

ALGORYTMICZNA I STATYSTYCZNA ANALIZA DANYCH

13/11/2014

WFAiS UJ, Informatyka Stosowana
II stopień studiów

2

Regresja liniowa

Korelacja

Modelowanie

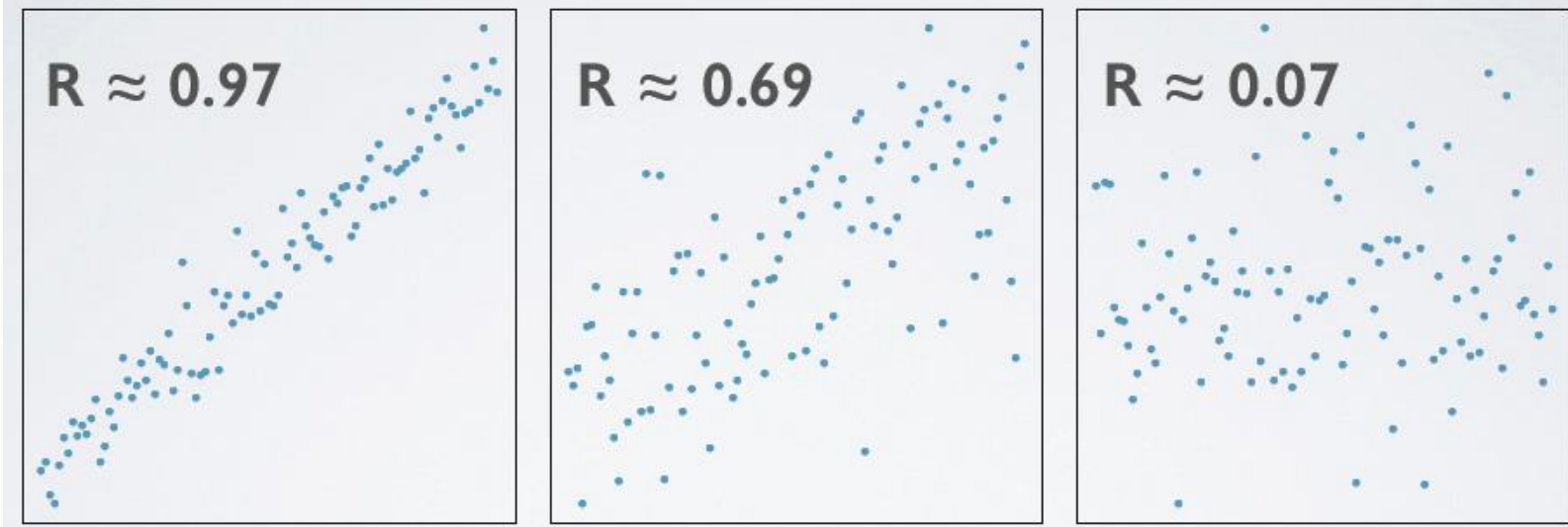
Analiza modelu

Wnioskowanie

Korelacja: własności

4

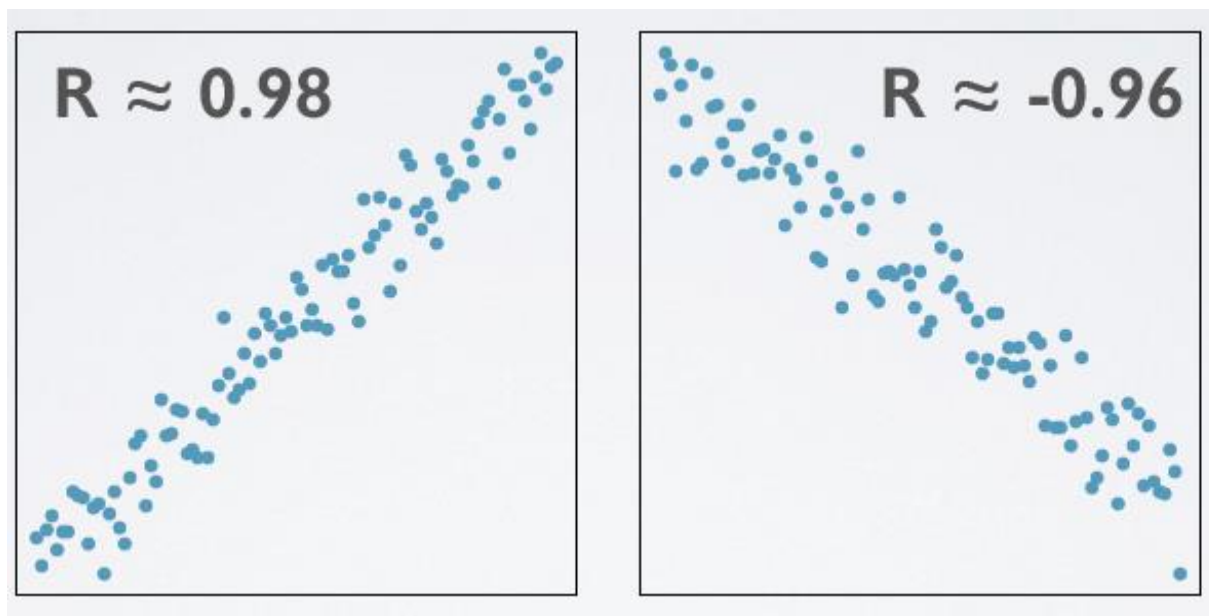
- Wielkość (absolutna) współczynnika korelacji mierzy siłę liniowej asocjacji pomiędzy dwoma zmiennymi numerycznymi



Korelacja: własności

5

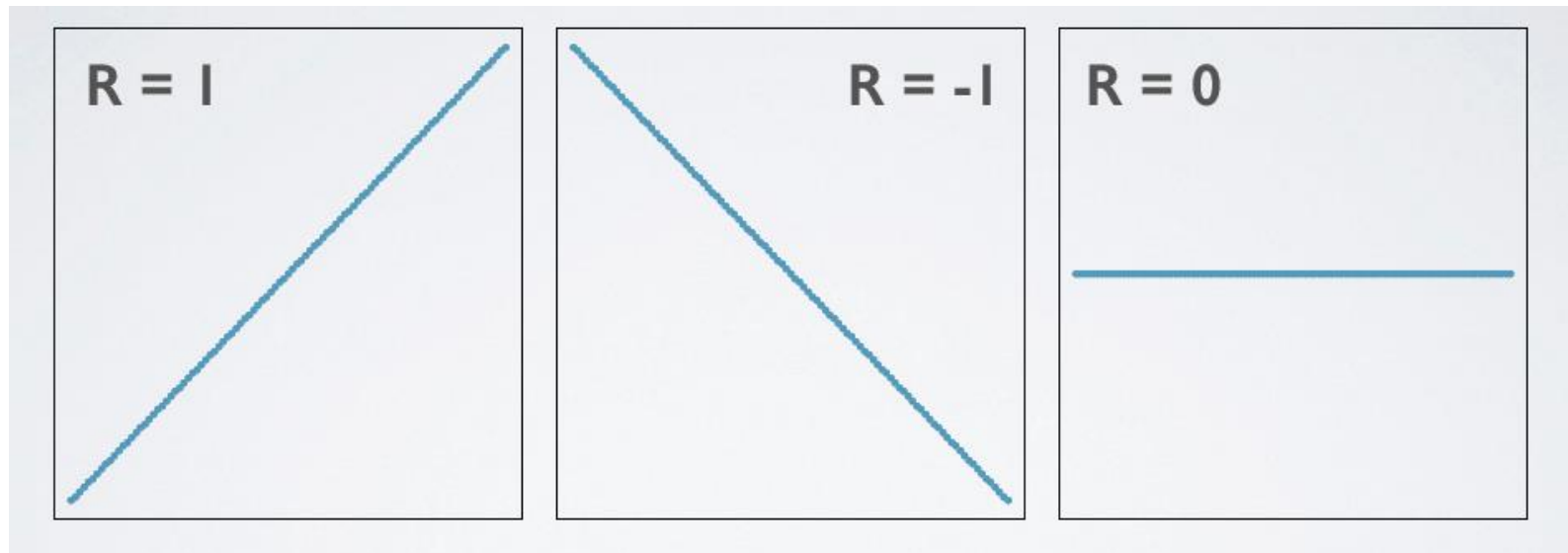
- Znak korelacji wskazuje na kierunek asocjacji



