

Poprzednim razem

- Statystyka (w dwóch znaczeniach)
- Estymatory punktowe

$$T_n(E(X)) \equiv \bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$T_n(\sigma(X)) \equiv S(X) = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2}$$
- Estymatory przedziałowe
 - Przedział ufności, poziom ufności
 - Dla E(X) przy znanym $\sigma(X)$

$$P\left(\bar{X} - \frac{1}{\sqrt{n}} \sigma(X) Z_{\frac{1-\gamma}{2}} \leq E(X) \leq \bar{X} + \frac{1}{\sqrt{n}} \sigma(X) Z_{\frac{1-\gamma}{2}}\right) = \gamma$$

1

Estymacja punktowa

- Nieobciążony
Dow: niech wszystkie X_i pochodzą z populacji o $E(X)=X_0$

$$E\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n} \sum_{i=1}^n E(x_i) = \frac{1}{n} \sum_{i=1}^n E(X) = \frac{1}{n} \sum_{i=1}^n X_0 = \frac{1}{n} n X_0 = X_0$$
- Zgodny

$$P(|\bar{X} - E(X)| > \varepsilon) = P(|\bar{X} - E(\bar{X})| > \varepsilon) \leq \frac{\text{var}(\bar{X})}{\varepsilon^2} = \frac{\text{var}(X)}{n\varepsilon^2}$$

$$\lim_{n \rightarrow \infty} P(|\bar{X} - E(X)| > \varepsilon) = \lim_{n \rightarrow \infty} \frac{\text{var}(X)}{n\varepsilon^2} = 0$$
- Najbardziej efektywny

$$\forall a > 0: P(|X - \mu| \geq a) \leq \frac{\sigma^2}{a^2}$$

Nierówność Czebyszewa-Bienymy

RPIS 2024/2025 2

2

Estymacja punktowa

$$\text{var}(\bar{X}) = E[(\bar{X} - E(\bar{X}))^2] = E[\bar{X}^2] - (E(\bar{X}))^2 = \frac{1}{n^2} E\left[\left(\sum_{i=1}^n x_i\right)^2\right] - X_0^2 =$$

$$= \frac{1}{n^2} E\left[\left(\sum_{i=1}^n x_i\right)\left(\sum_{j=1}^n x_j\right)\right] - X_0^2 = \frac{1}{n^2} E\left[\sum_{i=1, \dots, n, j=1}^n x_i x_j + \sum_{i=1}^n x_i^2\right] - X_0^2 =$$

$$= \frac{1}{n^2} E\left[\sum_{i=1, \dots, n, j=1}^n x_i x_j\right] + \frac{1}{n^2} E\left[\sum_{i=1}^n x_i^2\right] - X_0^2 = \frac{1}{n^2} \sum_{i=1, \dots, n, j=1}^n E(x_i)E(x_j) + \frac{n}{n^2} E(X^2) - X_0^2 =$$

$$= \frac{n(n-1)}{n^2} (E(X))^2 + \frac{1}{n} E(X^2) - X_0^2 = \frac{1}{n} [(n-1)(E(X))^2 + E(X^2) - n(E(X))^2] =$$

$$= \frac{1}{n} [E(X^2) - (E(X))^2] = \frac{1}{n} \text{var}(X)$$

- Czyli $\text{var}(\bar{X}) = \frac{1}{n} \text{var}(X) \Rightarrow \sigma(\bar{X}) = \frac{1}{\sqrt{n}} \sigma(X)$
- Ostatecznie przyjmujemy $X_0 = \bar{X} \pm \sigma(\bar{X})$

RPIS 2024/2025 3

3

Estymacja przedziałowa wariancji

- Twierdzenie: Statystyka

$$Y = \frac{(n-1)S^2(X)}{\sigma^2(X)} \quad S(X) = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2}$$
- ma rozkład χ^2 (Chi-kwadrat) o n-1 stopniach swobody. Ten rozkład nie jest symetryczny względem x=0.

