



Elementy statystyki

STA - Wykład 3

dr hab. Waldemar Wołyński
Wydział Matematyki i Informatyki
Uniwersytet im. Adama Mickiewicza



Prosta regresja liniowa

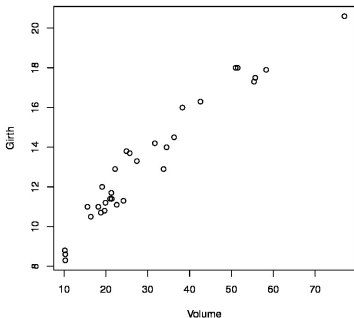
Model: prosta regresja liniowa



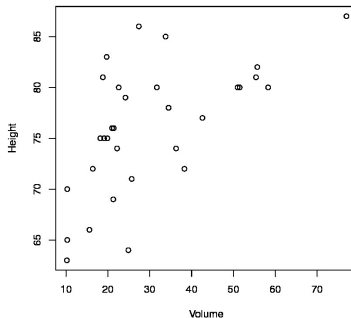
Modelujemy wyniki doświadczenia w którym dokonujemy n -niezależnych obserwacji zmiennej (zależnej) Y (typu ilościowego ciągłego) oraz niezależnej zmiennej X (typu ilościowego), związanych zależnością liniową na losowo wybranych z populacji jednostkach eksperymentalnych.

Przykład Staramy się oszacować objętość tarcicy (zmienna zależna) jaką można pozyskać z drzew czarnej wiśni. Objętość tarcicy zależy od wielkości drzewa opisanego przez jego średnicę i wysokość (zmienne niezależne). Dane dla 31 drzew zawarte są w ramce "trees" dostępnej w programie **R**. Budowę modelu rozpoczynamy od wykonania diagramów korelacyjnych.

Model: prosta regresja liniowa



objętość – średnica



objętość – wysokość

Jako zmienną niezależną przyjmujemy średnicę (dodatkowo przyjmujemy liniową zależność pomiędzy objętością a średnicą).

Model: prosta regresja liniowa



Wybieramy model prostej regresji liniowej postaci:

$$Y_i = a + bx_i + \varepsilon_i, \quad i = 1, \dots, n,$$

gdzie

Y_i – i -ta obserwacja zmiennej zależnej Y ,
 x_i – i -ta wartość zmiennej niezależnej X ,
 a, b – parametry liniowej funkcji regresji,
 ε_i – błędy (reszty).

Uwaga:

W rozważanym przykładzie zmienna niezależna to średnica drzewa, zmienna zależna to objętość drzewa.



O błędach zakładamy, że:

- ▶ są niezależne (dokładnie: są niezależnymi zmiennymi losowymi),
- ▶ mają wartość oczekiwaną równą zero (nie ma błędu systematycznego), tzn.

$$E(\varepsilon_i) = 0, \quad i = 1, \dots, n,$$

- ▶ mają jednakową, stałą i niezerową wariancję, tzn.

$$\text{Var}(\varepsilon_i) = \sigma^2, \quad i = 1, \dots, n.$$

Uwaga:

Model prostej regresji liniowej ma trzy parametry: a , b i σ^2 .

Estymacja parametrów a i b liniowej funkcji regresji



Do estymacji parametrów a i b używamy metody **najmniejszych kwadratów**, polegającej na minimalizacji sumy kwadratów błędów, tzn.

$$S(a, b) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - a - bx_i)^2.$$

Fakt

Estymatorami najmniejszych kwadratów (ENK) parametrów a i b liniowej funkcji regresji są statystyki:

$$\hat{a} = \bar{Y} - \hat{b}\bar{x}, \quad \hat{b} = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$



Twierdzenie

W modelu prostej regresji liniowej, statystyki \hat{a} i \hat{b} są nieobciążonymi estymatorami parametrów a i b .

Ponadto, statystyka

$$\hat{\sigma}^2 = S^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{a} - \hat{b}x_i)^2$$

jest nieobciążonym estymatorem parametru σ^2 .



Twierdzenie

Przy dodatkowym założeniu normalności rozkładu błędów, w modelu prostej regresji liniowej, statystyki

$$\hat{a} \text{ i } S^2 \text{ oraz } \hat{b} \text{ i } S^2$$

są niezależnymi zmiennymi losowymi.

Ponadto

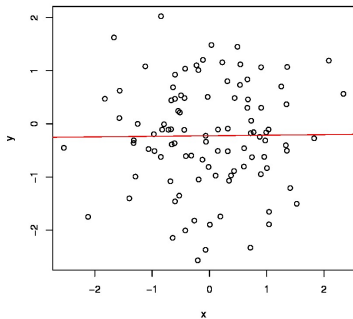
$$\hat{a} \sim N \left(a, \sigma^2 \left(\frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} \right) \right),$$

$$\hat{b} \sim N \left(b, \sigma^2 \left(\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \right)$$