Korelacja: własności

6

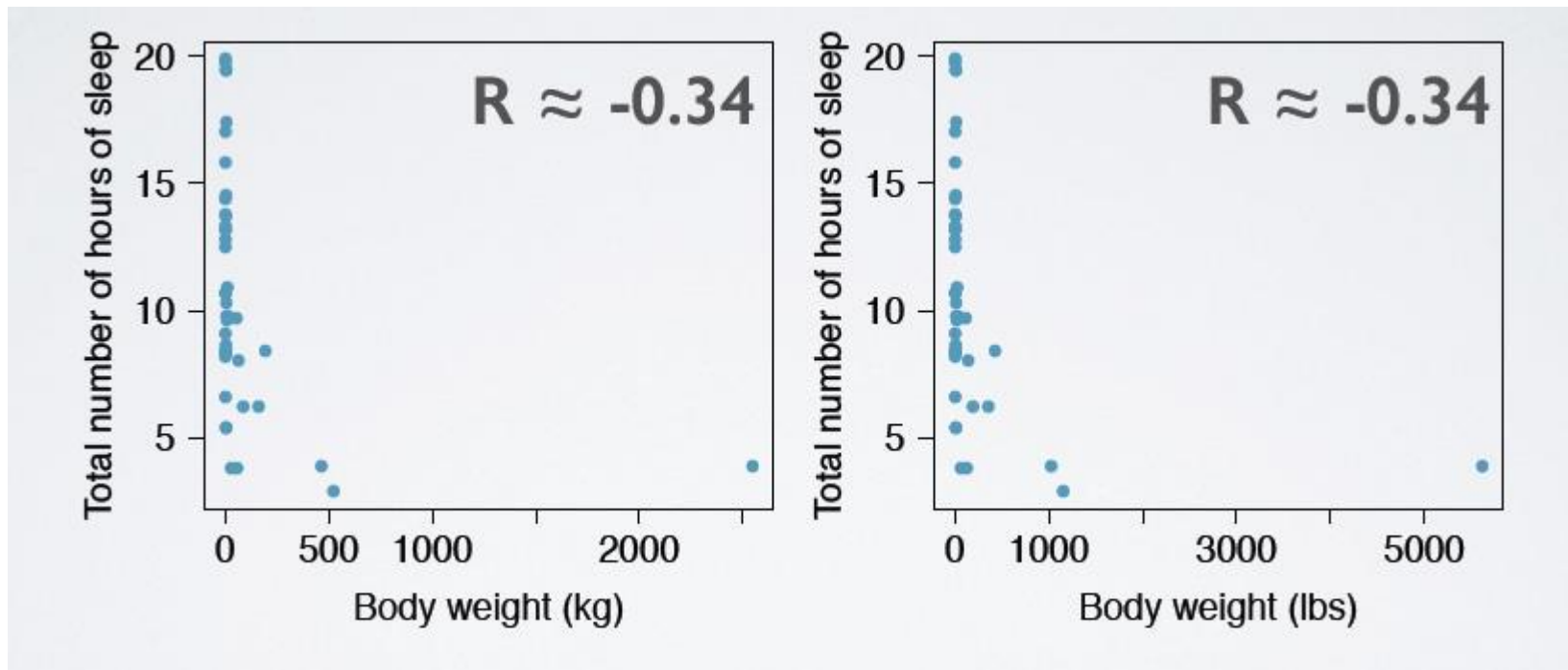
- Z definicji, wartość R jest pomiędzy -1 (idealna negatywna asocjacja) do 1 (idealna pozytywna asocjacja)
- $R = 0$ wskazuje na brak zależności pomiędzy zmiennymi



Korelacja: własności

7

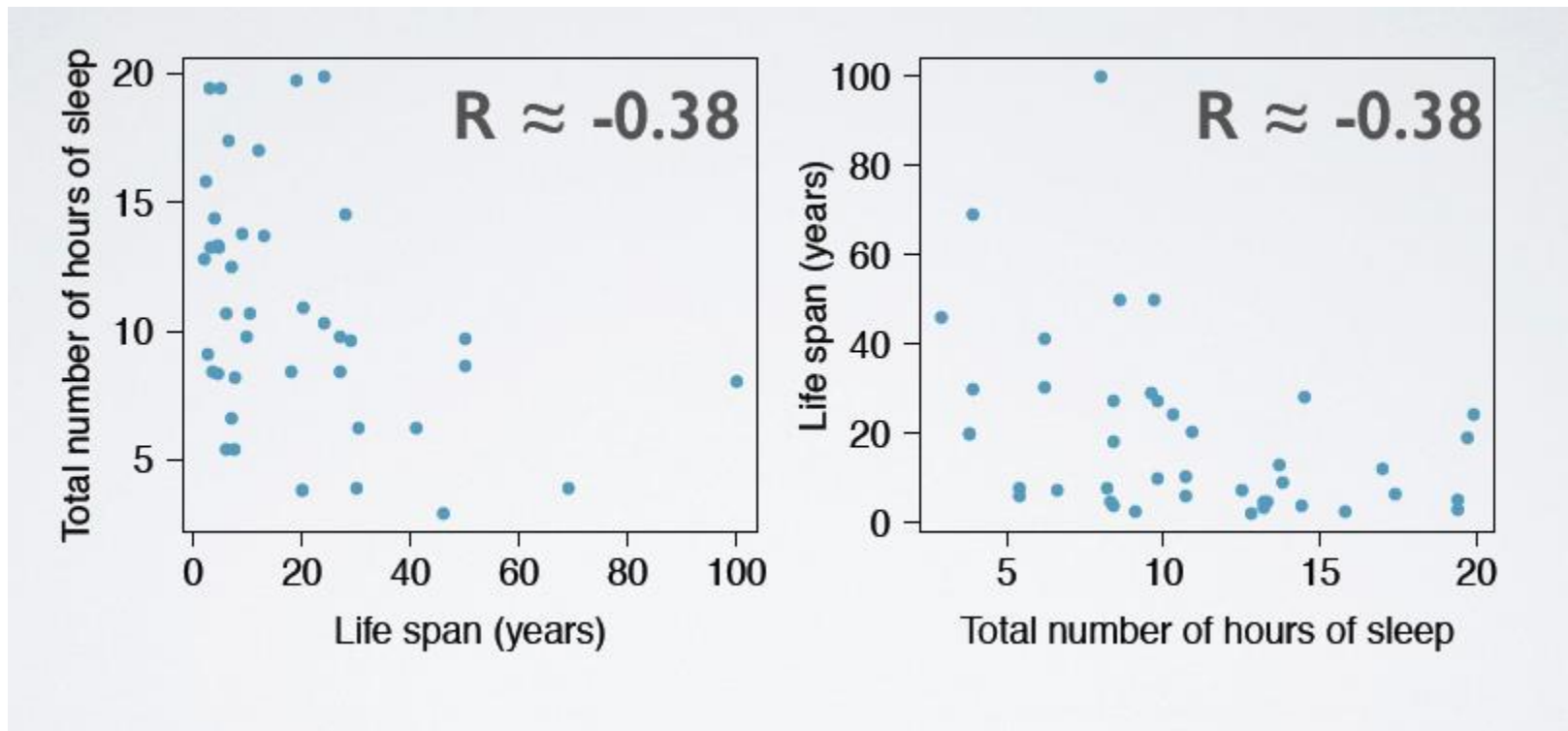
- Współczynnik jest bezwymiarowy, nie zależy od skali żadnej ze zmiennych lub przesunięcia