RPIS 2024/2025 4

4

Estymacja przedziałowa wariancji

$$P\left(\chi_{\frac{1-\gamma}{2}}^{2, n-1} \leq Y \leq \chi_{\frac{1-\gamma}{2}}^{2, n-1}\right) = \gamma$$

$$P\left(\chi_{\frac{1-\gamma}{2}}^{2, n-1} \leq \frac{(n-1)S^2(X)}{\sigma^2(X)} \leq \chi_{\frac{1-\gamma}{2}}^{2, n-1}\right) = \gamma$$

$$P\left(\frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}} \leq \sigma^2(X) \leq \frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}}\right) = \gamma$$

$$P\left(\sqrt{\frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}}} \leq \sigma(X) \leq \sqrt{\frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}}}\right) = \gamma$$

$$T_n^L(\text{var}(X)) = \frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}} \quad T_n^p(\text{var}(X)) = \frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}}$$

RPIS 2024/2025 5

5

Estymacja przedziałowa odchylenia standardowego

Z estymacji przedziałowej wariancji otrzymujemy

$$P\left(\sqrt{\frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}}} \leq \sigma(X) \leq \sqrt{\frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}}}\right) = \gamma$$

Czyli przedział ufności dla odchylenia standardowego dany jest przez:

$$T_n^L(\sigma(X)) = \sqrt{\frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}}} \quad T_n^p(\sigma(X)) = \sqrt{\frac{(n-1)S^2(X)}{\chi_{\frac{1-\gamma}{2}}^{2, n-1}}}$$

RPIS 2024/2025 6

6

Znajdowanie estymatorów

– metoda największej wiarygodności

- Idea: To co zaobserwowaliśmy byłoby najbardziej prawdopodobne spośród wszystkich możliwych prób.
- Prawdopodobieństwo wylosowania próby x_1, x_2, \dots, x_n

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i; \theta_1, \theta_2, \dots, \theta_p) dx_i$$
- Funkcja (największej) wiarygodności**

$$\tilde{L}(\theta_1, \theta_2, \dots, \theta_p) = \prod_{i=1}^n f(x_i; \theta_1, \theta_2, \dots, \theta_p)$$
- Metoda największej wiarygodności** polega na szukaniu wartości parametrów $\theta_1, \theta_2, \dots, \theta_p$, dla których funkcja największej wiarygodności osiąga maksimum. Tak otrzymane wartości przyjmujemy za estymatory parametrów $\theta_1, \theta_2, \dots, \theta_p$.

RPIS 2024/2025 7

7

Metoda największej wiarygodności

- Uwaga: często wygodniej szukać maksimum funkcji

$$L(\theta_1, \theta_2, \dots, \theta_p) = \ln(\tilde{L}(\theta_1, \theta_2, \dots, \theta_p))$$

$$L(\theta_1, \theta_2, \dots, \theta_p)$$
 nazywamy **(logarytmiczną) funkcją (największej) wiarygodności**.
- Estymatory otrzymane metodą największej wiarygodności są zgodne, asymptotycznie nieobciążone i asymptotycznie najbardziej efektywne.
- Przykład: Estymatory parametrów rozkładu normalnego.

RPIS 2024/2025 8

8

Metoda największej wiarygodności - przykład: estymatory parametrów rozkładu normalnego

- Szukamy estymatorów parametrów θ_1 i θ_2 rozkładu $f_X(x) = \frac{1}{\sqrt{2\pi}\theta_2} e^{-\frac{(x-\theta_1)^2}{2\theta_2^2}}$
- Mając do dyspozycji próbę (x_1, x_2, \dots, x_n) budujemy logarytmiczną funkcję największej wiarygodności L

$$\tilde{L}(\theta_1, \theta_2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\theta_2} e^{-\frac{(x_i-\theta_1)^2}{2\theta_2^2}} = \frac{1}{(2\pi)^{\frac{n}{2}} \theta_2^n} e^{-\frac{1}{2\theta_2^2} \sum_{i=1}^n (x_i-\theta_1)^2}$$

$$L = \ln(\tilde{L}) = \ln((2\pi)^{-\frac{n}{2}} \theta_2^{-n} e^{-\frac{1}{2\theta_2^2} \sum_{i=1}^n (x_i-\theta_1)^2}) = -\frac{n}{2} \ln(2\pi) - n \ln(\theta_2) - \frac{1}{2\theta_2^2} \sum_{i=1}^n (x_i - \theta_1)^2$$
- I szukamy jej maksimum $\frac{\partial L}{\partial \theta_1} = 0$ $\frac{\partial L}{\partial \theta_2} = 0$

$$\frac{\partial L}{\partial \theta_1} = -\frac{1}{2\theta_2^2} \sum_{i=1}^n 2(x_i - \theta_1)(-1) = \frac{1}{\theta_2^2} \sum_{i=1}^n (x_i - \theta_1)$$

$$\frac{\partial L}{\partial \theta_2} = -n \frac{1}{\theta_2} - \frac{(-1)}{2\theta_2^3} 2\theta_2 \sum_{i=1}^n (x_i - \theta_1)^2 = -\frac{n}{\theta_2} + \frac{1}{\theta_2^3} \sum_{i=1}^n (x_i - \theta_1)^2$$

