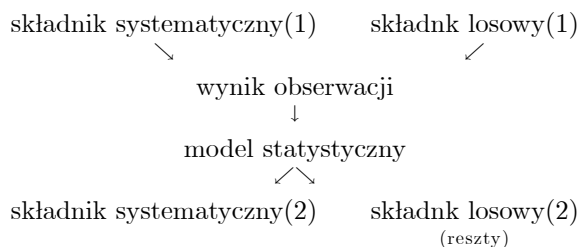


1 Regresja liniowa cz. I

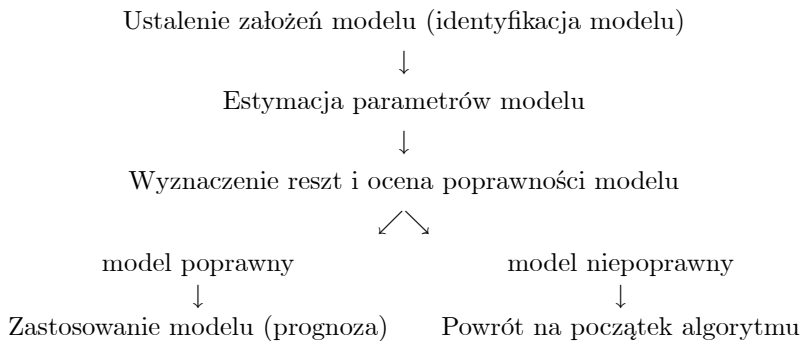
1.1 Model statystyczny

Model statystyczny to zbiór założeń. Wprowadzamy model, który możliwie najlepiej opisuje interesujący nas fragment rzeczywistość. Błędy modelu wynikają z nieuwzględnienia wszystkich czynników, błędnego uwzględnienia czynników oraz losowego zaburzenia. Model statystyczny ma opisywać składnik systematyczny, to co pozostaje to losowe reszty.



Zakładamy, że składnik losowy (1)(zaburzenie losowe, błąd losowy) jest wynikiem działania mało istotnych zdarzeń. Czego spodziewamy się od reszt? Powinny być niezależne o średniej 0, i stałej wariancji (równomiernie rozłożona chmurka punktów), często zakłada się, że mają rozkład normalny (pomocne założenie w technikach obliczeniowych).

Budowa modelu statystycznego:



1.2 Model regresji liniowej

I ETAP:Ustalenie założeń modelu

Ze względu na przyjęte założenie o liniowej zależności między zmiennymi nasz model jest postaci:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Y jest zmienną zależną (inaczej zmienną objaśnianą), której wartości chcemy wyjaśnić lub przewidzieć

X jest zmienną niezależną (inaczej zmienną objaśniającą) nazywaną też predyktorem, zakładamy, że zmienna ta nie jest zdegenerowana do stałej. W przeciwnym przypadku problemu regresji nie byłoby sensu rozważać.

Uwaga 1 Umawiamy się, że wartości X są ustalone (brak losowości).

ε jest błędem losowym (inaczej zakłóceniem, szumem), jedynym źródłem losowości

β_0 to wyraz wolny będący punktem przecięcia linii $Y = \beta_0 + \beta_1 X$ z osią rzędnych.

β_1 jest współczynnikiem kierunkowym, czyli tangensem kąta pod którym linia $Y = \beta_0 + \beta_1 X$ nachylona jest do osi odciętych.

$$\begin{aligned} Y &= \underbrace{\beta_0 + \beta_1 X}_{\text{składnik systematyczny}} + \underbrace{\varepsilon}_{\text{składnik losowy}} \\ y_i &= \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n \end{aligned}$$

Dodatkowe założenia modelu:

ε_i są nieskorelowane, o tym samym rozkładzie

$$\begin{aligned} E(\varepsilon) &= 0 \\ \text{Var}(\varepsilon) &= \sigma^2 \end{aligned}$$

Jeśli dodatkowo

$$\varepsilon \sim N(0, \sigma^2)$$

to mówimy o modelu normalnej regresji liniowej.

Przykład 1 RYSUNEK!!!

Uwaga 2 W modelu regresji obie zmienne są losowe, a w regresji tylko Y jest losowa.