Korelacja: własności

8

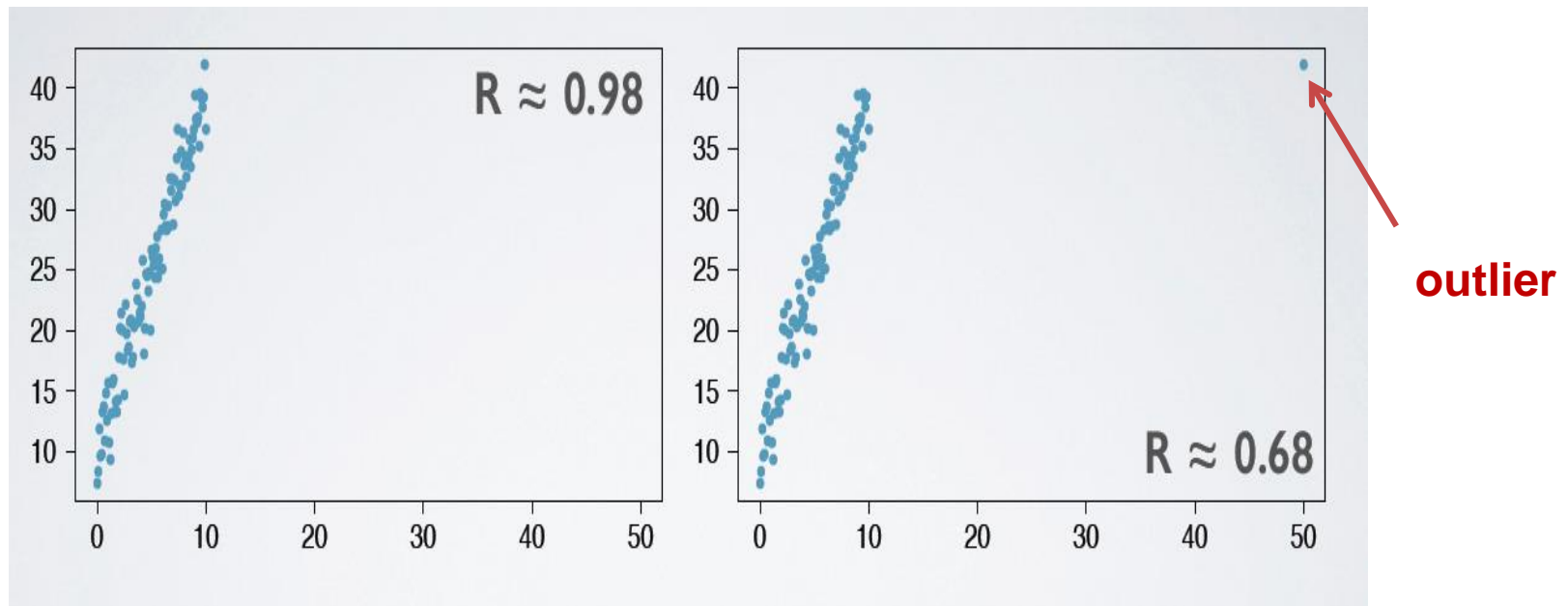
- Korelacja pomiędzy X i Y jest taka sama jak pomiędzy Y i X



Korelacja: własności

9

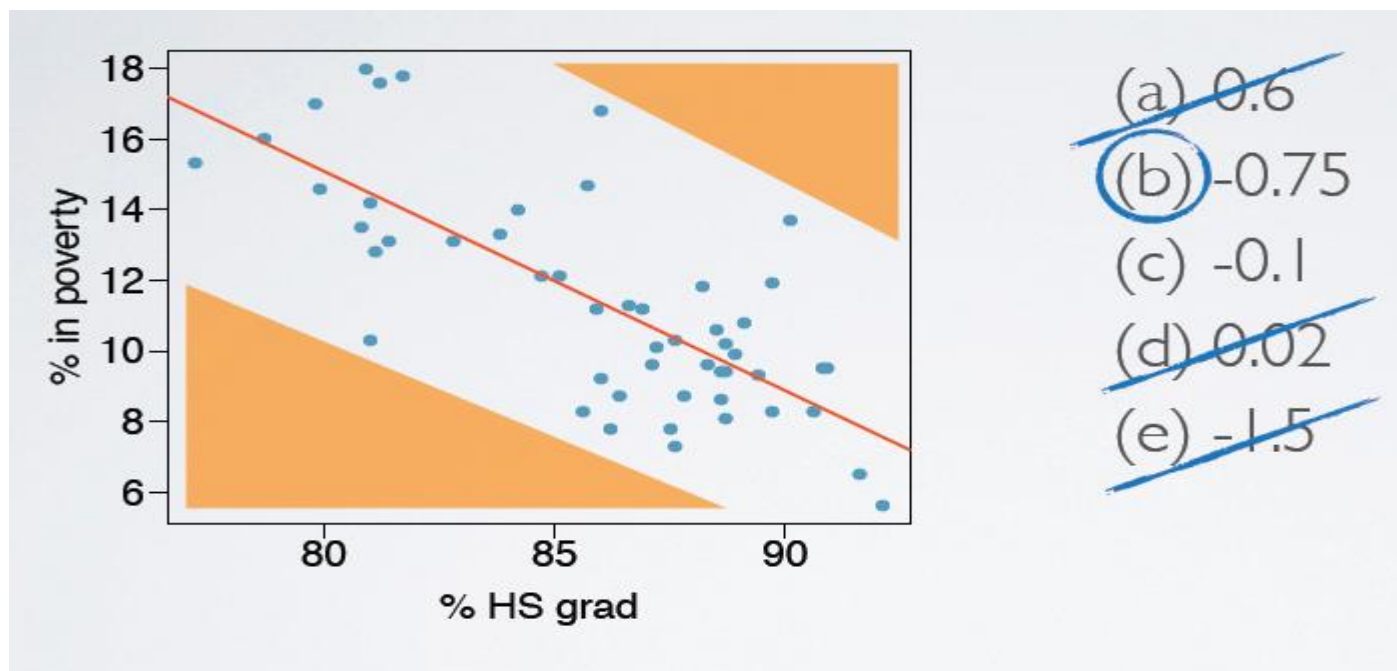
- Korelacja R jest zależna od występowania tzw. „outliers”



Korelacja

10

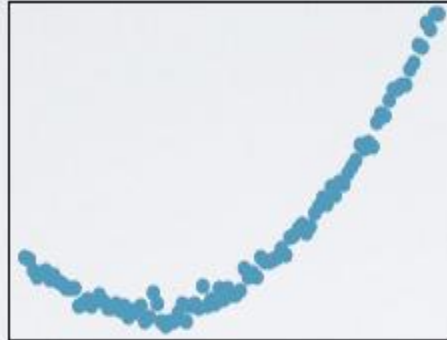
- Jaki jest R w naszym przykładzie?



Korelacja

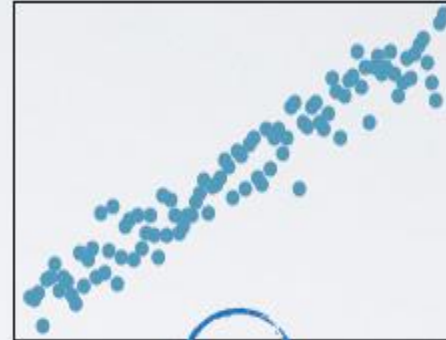
11

Bardzo silna
nieliniowa



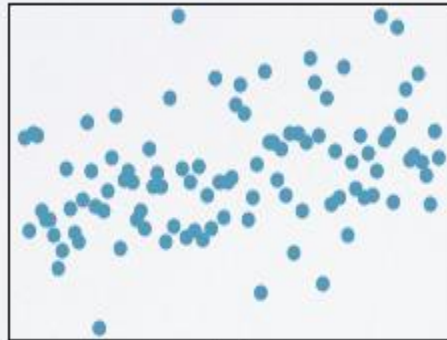
(a)

Silna
liniowa



(b)