(proszę sprawdzić drugie pochodne)

RPIS 2024/2025 9

9

Metoda największej wiarygodności - przykład: estymatory parametrów rozkładu normalnego

- Porównujemy pochodne do zera, wiemy, że $\theta_2 \neq 0$

$$\frac{\partial L}{\partial \theta_1} = 0 \rightarrow \sum_{i=1}^n (x_i - \theta_1) = 0 \rightarrow (\sum_{i=1}^n x_i) - n\theta_1 = 0 \rightarrow \theta_1 = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial L}{\partial \theta_2} = 0 \rightarrow -\frac{n}{\theta_2} + \frac{1}{\theta_2^3} \sum_{i=1}^n (x_i - \theta_1)^2 = 0 \quad / \cdot \theta_2^3 \rightarrow \frac{1}{\theta_2^2} \sum_{i=1}^n (x_i - \theta_1)^2 = n \rightarrow \frac{1}{n} \sum_{i=1}^n (x_i - \theta_1)^2 = \theta_2^2$$
- Ostatecznie przyjmujemy

$$T_1(\theta_1) = \frac{1}{n} \sum_{i=1}^n x_i \equiv \bar{x}$$

$$T_1(\theta_2^2) = \frac{1}{n} \sum_{i=1}^n (x_i - \theta_1)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - T_1(\theta_1))^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$T_1(\theta_2) = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$
- Są to te same estymatory jakie wynikają z metody momentów; estymator wariancji jest obciążony.

RPIS 2024/2025 10

10

Znajdowanie estymatorów

– metoda najmniejszych kwadratów

- Idea: Szukamy estymatora parametrów $\theta_1, \theta_2, \dots, \theta_p$ występujących w równaniu $g(y_1, y_2, \dots, y_n; \theta_1, \theta_2, \dots, \theta_p) = 0$, gdzie (y_1, y_2, \dots, y_n) to wynik n-elementowej próby.
- Metoda najmniejszych kwadratów** polega na szukaniu wartości parametrów $\theta_1, \theta_2, \dots, \theta_p$, dla których funkcja

$$\sum_{i=1}^n w_i (g(y_i; \theta_1, \theta_2, \dots, \theta_p))^2$$
 osiąga minimum. Współczynniki liczbowe w_i określają wagę jaką przykładamy do kolejnych wartości y_i , mogą być to na przykład odwrotności kwadratów błędów pomiaru y_i . Tak otrzymane wartości przyjmujemy za estymatory parametrów $\theta_1, \theta_2, \dots, \theta_p$.

RPIS 2024/2025 11

11

Metoda najmniejszych kwadratów

- Np. zmierzono n razy zmienną Y uzyskując wartości y_1, y_2, \dots, y_n . Ile wynosi prawdziwa wartość zmiennej Y (oznaczamy ją przez θ)? Rozpatrzmy $g(Y, \theta) = y_i - \theta$
- Gdyby pomiar był dokładny to $g(Y, \theta) = 0$.

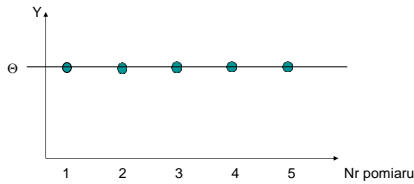
RPIS 2024/2025 12

12

Metoda najmniejszych kwadratów

- Np. zmierzono n razy zmienną Y uzyskując wartości y_1, y_2, \dots, y_n . Ile wynosi prawdziwa wartość zmiennej Y (oznaczamy ją przez θ)? Rozpatrzmy $g(Y, \theta) = y_i - \theta$

Gdyby pomiar był dokładny to $g(Y, \theta) = 0$.



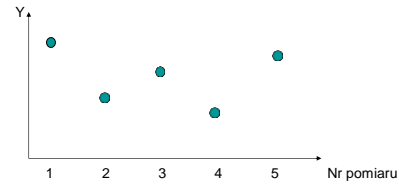
RPiS 2024/2025 13

13

Metoda najmniejszych kwadratów

- Np. zmierzono n razy zmienną Y uzyskując wartości y_1, y_2, \dots, y_n . Ile wynosi prawdziwa wartość zmiennej Y (oznaczamy ją przez θ)? Rozpatrzmy $g(Y, \theta) = y_i - \theta$

Gdyby pomiar był dokładny to $g(Y, \theta) = 0$.



