

Analiza korelacji i regresji

NIEZALEŻNOŚĆ ZMIENNYCH

Zależność stochastyczna dwóch zmiennych losowych polega na tym, że zmiana jednej z nich zmienia rozkład prawdopodobieństwa drugiej zmiennej

Cecha X jest stochastycznie niezależna od cechy Y gdy:

$$\bar{x}_1 = \bar{x}_2 = \dots \bar{x}_k \quad \text{oraz} \quad s_1^2(x) = s_2^2(x) = \dots = s_k^2(x)$$

Cecha Y jest stochastycznie niezależna od cechy X gdy:

$$\bar{y}_1 = \bar{y}_2 = \dots \bar{y}_r \quad \text{oraz} \quad s_1^2(y) = s_2^2(y) = \dots = s_r^2(y)$$

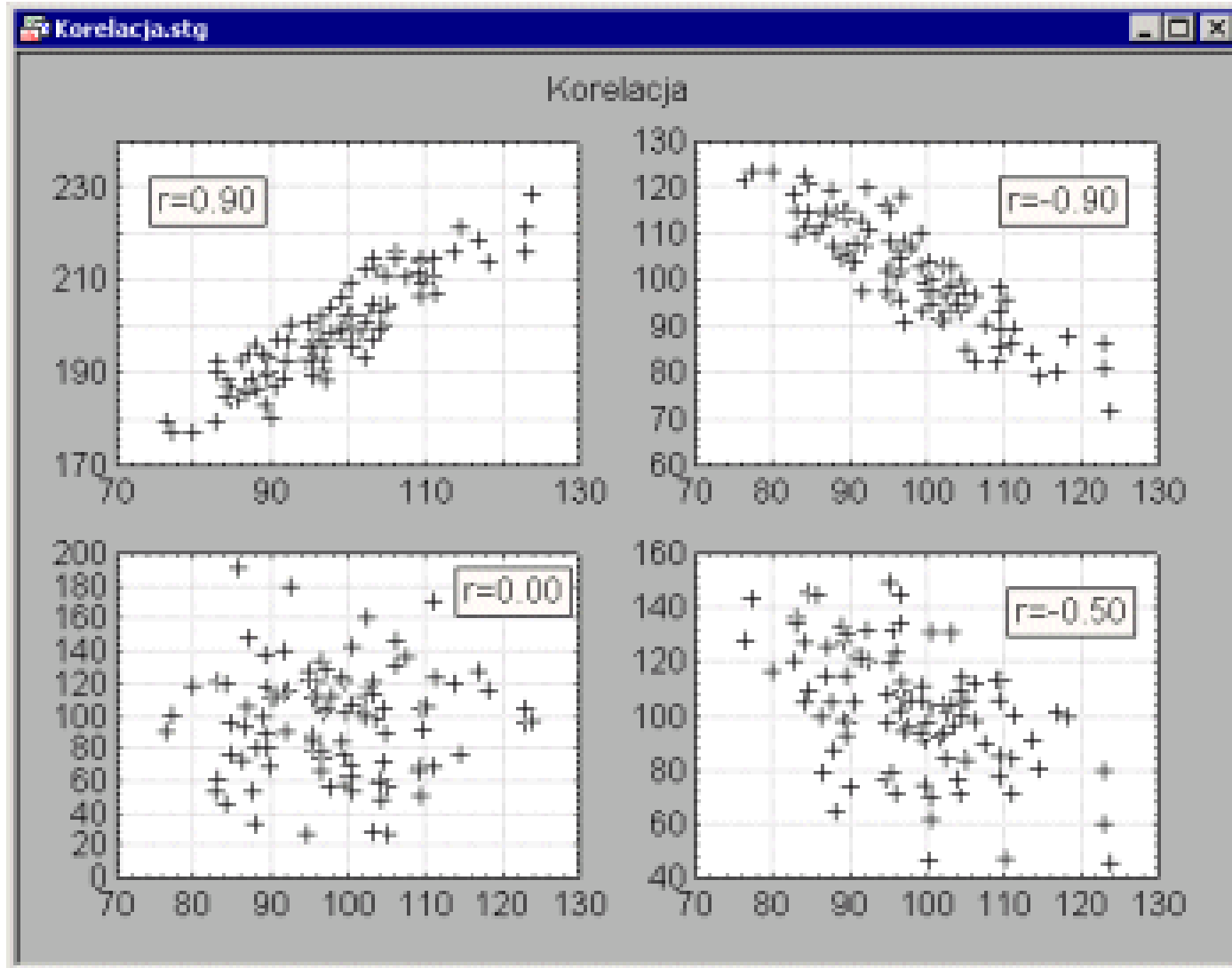
Zależność korelacyjna dwóch zmiennych losowych polega na tym, że zmiana jednej z nich pociąga za sobą zmianę średnich wartości drugiej cechy