Bardzo
słaba



(c)

słaba



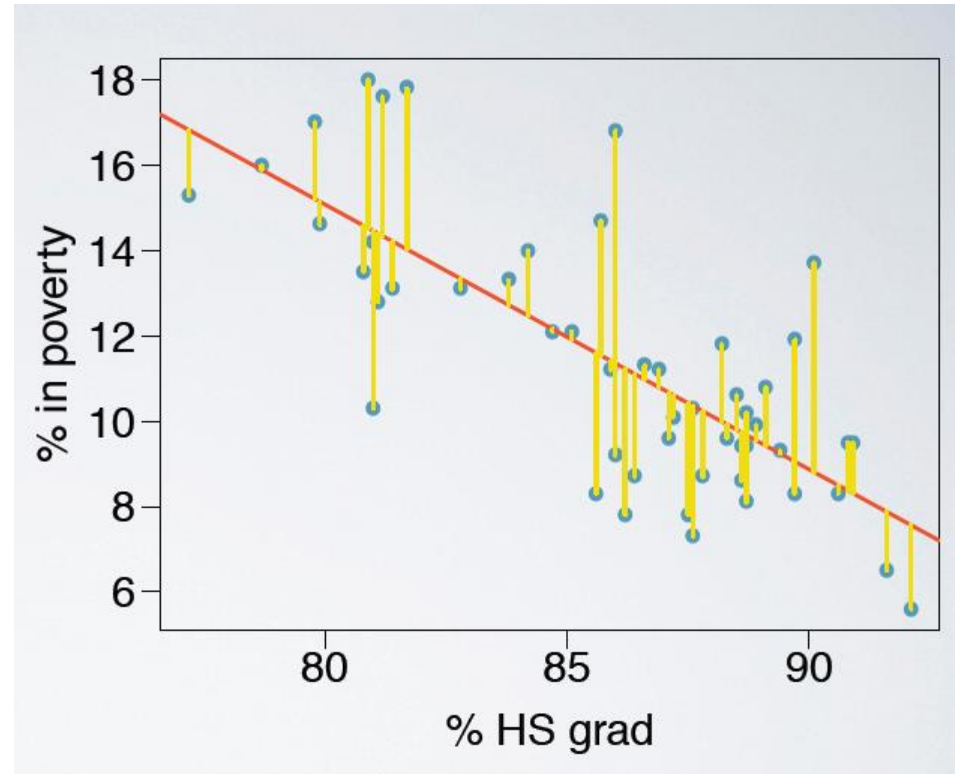
(d)

Residuals

12

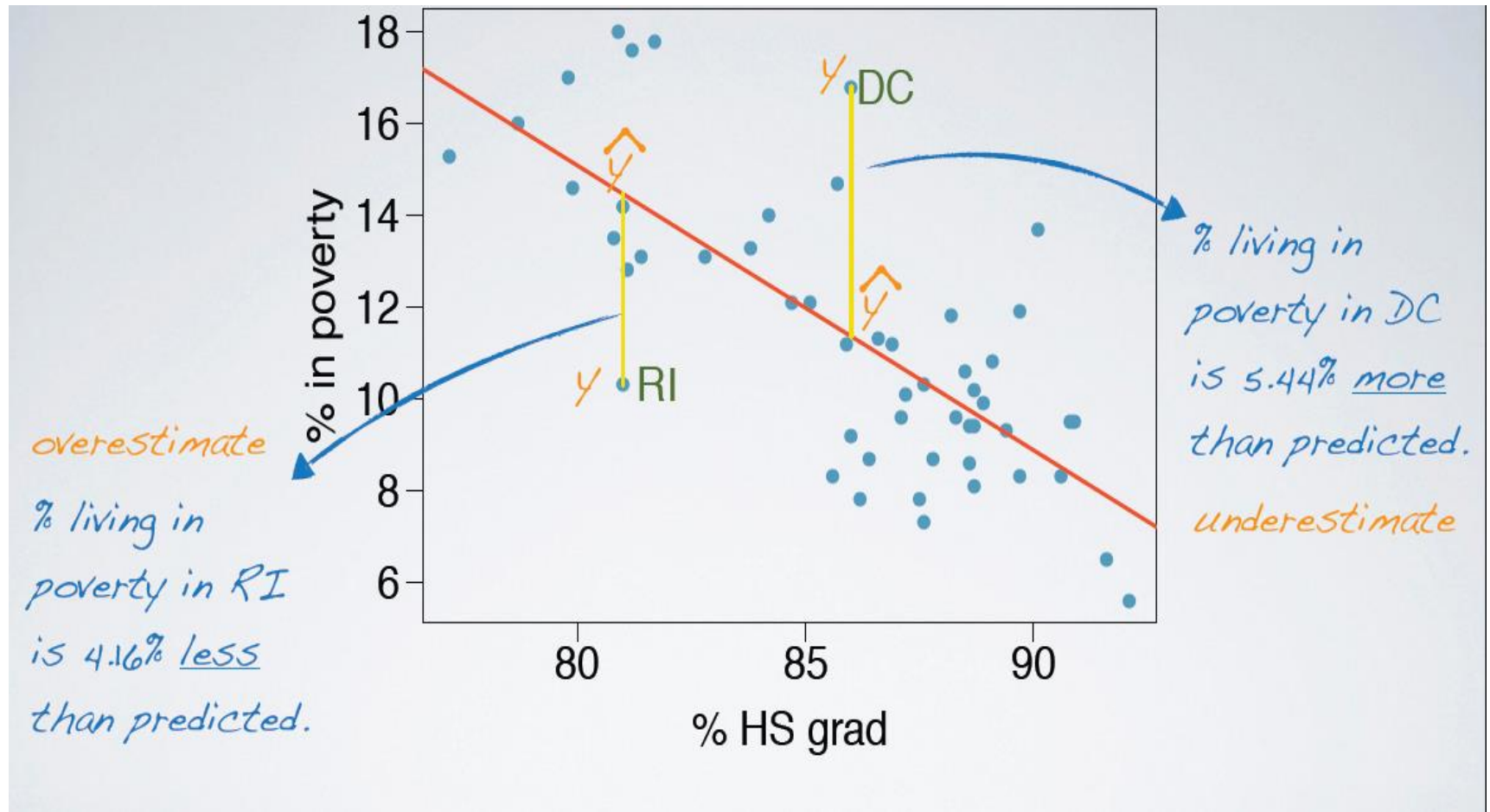
- Odchylenia od wartości modelowej
- Dane = fit + residuals
- Różnica pomiędzy przewidywane i obserwowana wartością y

$$e_i = y_i - \hat{y}_i$$



Residuals

13



Residuals

14

- W jaki sposób wyznaczymy model:
 - ▣ Minimalizacja sumy wartości odchyłeń?

$$|e_1| + |e_2| + \dots + |e_n|$$

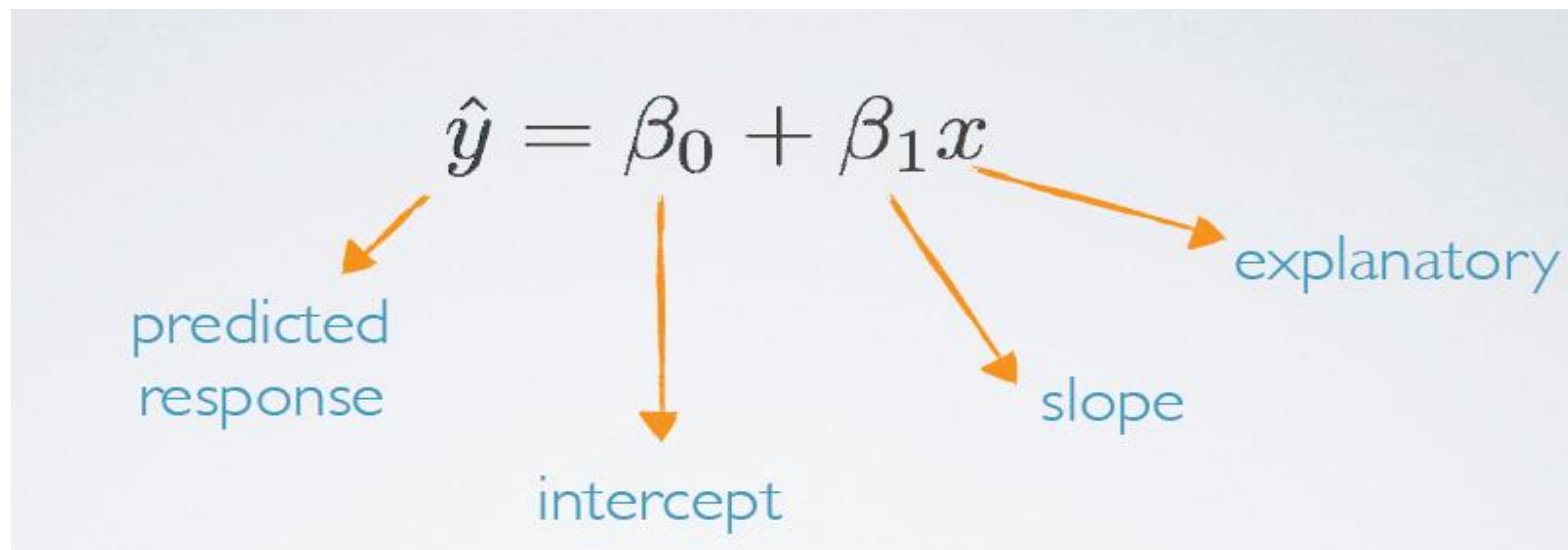
- ▣ Minimalizacja sumy kwadratów odchyłeń?