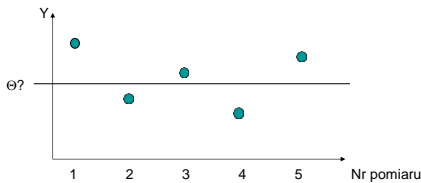
RPiS 2024/2025 14

14

Metoda najmniejszych kwadratów

- Np. zmierzono n razy zmienną Y uzyskując wartości y_1, y_2, \dots, y_n . Ile wynosi prawdziwa wartość zmiennej Y (oznaczamy ją przez θ)? Rozpatrzmy $g(Y, \theta) = y_i - \theta$

Gdyby pomiar był dokładny to $g(Y, \theta) = 0$.



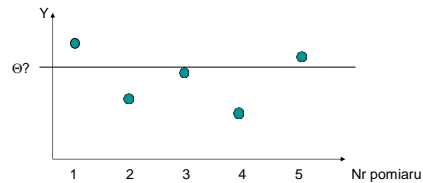
RPiS 2024/2025 15

15

Metoda najmniejszych kwadratów

- Np. zmierzono n razy zmienną Y uzyskując wartości y_1, y_2, \dots, y_n . Ile wynosi prawdziwa wartość zmiennej Y (oznaczamy ją przez θ)? Rozpatrzmy $g(Y, \theta) = y_i - \theta$

Gdyby pomiar był dokładny to $g(Y, \theta) = 0$.



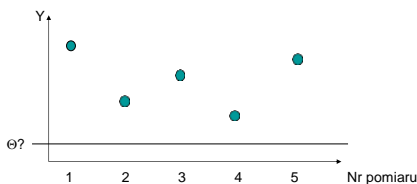
RPiS 2024/2025 16

16

Metoda najmniejszych kwadratów

- Np. zmierzono n razy zmienną Y uzyskując wartości y_1, y_2, \dots, y_n . Ile wynosi prawdziwa wartość zmiennej Y (oznaczamy ją przez θ)? Rozpatrzmy $g(Y, \theta) = y_i - \theta$

Gdyby pomiar był dokładny to $g(Y, \theta) = 0$.



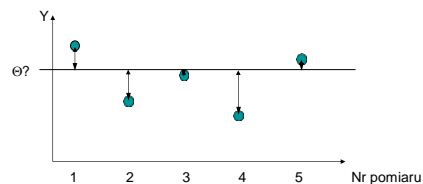
RPiS 2024/2025 17

17

Metoda najmniejszych kwadratów

- Np. zmierzono n razy zmienną Y uzyskując wartości y_1, y_2, \dots, y_n . Ile wynosi prawdziwa wartość zmiennej Y (oznaczamy ją przez θ)? Rozpatrzmy $g(Y, \theta) = y_i - \theta$

Gdyby pomiar był dokładny to $g(Y, \theta) = 0$.



RPiS 2024/2025 18

18

Metoda najmniejszych kwadratów

- Np. zmierzono n razy zmienną Y uzyskując wartości y_1, y_2, \dots, y_n . Ile wynosi prawdziwa wartość zmiennej Y (oznaczamy ją przez θ)? Rozpatrzmy

$$g(Y, \theta) = y_i - \theta$$

Gdyby pomiar był dokładny to $g(Y, \theta) = 0$.

W praktyce minimalizujemy

$$\sum_{i=1}^n (y_i - \theta)^2 \rightarrow \frac{d}{d\theta} \left(\sum_{i=1}^n (y_i - \theta)^2 \right) = -2 \sum_{i=1}^n (y_i - \theta)$$

$$-2 \sum_{i=1}^n (y_i - \theta) = 0 \quad /: (-2)$$

$$\sum_{i=1}^n y_i - \sum_{i=1}^n \theta = 0$$

$$\sum_{i=1}^n y_i = n\theta \rightarrow T_n(\theta) = \frac{1}{n} \sum_{i=1}^n y_i$$

RPIS 2024/2025 19

19

Metoda najmniejszych kwadratów

- W ogólnym przypadku estymatory pochodzące z metody najmniejszych kwadratów nie mają optymalnych własności (nawet asymptotycznie).

Wyjątki:

- Wyniki pomiaru y_i mają rozkład normalny i są nieskorelowane: wtedy metoda najmniejszych kwadratów jest równoważna metodzie największej wiarygodności.
- Szukane parametry są liniowymi współczynnikami w funkcji regresji.