Cecha X jest korelacyjnie niezależna od cechy Y gdy: $\bar{x}_1 = \bar{x}_2 = \dots \bar{x}_k$

Cecha Y jest korelacyjnie niezależna od cechy X gdy: $\bar{y}_1 = \bar{y}_2 = \dots \bar{y}_r$

Zależność funkcyjna dwóch zmiennych losowych polega na tym, że zmiana wartości jednej z nich pociąga za sobą ściśle określoną zmianę wartości drugiej cechy

Korelacyjny diagram rozrzutu

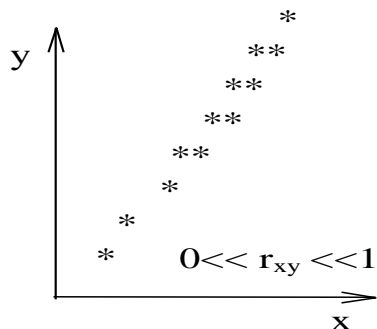


Korelacyjny diagram rozrzutu

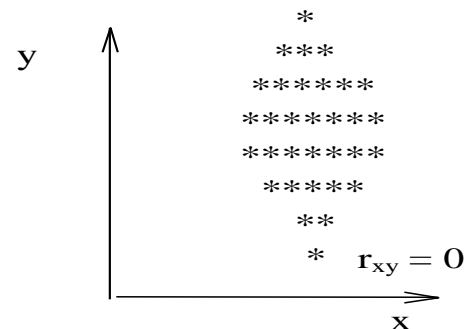
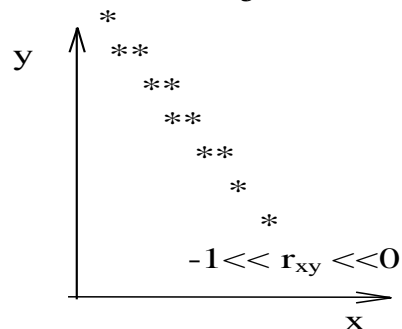
Korelacja liniowa

Brak korelacji

Silna dodatnia

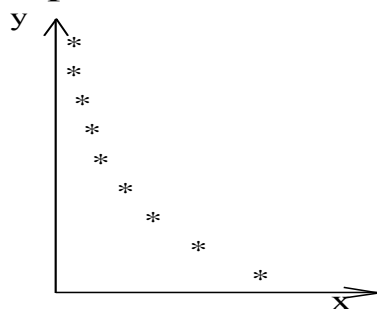


Silna ujemna

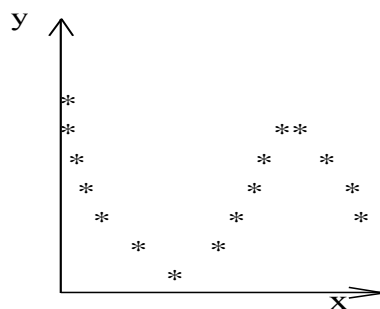


Korelacja nieliniowa

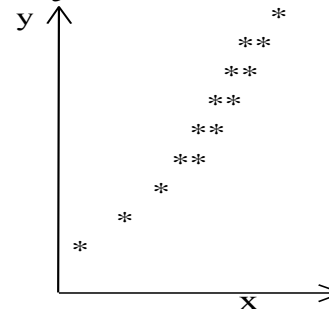
hiperbola



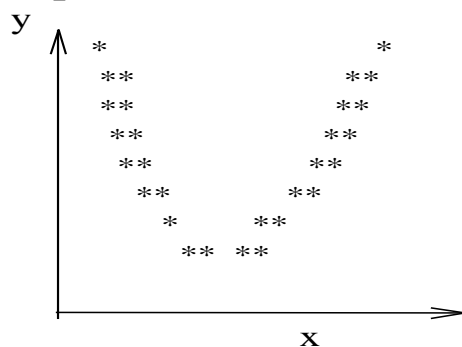
wielomian (3°)