Własności 1 W modelu normalnej regresji liniowej: $\varepsilon_i \perp \varepsilon_j, i \neq j$

Uwaga 3 Własności 2 Błędy (ε) są jednakowo rozproszone wokół linii regresji tzn. ich rozkład jest identyczny (w tym identyczna średnia i odchylenie)

Własności 3 1. $E(Y) = \beta_0 + \beta_1 X, (E(Y|X) = \beta_0 + \beta_1 X)$

2. $\text{Var}(Y) = \sigma^2, (\text{Var}(Y|X) = \sigma^2)$

3. $\text{Cov}(\varepsilon_j, y_i) = \begin{cases} 0, & i \neq j \\ \sigma^2, & i = j \end{cases}$

II Etap: Estymacja parametrów modelu (obecny wykład skoncentruje się wokół tego etapu)

Estymatory parametrów β_0 i β_1 oznaczamy odpowiednio $\hat{\beta}_0$ i $\hat{\beta}_1$. Na podstawie n -elementowej próby $(x_i, y_i), i = 1, \dots, n$ oszacowany model regresji jest postaci:

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X + e$$

gdzie e reprezentuje zaobserwowane błędy, czyli reszty z dopasowania linii regresji $\hat{\beta}_0 + \hat{\beta}_1 X$ do zbioru obserwacji Y czyli

$$e = Y - \hat{\beta}_0 - \hat{\beta}_1 X$$

Dla poszczególnych par punktów mamy związki

$$\begin{aligned} y_i &= \hat{\beta}_0 + \hat{\beta}_1 x_i + e_i, \quad i = 1, \dots, n \\ e_i &= y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \end{aligned}$$

Linia regresji:

$$\begin{aligned} \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X \\ \hat{y}_i &= \hat{\beta}_0 + \hat{\beta}_1 x_i, \quad i = 1, \dots, n \end{aligned}$$

Wartości \hat{y}_i nazywa się estymatorami y_i , wartościami teoretycznymi lub prognozowanymi wartościami Y .

-Metoda najmniejszych kwadratów:

$$\begin{aligned} SSE &= \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ SSE(\hat{\beta}_0, \hat{\beta}_1) &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\ &= \min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \end{aligned}$$

Szukamy $\hat{\beta}_0, \hat{\beta}_1$ realizującymi minimum:

$$\begin{aligned} \frac{\partial SSE(\hat{\beta}_0, \hat{\beta}_1)}{\partial \hat{\beta}_0} &= 0 \\ \frac{\partial SSE(\hat{\beta}_0, \hat{\beta}_1)}{\partial \hat{\beta}_1} &= 0 \end{aligned}$$

Ostatecznie

$$\begin{aligned} \hat{\beta}_1 &= \frac{SS_{xy}}{SS_x} \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \end{aligned}$$

gdzie

$$\begin{aligned} SS_x &= \sum_{i=1}^n (x_i - \bar{x})^2 \\ SS_{xy} &= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \end{aligned}$$

Uwaga 4 Terminologia: układ równań normalnych:

$$\begin{cases} 0 = \sum_{i=1}^n y_i - n\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i \\ 0 = \sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 \end{cases}$$

Własności 4 1. (\bar{x}, \bar{y}) należy do linii regresji

2. Układ równań normalnych jest równoważny z

$$\begin{cases} 0 = \sum_{i=1}^n e_i \\ 0 = \sum_{i=1}^n x_i e_i \end{cases}$$

3.

$$\hat{\beta}_1 = \sqrt{\frac{SS_y}{SS_x}} \hat{\rho}_{XY}$$

4.

$$\hat{y}_i - \bar{y} = \hat{\beta}_1 (x_i - \bar{x})$$

Uwaga 5 $\hat{\beta}_0$ i $\hat{\beta}_1$ są wartościami (relizacjami) estymatorów prawdziwych parametrów regresji.

1.3 Twierdzenie Gaussa-Markowa (przypadek jednowymiarowy)

Lemat 1

$$Var(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), Var(\hat{\beta}_1) = \frac{\sigma^2}{S_{XX}}$$

Twierdzenie 1 Gaussa i Markowa

W modelu regresji liniowej nieobciążonymi, liniowymi estymatorami parametrów β_0 i β_1 o najmniejszej wariancji w klasie estymatorów liniowych są estymatory wyznaczone metodą najmniejszych kwadratów (**BLUE** = best linear unbiased estimator).

