepsilon} \geq 0$  oraz  $\mathbf{A}^2 \boldsymbol{\varepsilon} \geq 0$ . Stąd na mocy lematu 2. dostajemy, że rozkładem warunkowym  $\frac{\|\hat{\boldsymbol{\varepsilon}}\|^2}{\sigma^2}$  przy zadanym  $J$  jest  $\chi_d^2$ . Jako że  $R_1 = \|\hat{\boldsymbol{\rho}}\|^2$  a przy założeniu prawdziwości hipotezy zerowej  $\|\hat{\boldsymbol{\rho}}\|^2$  jest równa  $\frac{\|\hat{\boldsymbol{\varepsilon}}\|^2}{\sigma^2}$  przy ustalonym zbiorze  $J$  możemy napisać następujący wniosek

**Wniosek 1.** *Jeśli hipoteza zerowa  $\boldsymbol{\theta} \in \text{span}\{\boldsymbol{\gamma}_{n-1}, \boldsymbol{\gamma}_n\}$  jest prawdziwa to rozkładem warunkowym  $\frac{R_1}{\sigma^2}$  przy ustalonym  $\mathbf{y} \in C_J$  jest  $\chi_d^2$ , gdzie  $d = |L \setminus J|$ .*

Zmienna losowa  $\frac{R_0}{\sigma^2}$  ma oczywiście rozkład  $\chi_{n-2}^2$ . Niech  $D$  będzie zmienną losową reprezentującą licznosc zbioru  $L \setminus J$ . Z Wniosku 1. mamy, że rozkładem warunkowym  $\frac{R_1}{\sigma^2}$  pod warunkiem  $D = d$  jest  $\chi_d^2$ . Z lematu 4. dostajemy zatem, że rozkładem warunkowym  $M$  jest  $\chi_{n-d-2}^2$  pod warunkiem  $D = d$ . Stąd możemy zapisać następujący wniosek

**Wniosek 2.** *Przy założeniu prawdziwości hipotezy zerowej postawionego problemu mamy*

$$P(M \leq a) = \sum_{d=0}^{n-2} P(\chi_{n-d-2}^2 \leq a) P(D = d), \quad (3.11)$$

gdzie  $\chi_0^2 \equiv 0$ .

Wartości prawdopodobieństw  $P(D = d), d = 0, 1, \dots, n - 2$  jest wyliczane na podstawie względnych objętości zbiorów  $C_J, J \in \mathcal{P}(L)$ . Prawdopodobieństwo, że  $\mathbf{y} \in C_J$ , gdy hipoteza zerowa jest prawdziwa, jest równoważne prawdopodobieństwu, że wektor losowy o  $n$ - wymiarowym standardowym rozkładzie normalny wpada do zbioru  $C_J$ .

*Kropka nie oznacza końca zdania.  
Ona daje możliwość coraz to lepszej kontynuacji.*

## Bibliografia

- [1] Fraser D.A.S., Massam H., *A Mixed Primal- Dual Bases Algorithm for Regression under Inequality Constraints. Application to Concave regression*, Scand J. Statist, **16** 65-74, 1989
- [2] Meyer Mary C., *A test for linear vs convex regression function using shape-restricted regression*, Stanford University, Technical Report No. 2001-20, sierpień 2001