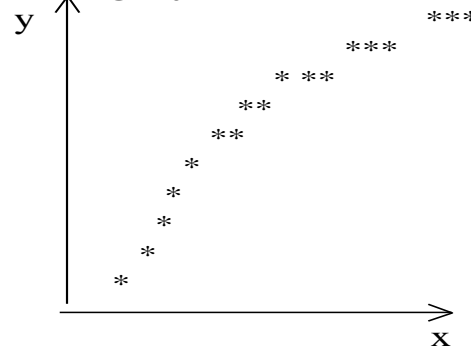
wykładnicza



parabola (2°)



logarytmiczna



Współczynnik korelacji liniowej Pearsona

- $r_{xy} = r_{yx}$ - mierzy siłę i kierunek związku prosto-linowego pomiędzy dwoma cechami
- *współczynnik korelacji* - standaryzacja kowariancji
- *kowariancja* - średnia arytmetyczna iloczynu odchyleń zmiennych X i Y od wartości średnich arytmetycznych
- *związek liniowy* - jednostkowym przyrostom przyczyny (zmiennej niezależnej) towarzyszy stały przyrost skutku (zmiennej zależnej)
- kowariancja informuje o:
 - $\text{cov}(X,Y)=0$ braku zależności korelacyjnej
 - $\text{cov}(X,Y)<0$ ujemnej zależności korelacyjnej
 - $\text{cov}(X,Y)>0$ dodatniej zależności korelacyjnej
- Unormowanym miernikiem natężenia siły związku dwóch zmiennych jest *współczynnik korelacji liniowej Pearsona*

$$r_{xy} = r_{yx} = \frac{\text{cov}(X,Y)}{s(x)s(y)}$$

$$r \in \langle -1;1 \rangle$$

$$r = \frac{\text{cov}(X,Y)}{s(x) \cdot s(y)} = \frac{\sum_j (x_j - \bar{x})(y_j - \bar{y})}{n \cdot s(x) \cdot s(y)} = \frac{\frac{1}{n} \sum_j x_j y_j - \bar{x} \cdot \bar{y}}{s(x) \cdot s(y)}$$

Współczynnik korelacji liniowej Pearsona

$$r = 1$$

- korelacja dodatnia - związek funkcyjny

$$0 < r < 1$$

- korelacja dodatnia niedoskonała

$$r = 0$$

- brak korelacji prostoliniowej

$$-1 < r < 0$$

- korelacja ujemna niedoskonała

$$r = -1$$

- korelacja ujemna - związek funkcyjny

$$r_{xy} = r_{yx} = \frac{\text{cov}(X, Y)}{s(x)s(y)}$$

$$r \in \langle -1 ; 1 \rangle$$

Funkcja regresji

- Analityczny wyraz przyporządkowania średnich wartości zmiennej objaśnianej konkretnym wartościom zmiennych objaśniających.
- Modele regresji I-szego rodzaju:

$$Y = \beta_{1y} \cdot X + \beta_{0y} + \xi$$

$$X = \beta_{1x} \cdot Y + \beta_{0x} + Z$$

- **Składnik losowy ξ (Z)** ma swoje źródło w mechanizmach losowych
 - badana zależność ma charakter losowy,
 - nie jest możliwe uwzględnienie wszystkich czynników wpływających na poziom badanego zjawiska,
 - nie znamy dokładnej postaci analitycznej funkcji
 - pomiar danych liczbowych nie jest dokładny