RPIS 2024/2025 20

20

Funkcja regresji

- Funkcją regresji I rodzaju** zmiennej Y względem zmiennej X nazywamy warunkową wartość oczekiwaną $E(Y|X)$ traktowaną jako funkcję zmiennej X .

$$\left(\begin{array}{l} E(Y|X=x) = \int_{-\infty}^{\infty} y f_{Y|X}(y|x) dy \\ f_{Y|X}(y|x) = f_{X,Y}(x,y) / f_X(x) \end{array} \right)$$

- Tw: Wartość oczekiwana kwadratu odchyłań zmiennej losowej Y od dowolnej funkcji $u(X)$ jest minimalna gdy $u(X) = E(Y|X)$

$$E[(Y - u(X))^2] \geq E[(Y - E(Y|X))^2]$$

- Dow: $E[(Y - u(X))^2] = \iint dx dy f_{X,Y}(x,y) (y - u(x))^2 =$
 $= \iint dx dy f_X(x) f_{Y|X}(y|x) (y - u(x))^2 =$
 $= \int dx f_X(x) \int dy f_{Y|X}(y|x) (y - u(x))^2 =$

RPIS 2024/2025 21

21

Funkcja regresji

$$\begin{aligned} & \int dy f_{Y|X}(y|x) (y - u(x))^2 = \\ & = \int dy f_{Y|X}(y|x) ((y - E(Y|X) + E(Y|X) - u(x))^2) = \\ & = \int dy f_{Y|X}(y|x) [(y - E(Y|X))^2 + \\ & \quad + 2(y - E(Y|X))(E(Y|X) - u(x)) + (E(Y|X) - u(x))^2] = \\ & = \int dy f_{Y|X}(y|x) [(y - E(Y|X))^2 + (E(Y|X) - u(x))^2] + \\ & \quad + 2(E(Y|X) - u(x)) \int dy f_{Y|X}(y|x) (y - E(Y|X)) = \\ & = \int dy f_{Y|X}(y|x) [(y - E(Y|X))^2 + (E(Y|X) - u(x))^2] + \\ & \quad + 2(E(Y|X) - u(x)) \left[\int dy f_{Y|X}(y|x) y - E(Y|X) \int dy f_{Y|X}(y|x) \right] = \\ & = \int dy f_{Y|X}(y|x) [(y - E(Y|X))^2 + (E(Y|X) - u(x))^2] + \\ & \quad + 2(E(Y|X) - u(x)) [E(Y|X) - E(Y|X) \cdot 1] = \\ & = \int dy f_{Y|X}(y|x) [(y - E(Y|X))^2 + (E(Y|X) - u(x))^2] \end{aligned}$$

RPIS 2024/2025 22

22

Funkcja regresji

- Zatem

$$\begin{aligned} E[(Y - u(X))^2] &= \int dx f_X(x) \int dy f_{Y|X}(y|x) (y - u(x))^2 = \\ &= \int dx f_X(x) \int dy f_{Y|X}(y|x) [(y - E(Y|X))^2 + (E(Y|X) - u(x))^2] = \\ &= \iint dx dy f_X(x) f_{Y|X}(y|x) [(y - E(Y|X))^2] + \\ &+ \iint dx dy f_X(x) f_{Y|X}(y|x) [E(Y|X) - u(x)]^2 \geq \\ &\geq \iint dx dy f_X(x) f_{Y|X}(y|x) [(y - E(Y|X))^2] = \\ &= \iint dx dy f_{X,Y}(x,y) [(y - E(Y|X))^2] = E[(Y - E(Y|X))^2] \end{aligned}$$

- Pozwala to szukać funkcji $u(x)$, zależnej od pewnych parametrów, aproksymującej funkcję regresji I rodzaju $E(Y|X)$ przez minimalizowanie kwadratów odchyłań $u(x)$ od y : $\sum_{i=1}^n (y_i - u(x_i))^2$

RPIS 2024/2025 23

23

Liniowa funkcja regresji II rodzaju

- Liniowa funkcja regresji II rodzaju** przybliża liniowo $E(Y|X)$

$$E(Y|X) \approx u(x) = ax + b$$

- Regresja krzywoliniowa** to nieliniowa funkcja $u(x)$ przybliżająca $E(Y|X)$, której parametry znajdujemy korzystając z metody najmniejszych kwadratów.

- Przypadek I:

Parametry $\theta_1, \theta_2, \dots, \theta_p$ wchodzą do $u(x)$ liniowo, np. jako współczynniki wielomianów, wtedy dostajemy układ równań liniowych.

Estymatory $\theta_1, \theta_2, \dots, \theta_p$ pochodzące z metody najmniejszych kwadratów są nieobciążone, najbardziej efektywne i są liniowymi funkcjami y_i . Te własności nie zależą od rozkładu zmiennej Y i są spełnione nawet dla niewielkich prób.

- Przypadek II:

Parametry $u(x)$ wchodzą do niej nieliniowo, wtedy dostajemy układ równań nieliniowych.

RPIS 2024/2025 24

24