Wniosek 1 $E(\hat{\beta}_0) = \beta_0$, $E(\hat{\beta}_1) = \beta_1$

Ćwiczenie 1 1. Zaproponować (sztuczny) model regresji liniowej.

2. Wygenerować pary punktów z tego modelu.

3. Przeprowadzić estymację parametrów

4. Narysować wygenerowane pary punktów wraz z linią regresji

5. Interpretacja

Solution 1 $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, $i = 1, \dots, 10$

1. Narzędzia>Analiza Danych>Generowanie Liczb Pseudolosowych (normalny)
2. Zaznacz obszar 1x2, naciśnij F2 wpisz =REGLINP(Zakres_Y;Zakres_X;1;1); CTRL+SHIFT+ENTER
3. Na wykresie rozrzutu punktów (x_i, y_i) zaznaczyć punkty>Prawy Przycisk Myszy>Dodaj linię trendu.
4. Jeśli wartość zmiennej objaśniającej X zwiększy się o jednostkę, to wartość zmiennej Y zmieni się o $\hat{\beta}_1$ (w zależności od znaku).

Ćwiczenie 2 Jak powyżej ale dla $n = 5000$. Wykorzystać generowanie liczb pseudolosowych dla X . Zbadaj wpływ wielkości odchylenia standardowego zmiennej ε na wyniki estymacji. Wnioski?

Ćwiczenie 3 Dla par $(x_i, y_i) = (x_i, x_{i+1})$, gdzie x_i jest i -tą obserwacją WIG20 przeprowadź estymację regresji (wyznacz estymatory, narysuj wykres rozrzutu wraz z linią regresji).

1.4 Badanie istotności parametrów

Twierdzenie 2 Jeśli ε mają rozkład normalny, to

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{SSE}{n-2},$$

gdzie $SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$, jest nieobciążonym estymatorem parametru¹ σ^2 .

Twierdzenie 3 Jeśli ε mają rozkład normalny, to

$$\frac{\hat{\beta}_0 - \beta_0}{SE(\hat{\beta}_0)} \sim t_{n-2}$$

$$\frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \sim t_{n-2}$$

gdzie² $SE(\hat{\beta}_0) = \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}$, $SE(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}$.

¹Dzielimy przez liczbę stopni swobody równą $n - 2$. Odejmujemy 2 bo stopnie swobody zostały związane przez estymatory współczynników.

² $SE(\hat{\beta}_0) = \sqrt{\widehat{Var}(\hat{\beta}_0)}$, $SE(\hat{\beta}_1) = \sqrt{\widehat{Var}(\hat{\beta}_1)}$

1.4.1 Przedziały ufności

Lemat 2 *Jeśli ε mają rozkład normalny, to*

1. $(1 - \alpha)$ 100% przedział ufności dla parametru β_0 jest postaci

$$\widehat{\beta}_0 \pm t_{n-2, \alpha/2} SE(\widehat{\beta}_0)$$

2. $(1 - \alpha)$ 100% przedział ufności dla parametru β_1 :

$$\widehat{\beta}_1 \pm t_{n-2, \alpha/2} SE(\widehat{\beta}_1)$$

Uwaga 6 *Wskazówka praktyczna: jeśli wartość zero nie znajduje się w przedziale ufności parametru β_1 to z prawdopodobieństwem $1 - \alpha$ zależność liniowa pomiędzy zmienną Y a zmienną X jest istotna (ściślej: współczynnik nachylenia linii regresji jest różny od zera).*

1.4.2 Testy istotności

Jeśli wartość zmiennej Y jest stała niezależnie od wartości X , to $\beta_1 = 0$.

Jeśli korelacja pomiędzy zmiennymi nie występuje to $\beta_1 = 0$.