Funkcja regresji

- Estymatorem funkcji regresji I rodzaju jest funkcja regresji II rodzaju :

$$\hat{y} = b_{1y} \cdot x + b_{0y}$$

$$\hat{x} = b_{1x} \cdot y + b_{0x}$$

- która spełnia warunek (funkcja kryterium KMNK) :

$$E\{[Y - \hat{Y}]^2\} = E\{[Y - (b_{1y} \cdot x + b_{0y})]^2\} \rightarrow \min$$

- Warunek ten oznacza, iż
$$E\{[X - \hat{X}]^2\} = E\{[X - (b_{1x} \cdot y + b_{0x})]^2\} \rightarrow \min$$
 - średnie odchylenie zmiennej losowej od prostej regresji jest równe zero,
 - prosta regresji przechodzi przez punkt o współrzędnych ($x = E(X)$, $y = E(Y)$)

- Model regresji II rodzaju zapiszemy jako :

$$y = b_{1y} \cdot x + b_{0y} + e$$

$$x = b_{1x} \cdot y + b_{0x} + z$$

- Wektory reszt określone są następująco :

$$e_j = y_j - \hat{y}_j$$

$$z_j = x_j - \hat{x}_j$$

Etapy estymacji parametrów modelu regresji

1. Dobór zmiennych do modelu:
2. Określenie postaci analitycznej modelu
3. Wybór metody szacunku parametrów strukturalnych modelu
4. Oszacowanie parametrów strukturalnych modelu.
5. Oszacowanie parametrów struktury stochastycznej modelu
6. Ocena jakości modelu.

Założenia KMNK

1. Postać modelu jest liniowa
2. Zmienne objaśniające są nielosowe
3. Składowik losowy ma nadzieję matematyczną równą zero $E(\xi) = 0$ i stałą wariancję $D^2(\xi) = \text{const.}$

Przez stałość wariancji rozumie się, że nie zależy ona od kolejnych realizacji zmiennych objaśniających modelu.

Homoscedastyczność - **heteroscedastyczność**

Scedastyczność - poważny problem, gdyż błędy losowe szacowane na podstawie różnych wartości zmiennej niezależnej X mogą się zmieniać. Gdy warunek ten nie jest spełniony należy zrezygnować z KMNK na rzecz metod alternatywnych.

4. Realizacje zmiennych objaśniających są niezależne, co sprawia, że ciąg $\{\xi_j\}$ jest ciągiem niezależnych zmiennych losowych
5. Składowik losowy ξ nie jest skorelowany ze zmiennymi objaśniającymi
6. Błędy losowe charakteryzują się rozkładem normalnym $N[0, D^2(\xi)]$

Szacowanie parametrów strukturalnych modelu

$$\hat{y}_j = b_{1y} \cdot x_j + b_{0y}$$

$$\begin{cases} \sum_{j=1}^n Y_j = b_{1y} \sum_{j=1}^n X_j + nb_{0y} \\ \sum_{j=1}^n X_j Y_j = b_{1y} \sum_{j=1}^n X_j^2 + b_{0y} \sum_{j=1}^n X_j \end{cases}$$

$$b_{1y} = \frac{\sum_{j=1}^n (X_j - \bar{X})(Y_j - \bar{Y})}{\sum_{j=1}^n (X_j - \bar{X})^2} = \frac{\sum_{j=1}^n X_j Y_j - n\bar{X} \bar{Y}}{\sum_{j=1}^n X_j^2 - n\bar{X}^2}$$

$$b_{0y} = \bar{Y} - b_{1y} \bar{X}$$

Szacowanie parametrów strukturalnych modelu

$$\hat{x}_j = b_{1x} \cdot y_j + b_{0x}$$

$$\begin{cases} \sum_{j=1}^n X_j = b_{1x} \sum_{j=1}^n Y_j + nb_{0x} \\ \sum_{j=1}^n X_j Y_j = b_{1x} \sum_{j=1}^n Y_j^2 + b_{0x} \sum_{j=1}^n Y_j \end{cases}$$

$$b_{1x} = \frac{\sum_{j=1}^n (X_j - \bar{X})(Y_j - \bar{Y})}{\sum_{j=1}^n (Y_j - \bar{Y})^2} = \frac{\sum_{j=1}^n X_j Y_j - n\bar{X} \bar{Y}}{\sum_{j=1}^n Y_j^2 - n\bar{Y}^2}$$

$$b_{0x} = \bar{X} - b_{1y} \bar{Y}$$

Szacowanie parametrów strukturalnych modelu

- b_{1y} - współczynnik regresji informuje o ile przeciętnie zmieni się wartość zmiennej zależnej gdy zmienna niezależna wzrośnie o jednostkę
- b_{1y} - tangens kąta nachylenia linii regresji do osi OX
- b_{0y} - wyraz wolny - rzędna punktu przecięcia linii regresji z osią OY