$$e_1^2 + e_2^2 + \dots + e_n^2$$

Metoda najmniejszych kwadratów

Model liniowy regresji

15



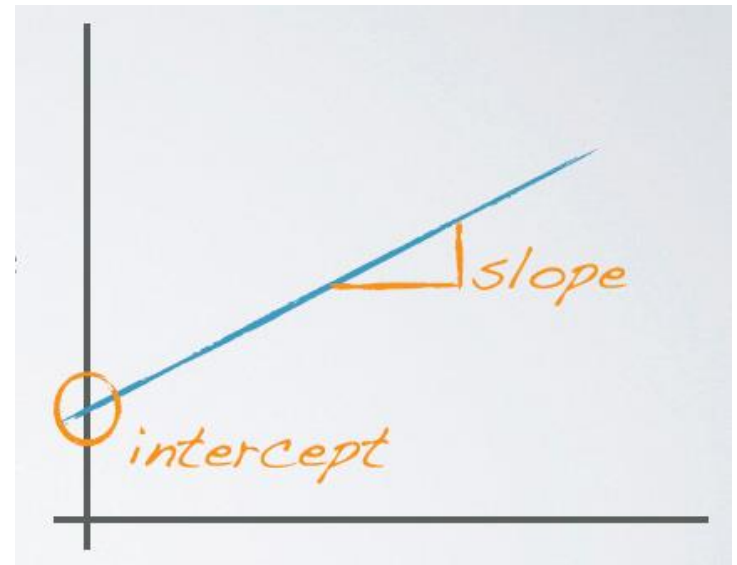
□ Notacja:

	parameter	point estimate
intercept	β_0	b_0
slope	β_1	b_1

Parametry modelu

16

- Punkt przecięcia (intercept): przewidywana wartość y dla $x=0$
- Współczynnik nachylenia (slope): dla zwiększenia x, y o jednostkę przewidywany wzrost/zmniejszenie x, y



Parametry modelu

17

□ Slope:

slope:

$$b_1 = \frac{s_y}{s_x} R$$

s_x : SD of x

s_y : SD of y

$R = \text{cor}(x, y)$

Dane:

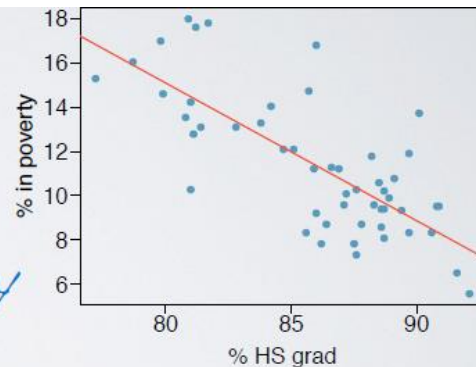
$$s_y = 3.1\%$$

$$s_x = 3.73\%$$

$$R = -0.75$$

$$b_1 = \frac{s_y}{s_x} R = \frac{3.1}{3.73} \times -0.75 \approx -0.62$$

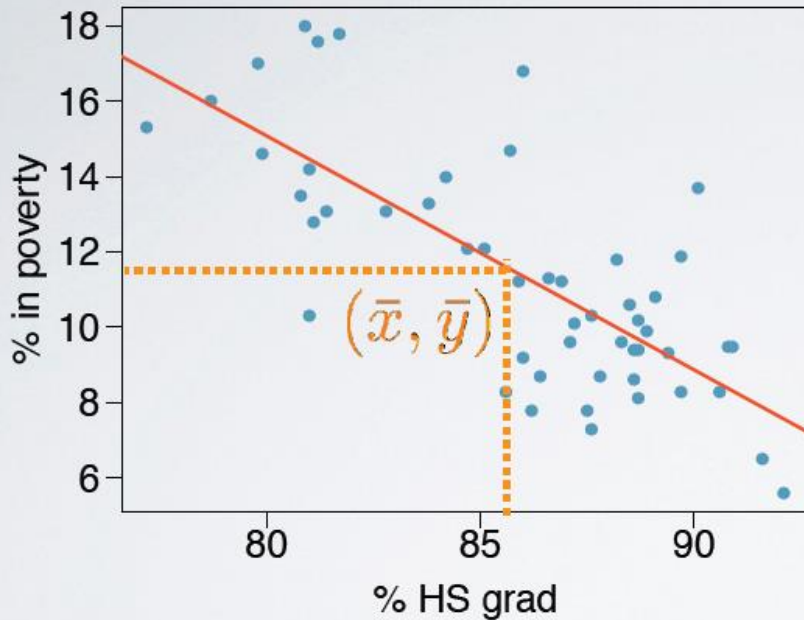
For each % point increase in HS graduate rate, we would expect the % living in poverty to be lower on average by 0.62% points.



Parametry modelu

18

□ Punkt przecięcia:



the least squares line always goes through (\bar{x}, \bar{y})

$$\bar{y} \hat{y} = b_0 + b_1 \cancel{x} \bar{x}$$

intercept: $b_0 = \bar{y} - b_1 \bar{x}$

Parametry modelu

19

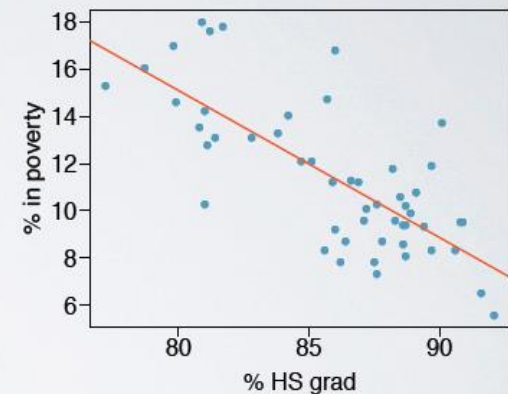
- Punkt przecięcia:
 - Często nie ma praktycznego znaczenie

Dane:

$$\bar{y} = 11.35\%$$
$$\bar{x} = 86.01\%$$

$$b_0 = \bar{y} - b_1 \bar{x} = 11.35 - (-0.62) 86.01 = 64.68$$

States with no HS graduates are expected on average to have 64.68% of their residents living below the poverty line.