Lemat 3 *Jeśli $\varepsilon \sim N(\cdot)$ oraz $c \in \mathbb{R}$ to*

$$H_0 : \beta_1 = c$$

$$H_0 : \beta_1 \neq c$$

$$\frac{\widehat{\beta}_1 - c}{SE(\widehat{\beta}_1)} \sim t_{n-2}$$

Jest to test dwustronny.

Wniosek 2 *Jeśli $\varepsilon \sim N(\cdot)$ to*

$$H_0 : \beta_1 = 0$$

$$H_0 : \beta_1 \neq 0$$

$$\frac{\widehat{\beta}_1}{SE(\widehat{\beta}_1)} \sim t_{n-2}$$

Jest to test dwustronny.

Uwaga 7 *Gdy współczynnik korelacji jest istotnie różny od zera lub gdy istnieje związek regresyjny pomiędzy zmiennymi to i tak nie możemy stwierdzić związku przyczynowo-skutowego między zmiennymi tzn. nie możemy stwierdzić, że jedna zmienna jest przyczyną drugiej.*

Ćwiczenie 4 Analiza regresji w oparciu o plik: *PilkaNozna.xls*.

1. Przeprowadzić estymację parametrów

2. Wyznaczyć błędy standardowe parametrów: $SE(\hat{\beta}_0) = \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}$,

$$SE(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}$$

3. Testuj:

$$H_0 : \beta_1 = 0$$

$$H_1 : \beta_1 \neq 0$$

Solution 2 Zaznacz obszar 2x2, naciśnij F2 wpisz =REGLINP(Zakres_Y;Zakres_X;1;1) oraz CTRL+SHIFT+ENTER. Otrzymasz

$$\begin{matrix} \hat{\beta}_1 & \hat{\beta}_0 \\ SE(\hat{\beta}_1) & SE(\hat{\beta}_0) \end{matrix}$$

$$t_{stat} = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1)}$$

$$p - value = ROZKLAD.T(\text{Moduł.Liczby}(t_{stat}), n - 2, 2)$$

$$U = ROZKLAD.T.ODW(\alpha/2; n - 2)$$

$$\text{Zbiór krytyczny} = (-\infty, -U) \cup (U, \infty)$$

α to prawdopodobieństwo zbioru krytycznego.

$$ROZKLAD.T(x, \nu, 2) = P(|X| > x), \text{gdzie } X \sim t(\nu)$$

$$ROZKLAD.T.ODW(p; \nu) = U, \text{gdzie } P(X > U) = p$$

Ćwiczenie 5 Dla danych z WIG20 oraz ich logarytmów sprawdzić normalność.

QQ-Plot: Chcemy ocenić czy x_1, \dots, x_n są realizacją rozkładu normalnego. Uporządkować próbkę x_1, \dots, x_n od najmniejszej do największej otrzymując $x_{(1)}, \dots, x_{(n)}$.

Zaznaczyć na układzie współrzędnym pary $(x_{(j)}, z_j)$, dla $j = 1, \dots, n$, gdzie z_j wyznaczamy z tożsamości:

$$\frac{j - 0.5}{n} = N\left(\frac{z_j - \hat{\mu}_n}{\hat{\sigma}_n}\right)$$

gdzie N jest dystrybuantą standardowego rozkładu normalnego $N(0, 1)$, $\hat{\mu}_n$, $\hat{\sigma}_n$ są odpowiednio oszacowaniem wartości oczekiwanej (μ) i odchylenia standardowego (σ) z próby. Należy narysować „prostą regresji liniowej” dla $(x_{(j)}, z_j)$. Zaznaczone punkty nie powinny „zbyt odstępować od równania prostej”.

Ćwiczenie 6 Zastosować regresję liniową do estymacji parametrów modelu dwumianowego

$$\begin{aligned}S(k) &= S(k-1)(1 + \xi_k U + (1 - \xi_k) D) \\S(0) &= S_0\end{aligned}$$

$$\xi_k = \begin{cases} 1 & p \\ 0 & 1-p \end{cases}$$

$$\begin{aligned}Y_k &= \frac{S(k)}{S(k-1)} = 1 + \xi_k U + (1 - \xi_k) D \\y_k &= Y_k - 1 = \xi_k U + (1 - \xi_k) D \\&= D + \xi_k (U - D)\end{aligned}$$