Relacje między współczynnikami regresji i współczynnikiem korelacji

$$r = \pm \sqrt{b_{1y} \cdot b_{1x}}$$

$$b_{1x} = r \frac{s(x)}{s(y)}$$

$$b_{1y} = r \frac{s(y)}{s(x)}$$

Miary jakości dopasowania funkcji regresji

1. **Wariancja składnika resztowego (nieobciążony estymator wariancji składnika losowego $D^2(\xi)$)**

$$s^2(e_j) = s_{y(x)}^2 = \frac{\sum_{j=1}^n e_j^2}{n-k} = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{n-k}$$

dla funkcji liniowej dwóch zmiennych :

$$s^2(e_j) = s_{y(x)}^2 = \frac{\sum_{j=1}^n (y_j - (b_{1y}x_j + b_{0y}))^2}{n-k} \approx s^2(y) \cdot (1 - r^2)$$

Miary jakości dopasowania funkcji regresji

2. Zgodnie z równością wariancyjną - **wariancja ogólna** $s^2(y)$ równa jest sumie **wariancji resztowej** *niewyjaśnionej regresją* i **wariancji teoretycznej** *wyjaśnionej regresją* :

$$s^2(y) = s^2(\hat{y}) + s^2(e)$$

$$\sum_{j=1}^n (y_j - \bar{y})^2 = \sum_{j=1}^n (\hat{y}_j - \bar{y})^2 + \sum_{j=1}^n (y_j - \hat{y}_j)^2$$

Miary jakości dopasowania funkcji regresji

3. Współczynnik zbieżności resztowej - *indeterminacji, alienacji, nieokresloności* φ^2

$$\varphi^2 = \frac{s^2(e)}{s^2(y)} = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