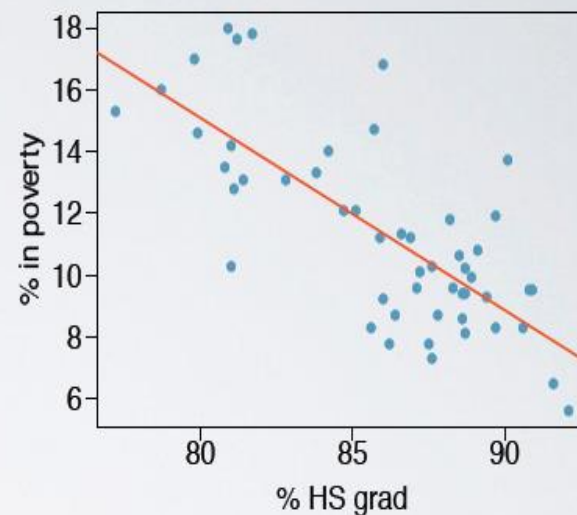


Parametry modelu liniowego

20

$$\widehat{\% \text{ in poverty}} = 64.68 - 0.62 \% \text{ HS grad}$$

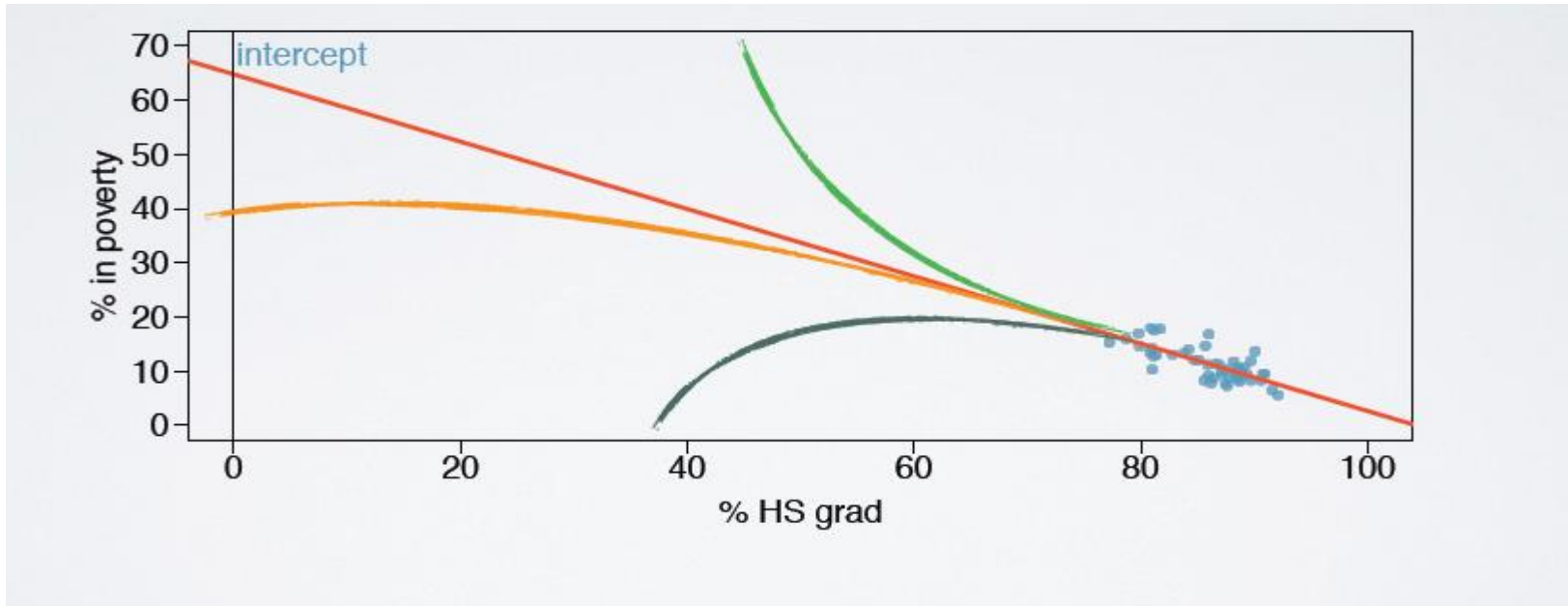
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	64.78	6.80	9.52	0.00
hsgrad	-0.62	0.08	-7.86	0.00



Extrapolacja

21

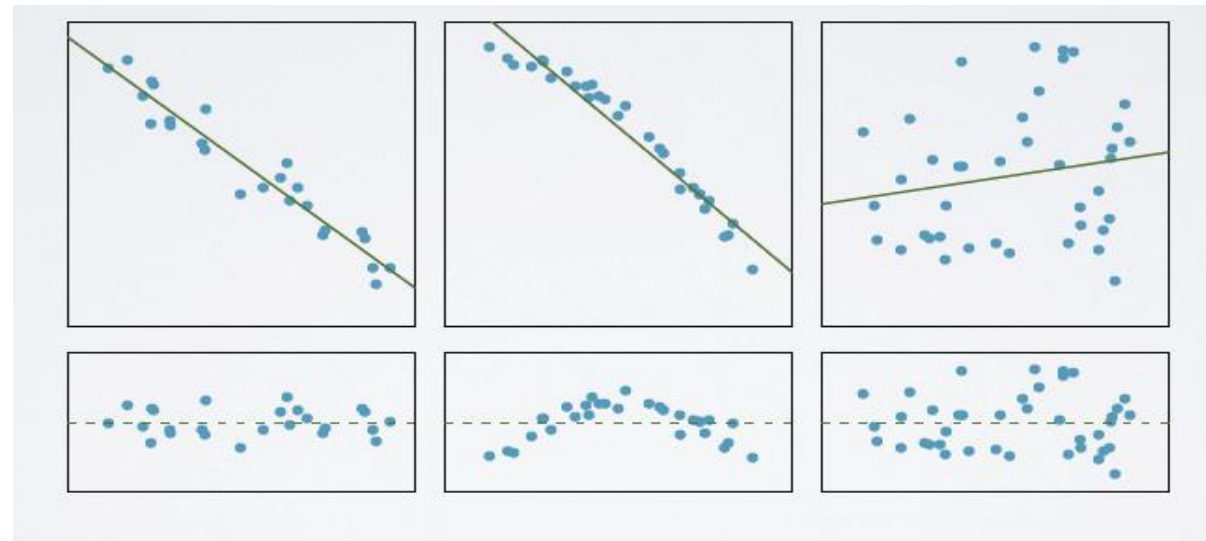
- Używamy model aby przewidzieć wartość x poza obszarem dla którego mamy punkty obserwacyjne



Warunki dla liniowej regresji

22

- Liniowość relacji pomiędzy x i y
 - ▣ Istnieją też metody aby fitować nieliniowe modele
 - ▣ Sprawdzamy używając scatter-plotu i residual plotu



Przykład:

23

* RI

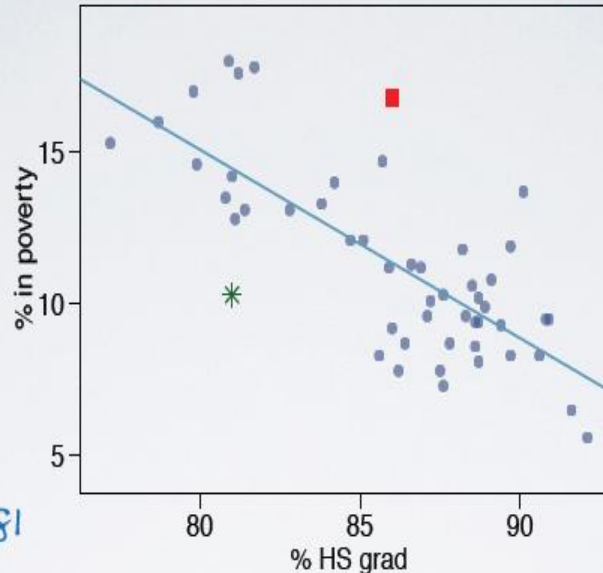
% HS grad = 81%
% in poverty = 10.3 %

$$\widehat{\%POV} = 64.68 - 0.62 \times 81$$

$$= 14.46\%$$

$$e = 10.3 - 14.46$$

$$= -4.16\%$$



■ DC

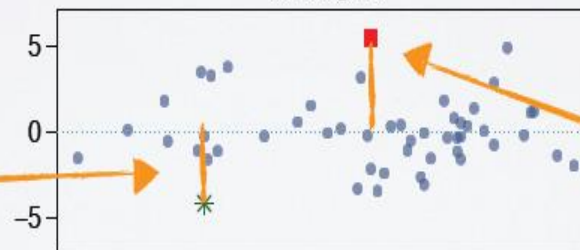
% HS grad = 86%
% in poverty = 16.8 %

$$\widehat{\%POV} = 64.68 - 0.62 \times 86$$

$$= 11.36\%$$

$$e = 16.8 - 11.36$$

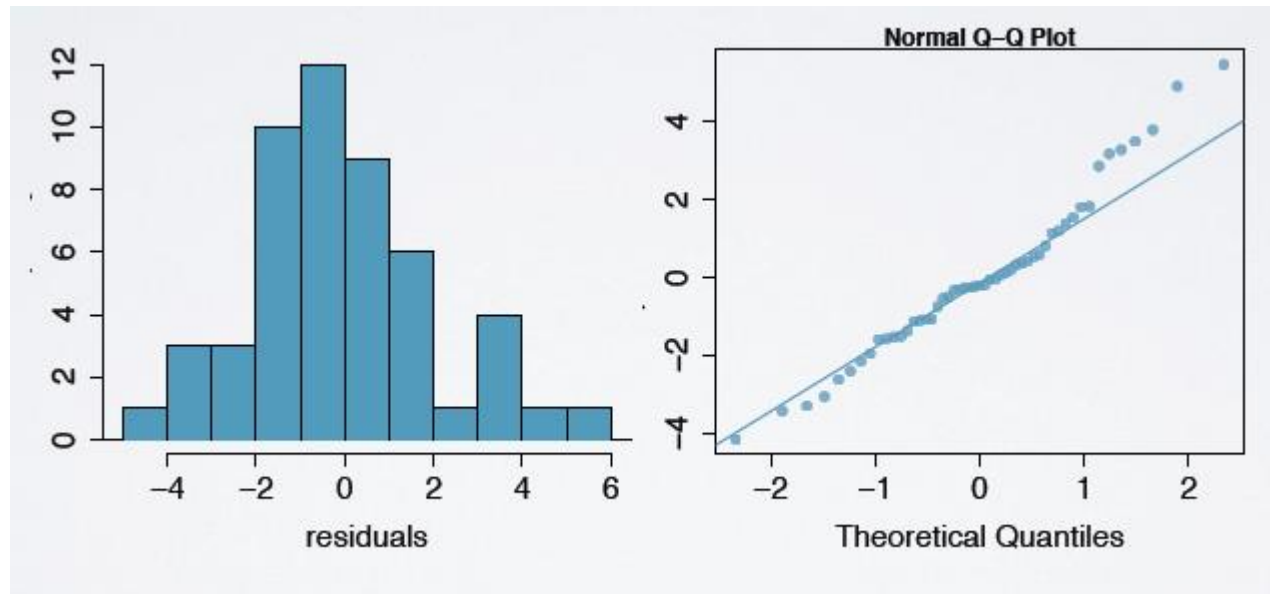
$$= 5.44\%$$



Rozkład punktów residuals

24

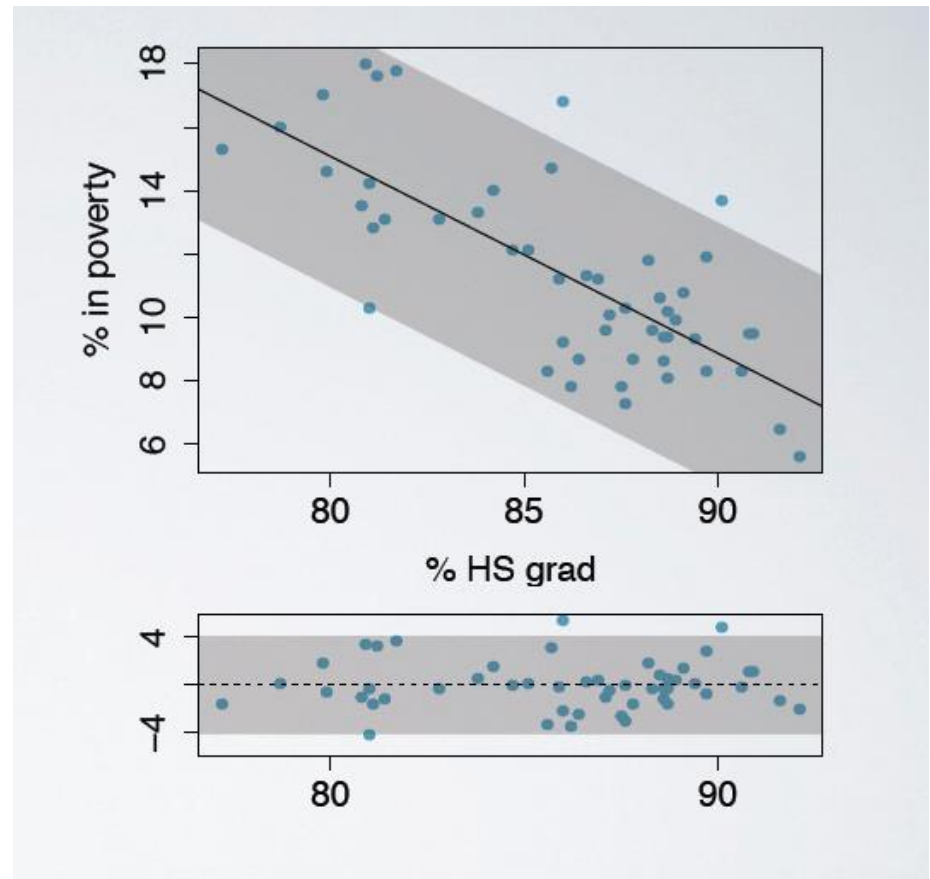
- Powinien być zbliżony do normalnego
- Może nie być spełniony jeżeli są obserwacje które mają inny trend niż reszta danych
- Możemy sprawdzić rysując histogram rozkładu percentili.



Stała zmienność

25

- Rozrzut punktów wokół linii modelu powinien być w miarę równomierny
- Sprawdzamy na płocie residuals



Możemy też użyć apletu

26

Diagnostics for simple linear regression

Select a trend:

- Linear up
- Linear down
- Curved up
- Curved down
- Fan-shaped

Show residuals

This applet uses ordinary least squares (OLS) to fit a regression line to the data with the selected trend. The applet is designed to help you practice evaluating whether or not the linear model is an appropriate fit to the data. The three diagnostic plots on the lower half of the page are provided to help you identify undesirable patterns in the residuals that may arise from non-linear trends in the data.

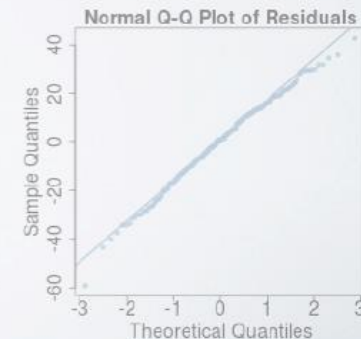
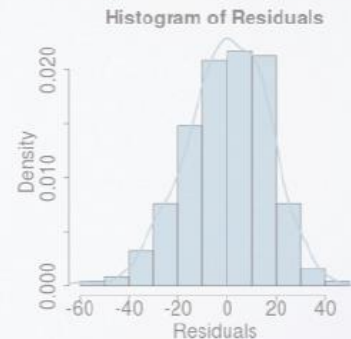
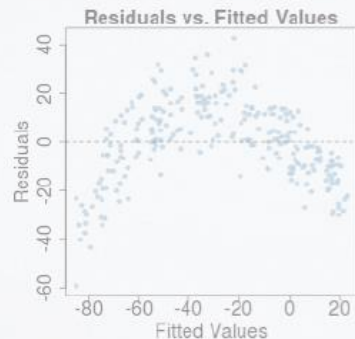
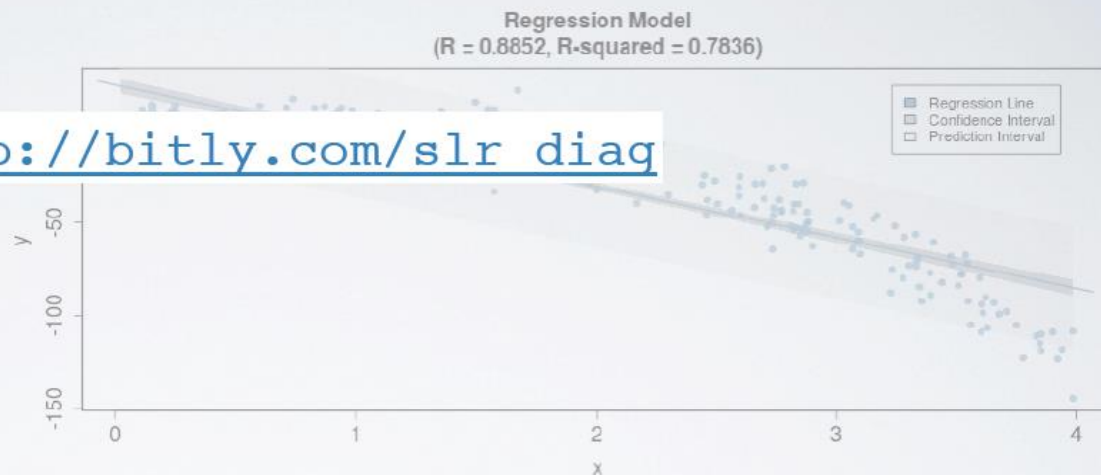
Rate this app!