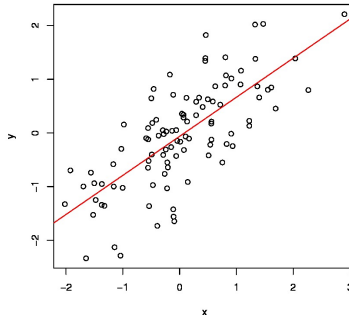
oraz

$$\frac{(n-2)S^2}{\sigma^2} \sim \chi^2(n-2).$$

Wpływ zmiennej niezależnej X na zmienną zależną Y



Brak istotnego wpływu
 $b = 0$



Istotny wpływ
 $b \neq 0$



Dopasowanie modelu

W modelu prostej regresji liniowej prawdziwa jest następująca zależność:

$$SST = SSR + SSE$$

gdzie

$$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2, \quad SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y}_i)^2,$$

$$SSE = \sum_{i=1}^n (\hat{Y}_i - Y_i)^2, \quad \hat{Y}_i = \hat{a} + \hat{b}x_i.$$

Liczbową miarą dopasowania prostej regresji do danych empirycznych jest **współczynnik determinacji** (podawany w %)

$$R^2 = 1 - \frac{SSE}{SST}, \quad 0 \leq R^2 \leq 1.$$



Model prostej regresji liniowej wykorzystujemy często do wyznaczania prognozy wartości zmiennej zależnej, przy ustalonej wartości zmiennej niezależnej. Niech x_p oznacza wartość zmiennej niezależnej X dla której uzyskać chcemy prognozę zmiennej zależnej Y równą Y_p .

Przyjmujemy:

$$\hat{Y}_p = \hat{a} + \hat{b}x_p.$$

Jako ocenę jakości prognozy przyjmujemy oszacowanie odchylenia standardowego prognozy (średni błąd prognozy) postaci:

$$S_p = S \sqrt{\frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$



Wielokrotna regresja liniowa



Model: wielokrotna (wieloraka) regresja liniowa

Model ten jest rozwinięciem modelu prostej regresji liniowej na przypadek wielu zmiennych niezależnych.

$$Y_i = a_0 + a_1 x_{i1} + a_2 x_{i2} + \dots + a_m x_{im} + \varepsilon_i, \quad i = 1, \dots, n,$$

gdzie

Y_i – i -ta obserwacja zmiennej zależnej Y ,

$x_{i1}, x_{i2}, \dots, x_{im}$ – i -te wartości zmiennych niezależnych

x_1, x_2, \dots, x_m ,

a_0, a_1, \dots, a_m – parametry liniowej funkcji regresji,

ε_i – błędy (reszty).

Uwagi:

Założenia dotyczące błędów są identyczne jak w modelu prostej regresji liniowej.

Model wielokrotnej regresji liniowej ma $m + 2$ parametry:

a_0, a_1, \dots, a_m i σ^2 .

Zapis macierzowy modelu



$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & X_{11} & \dots & X_{1m} \\ 1 & X_{21} & \dots & X_{2m} \\ \vdots & \vdots & & \vdots \\ 1 & X_{n1} & \dots & X_{nm} \end{bmatrix}, \quad \mathbf{a} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_m \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

Model wielokrotnej regresji liniowej: $\mathbf{Y} = \mathbf{X}\mathbf{a} + \boldsymbol{\varepsilon}$.

Uwaga: W modelu wielokrotnej regresji liniowej zakładamy, że $n > m$ oraz $\text{rzęd}(\mathbf{X}) = m + 1$.



Twierdzenie

W modelu wielokrotnej regresji liniowej, statystyka

$$\hat{\mathbf{a}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

jest nieobciążonym estymatorem parametru $\hat{\mathbf{a}}$.

Ponadto, statystyka

$$\hat{\sigma}^2 = S^2 = \frac{1}{n - m - 1} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2,$$

gdzie

$$\hat{Y}_i = \hat{a}_0 + \hat{a}_1 x_{i1} + \cdots + \hat{a}_m x_{im}, \quad i = 1, \dots, n,$$

jest nieobciążonym estymatorem parametru σ^2 .



W modelu wielokrotnej regresji liniowej (analogicznie jak w modelu prostej regresji liniowej) prawdziwa jest następująca zależność:

$$SST = SSR + SSE.$$

Liczbową miarą dopasowania funkcji regresji (hiperpłaszczyzny regresji) do danych empirycznych jest **poprawiony współczynnik determinacji** (podawany w %)

$$R_{pop}^2 = 1 - \frac{SSE/(n - m - 1)}{SST/(n - 1)}, \quad 0 \leq R_{pop}^2 \leq 1.$$



Niech $\mathbf{X}_p = [1, x_1^p, \dots, x_m^p]'$ oznacza wektor wartości zmiennych objaśniających X_1, X_2, \dots, X_m dla którego uzyskać chcemy prognozę zmiennej objaśnianej Y .

Przyjmujemy:

$$\hat{Y}_p = \mathbf{X}_p' \hat{\mathbf{a}}.$$

Jako ocenę jakości prognozy przyjmujemy oszacowanie odchylenia standardowego prognozy (średni błąd prognozy) postaci:

$$S_p = S \sqrt{\mathbf{X}_p' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}_p}.$$



Funkcje związane z modelem regresji liniowej (prostej i wielokrotnej):

lm – funkcja podstawowa,

summary – wartości estymatorów parametrów modelu regresji, wartość współczynnika determinacji, itp.

predict – wartości prognozowane,

abline – wykres prostej regresji.