$$0 \leq \varphi^2 \leq 1$$

Miary jakości dopasowania funkcji regresji

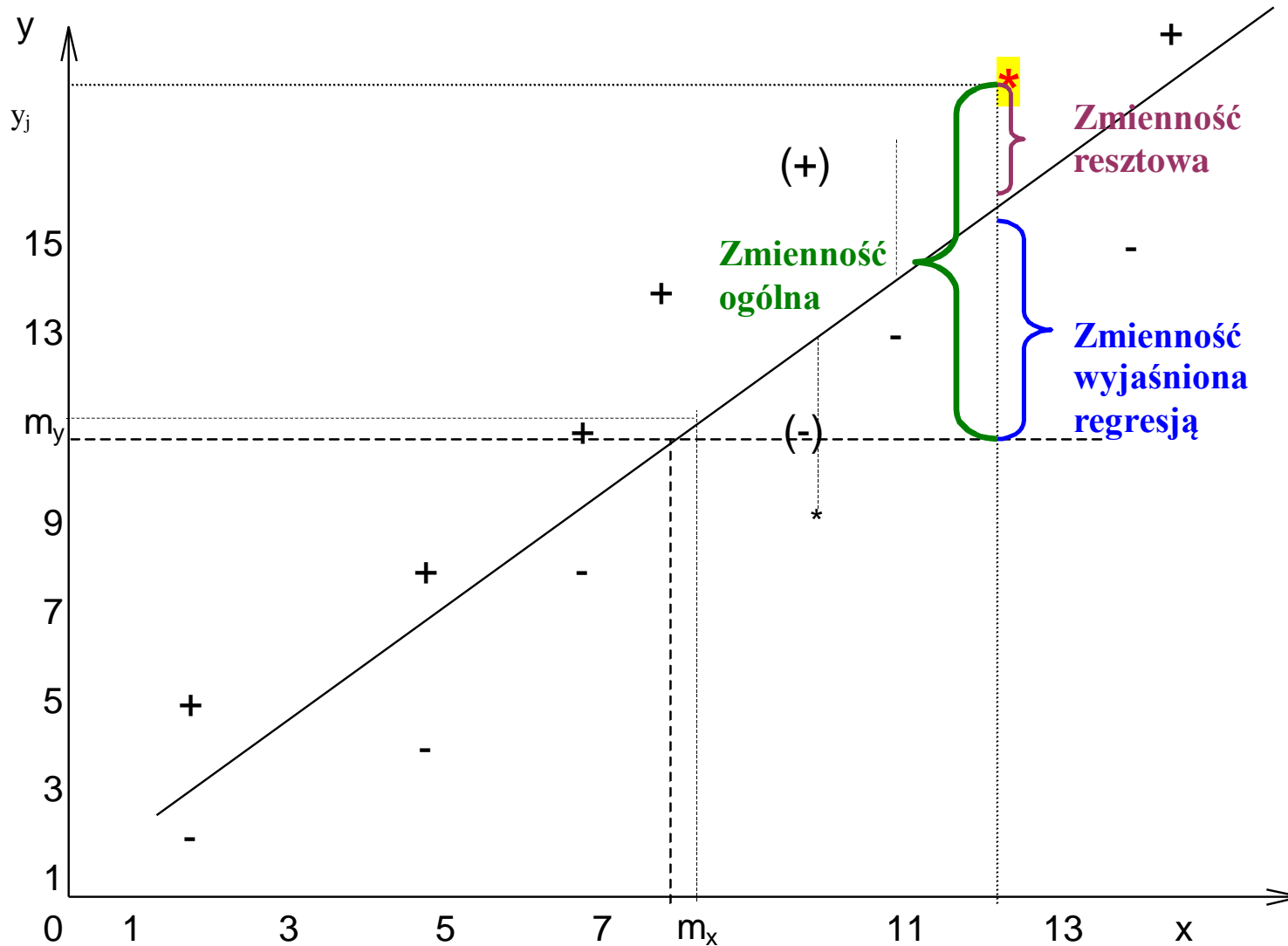
4. Współczynnik determinacji - *określoności*

$$R^2 = \frac{s^2(\hat{y})}{s^2(y)} = \frac{\sum_{j=1}^n (\hat{y}_j - \bar{y})^2}{\sum_{j=1}^n (y_j - \bar{y})^2} = 1 - \frac{s^2(e)}{s^2(y)} = 1 - \varphi^2$$

$$0 \leq R^2 \leq 1$$

$$R^2 + \varphi^2 = 1$$

Dekompozycja wariancji zmiennej zależnej



Błędy standardowe szacunku parametrów funkcji regresji

Współczynników regresji:

$$s(b_{1y}) = \frac{s_{y(x)}}{\sqrt{\sum x_j^2 - n\bar{x}^2}} = \frac{s_{y(x)}}{\sqrt{\sum (x_j - \bar{x})^2}}$$

Wyrazów wolnych:

$$s(b_{0y}) = \sqrt{\frac{s_{y(x)}^2 \sum x_j^2}{n \sum (x_j - \bar{x})^2}} = \sqrt{\frac{s_{y(x)}^2 \sum x_j^2}{n(\sum x_j^2 - n\bar{x}^2)}}$$

$$s(b_{1x}) = \frac{s_{x(y)}}{\sqrt{\sum y_j^2 - n\bar{y}^2}} = \frac{s_{x(y)}}{\sqrt{\sum (y_j - \bar{y})^2}}$$

$$s(b_{0x}) = \sqrt{\frac{s_{x(y)}^2 \sum y_j^2}{n \sum (y_j - \bar{y})^2}} = \sqrt{\frac{s_{x(y)}^2 \sum y_j^2}{n(\sum y_j^2 - n\bar{y}^2)}}$$

Pełny zapis modelu regresji

$$y_j = b_{1y}x + b_{0y} + e_j$$

$$(s_{(b_{1y})})(s_{(b_{0y})})(s_e)$$

- parametry strukturalne
- parametry struktury stochastycznej

$$x_j = b_{1x}y_i + b_{0x} + z_j$$

$$(s_{(b_{1x})})(s_{(b_{0x})})(s_z)$$

Analiza regresji w *Excelu*

1) *Funkcja statystyczna Excel* → REGLINP

Podaje wartości:

- b_{yi} - ocen parametrów strukturalnych modelu regresji,
- $s(b_{yi})$ - błędy standardowe ocen parametrów strukturalnych,
- R^2 - wartość współczynnika determinacji,
- Sey - standardowy błąd oceny y - błąd standardowy szacunku funkcji regresji - odchylenie standardowe składnika resztowego,
- F - statystykę F ,
- df - liczbę stopni swobody,
- $ssreg$ - sumę kwadratów wyjaśnioną regresją,
- $Ssresid$ - sumę kwadratów nie wyjaśnioną regresją.

W wyniku zastosowania funkcji REGLINP otrzymujemy:

1294.053321	25855.31826
37.75274933	2503.665761
0.646298476	21711.11342
1174.916964	643
5.53823E+11	3.03092E+11

co należy zapisać następująco:

$$y_i = 1294,05 \cdot x_i + 25855,52 + e_i$$

(37,75) (2503,67) (21711,11)

$$R^2 = 0,646$$

Analiza regresji w *Excelu*

2) Wybór metody: *Narzędzia* → *Analiza danych* Przygotowanie Dodatków

The screenshot shows the Microsoft Excel interface with the 'Dodatki' (Add-ins) dialog box open. The dialog box lists the following add-ins:

- ☒ Analysis ToolPak
- ☒ Analysis ToolPak - VBA
- ☒ Asystent internetowy VBA
- ☐ Dodatek Solver
- ☐ Kreator odnośników
- ☐ Kreator sum warunkowych
- ☐ Narzędzia do waluty euro

The background spreadsheet shows data for 'Powiat chodzieski' with columns for various categories and numerical values. The status bar at the bottom indicates 'Gotowy' and 'NUM'.

Analiza regresji w *Excelu*

2) Narzędzia → Analiza danych

Microsoft Excel - Powiat1a

zakres danych wejściowych

A1 Dane z mikrospisu NSP95

Analiza danych

Narzędzia analizy:

- Wyglądanie wykładowe
- Test F: z dwiema próbami dla wariancji
- Analiza Fouriera
- Histogram
- Średnia ruchoma
- Generowanie liczb pseudolosowych
- Ranga i percentyl
- Regresja**
- Próbkowanie
- Test t: par skojarzonych z dwiema próbami dla średniej

OK

Anuluj

Pomoc

chodziecki / czarnkowsko-trzcianecki / gnieznienski / gostynski

	A18	A26	A27
2	2	9	8
2	2	5	4
2	2	9	8
1	6	0	0
3	3	4	3
4	4	8	7
1	1	0	0
2	2	9	8
2	0	0	0
1	2	0	0
3	3	0	0

Analiza regresji w *Excelu*

2) Okno dialogowe *Regresja* — wypełnienie panelu

Microsoft Excel - Powiat1a

Plik Edycja Widok Wstaw Format Narzędzia Dane Okno Pomoc

zakres danych wejściowych

J14 Dane z mikrospisu NSP95

Regresja

Wejście

Zakres wejściowy Y: \$J\$14:\$J\$100

Zakres wejściowy X: \$E\$14:\$G\$100

☐ Tytuły ☐ Stała wynosi Zero

☐ Poziom ufności: 95 %

OK

Anuluj

Pomoc

Opcje wyjścia

☒ Zakres wyjściowy: \$L\$24

☐ Nowy arkusz:

☐ Nowy skoroszyt

Składniki resztowe

☐ Składniki resztowe ☐ Rozkład reszt

☐ Std. składniki resztowe ☐ Rozkład linii dopasowanej

Rozkład normalny

☐ Rozkład prawdopodobieństwa normalnego

chodziecki / czarnkowsko-trzcianecki / gnieźnieński / gostyniński

Wskaź NUM

11:41

Analiza regresji w *Excelu*

2) Okno dialogowe Regresja — wypełnienie panelu

Dane wejściowe - Dane wejściowe powinny być numeryczne w postaci wektorów kolumnowych o tej samej liczbie wierszy.

Zakres wejściowy Y - Wprowadź odwołanie do zakresu wejściowego zmiennej zależnej. Zakres musi składać się z pojedynczej kolumny danych.

Zakres wejściowy X - Wprowadź odwołanie do zakresu wejściowego zmiennych niezależnych. Program Microsoft Excel porządkuje zmienne niezależne z tego zakresu rosnąco od lewej do prawej. Maksymalna liczba zmiennych niezależnych jest równa 16.

Tytuły - Zaznacz to pole wyboru, jeżeli pierwszy wiersz albo pierwsza kolumna zakresów wejściowych zawiera etykiety. Wyczyść je, jeżeli zakres wejściowy nie zawiera etykiet; program *Excel* generuje odpowiednie etykiety danych w tabeli wyjściowej.

Poziom ufności - Zaznacz to pole, aby uwzględnić dodatkowy poziom w podsumowującej tabeli wyjściowej. W polu wprowadź wartość poziomu ufności, który będzie stosowany oprócz domyślnego poziomu 95 procent.

Stała wynosi zero - Zaznacz to pole, jeżeli chcesz wymusić, aby linia regresji przechodziła przez początek układu współrzędnych.

Zakres wyjściowy - Wprowadź odwołanie do lewej górnej komórki tabeli wyników.

Jeśli spełnione są założenia (3-6), to b_i są najlepszymi estymatorami parametrów β_i

Nowy arkusz - Kliknij, aby wstawić w bieżącym skoroszycie nowy arkusz i wkleić do niego wyniki, rozpoczynając od komórki A1. Nazwę nowego arkusza wpisz w polu.

Nowy skoroszyt - Kliknij, aby utworzyć nowy skoroszyt i wkleić wyniki do nowego arkusza w nowym skoroszycie.

Składniki resztkowe - Zaznacz to pole, aby uwzględnić składniki resztkowe w tabeli wyjściowej składników resztkowych.

Standaryzowane składniki resztkowe - Zaznacz to pole, aby uwzględnić standaryzowane składniki resztkowe w tabeli wyjściowej składników resztkowych.

Rozkład reszt - Zaznacz to pole wyboru, aby wygenerować wykres każdej zmiennej niezależnej w funkcji składnika resztkowego.

Rozkład linii dopasowanej (prognozowanych) - Zaznacz to pole, aby wygenerować wykres wartości teoretycznych

Rozkład prawdopodobieństwa normalnego - Zaznacz to pole, aby wygenerować wykres rozkładu prawdopodobieństwa normalnego.

Analiza regresji w *Excelu*

2) Wyniki analizy regresji

PODSUMOWANIE - WYJŚCIE								
<i>Statystyki regresji</i>								
Wielokrotność R	0,735786196							
R kwadrat	0,541381326							
Dopasowany R kwadrat	0,524804747							
Błąd standardowy	2,275715012							
Obserwacje	87							
ANALIZA WARIANCJI								
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Istotność F</i>			
Regresja	3	507,417426	169,1391	32,659413	4,86E-14			
Resztkowy	83	429,8469418	5,178879					
Razem	86	937,2643678						
	<i>Współczynniki</i>	<i>Błąd standardowy</i>	<i>t Stat</i>	<i>Wartość-p</i>	<i>Dolne 95%</i>	<i>Górne 95%</i>	<i>Dolne 95,0%</i>	<i>Górne 95,0%</i>
Przecięcie	-1,123984475	0,55911616	-2,01029	0,0476482	-2,236044	-0,011925	-2,236044331	-0,01192462
Zmienna X 1	2,282432125	0,25559222	8,929975	9,017E-14	1,7740693	2,7907949	1,774069342	2,790794909
Zmienna X 2	-0,443186869	0,726899877	-0,60969	0,5437297	-1,888962	1,0025881	-1,888961882	1,002588145
Zmienna X 3	1,186514967	0,320388406	3,703364	0,0003821	0,5492751	1,8237548	0,549275142	1,823754793