View code

Check out other apps

Want to learn more for free?

http://bitly.com/slr_diag



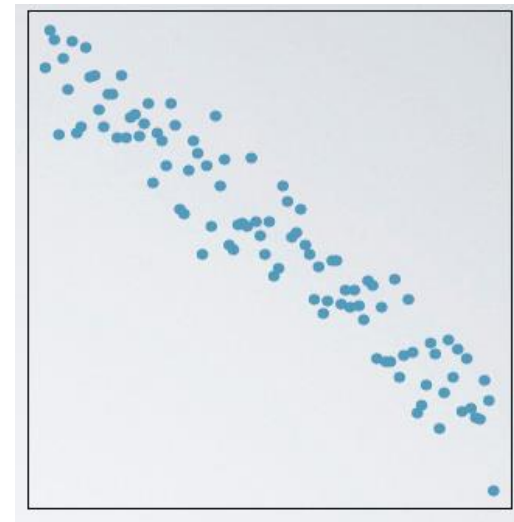
Miara zgodności modelu

27

- Jako miarę zgodności modelu często używamy R^2
- Ilościowa ocena jaki procent rozrzutu danych może być wyjaśnione przez model

$$R^2 = 92.16\%$$

$$\sqrt{0.9216} = 0.96 \rightarrow R = -0.96$$

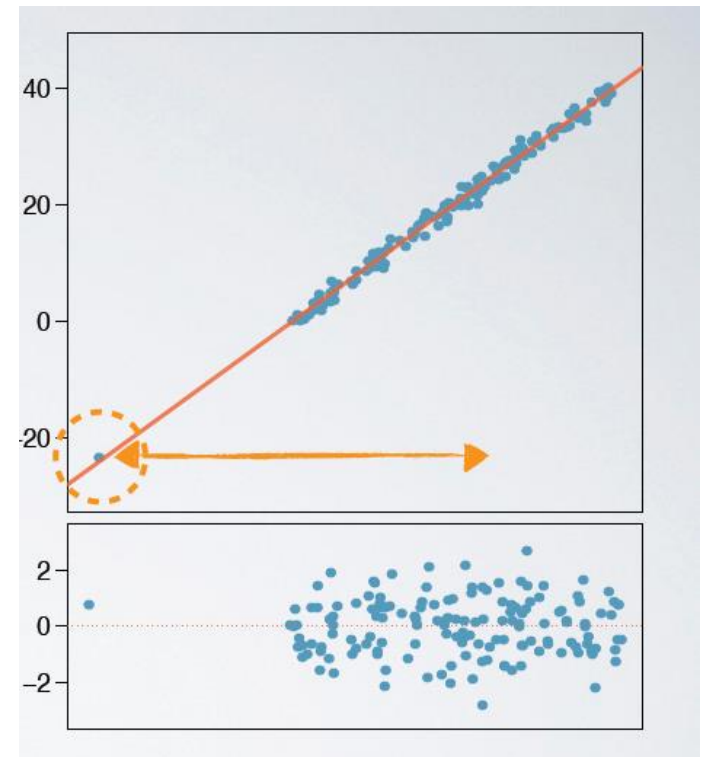


Outliers

28

- Outliers: punkty który leżą daleko od skupisk punktów
- Rozróżniemy dwie kategorie:
 - „leverage”: leżą wzdłuż osi poziomej na płocie residuals
 - „influential”: uwzględnienie silnie zmienia nachylenia linii

Leverage outlier

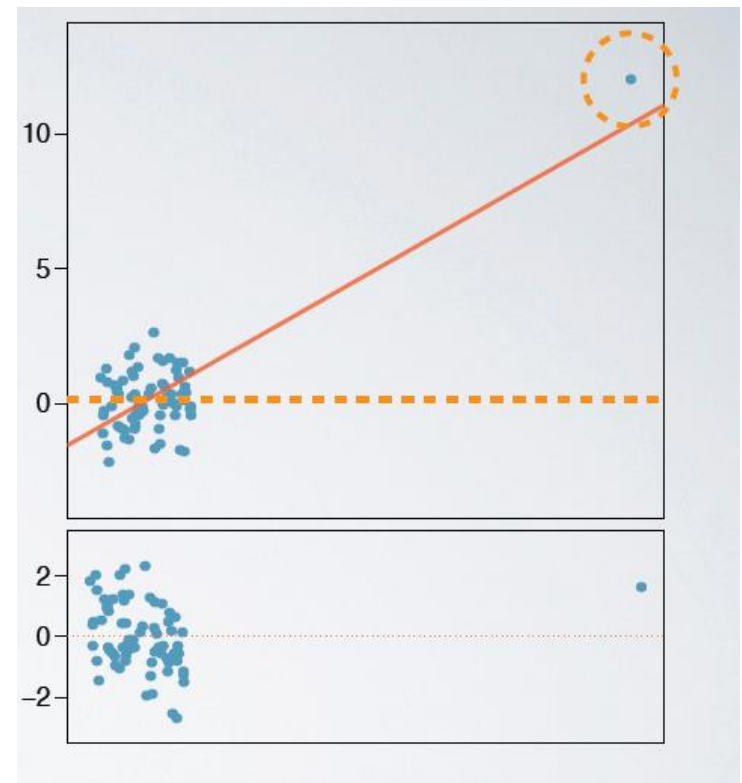


Outliers

29

- Outliers: punkty który leżą daleko od skupisk punktów
- Rozróżniemy dwie kategorie:
 - ▣ „leverage points”: leżą wzdłuż osi poziomej na płocie residuals
 - ▣ „influential points”: uwzględnienie silnie zmienia krzywą nachylenia lini
- W tym przykładzie bez punktu „outliers” nie byłoby relacji pomiędzy zmienną x i y .

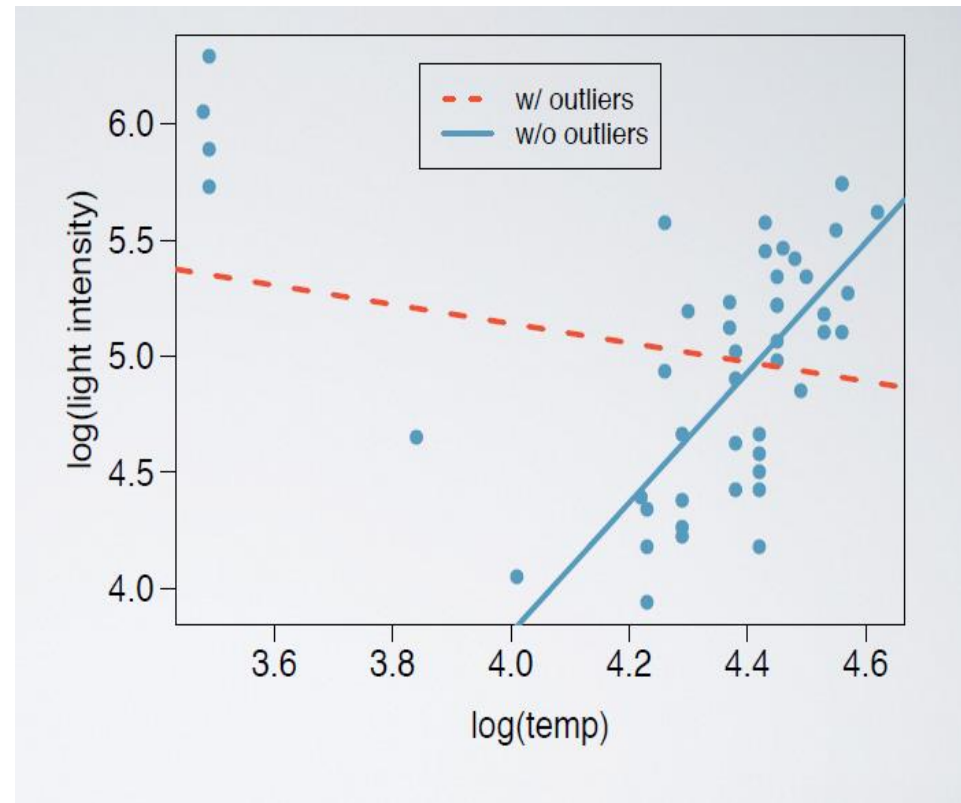
Influential outlier



Influential outliers

30

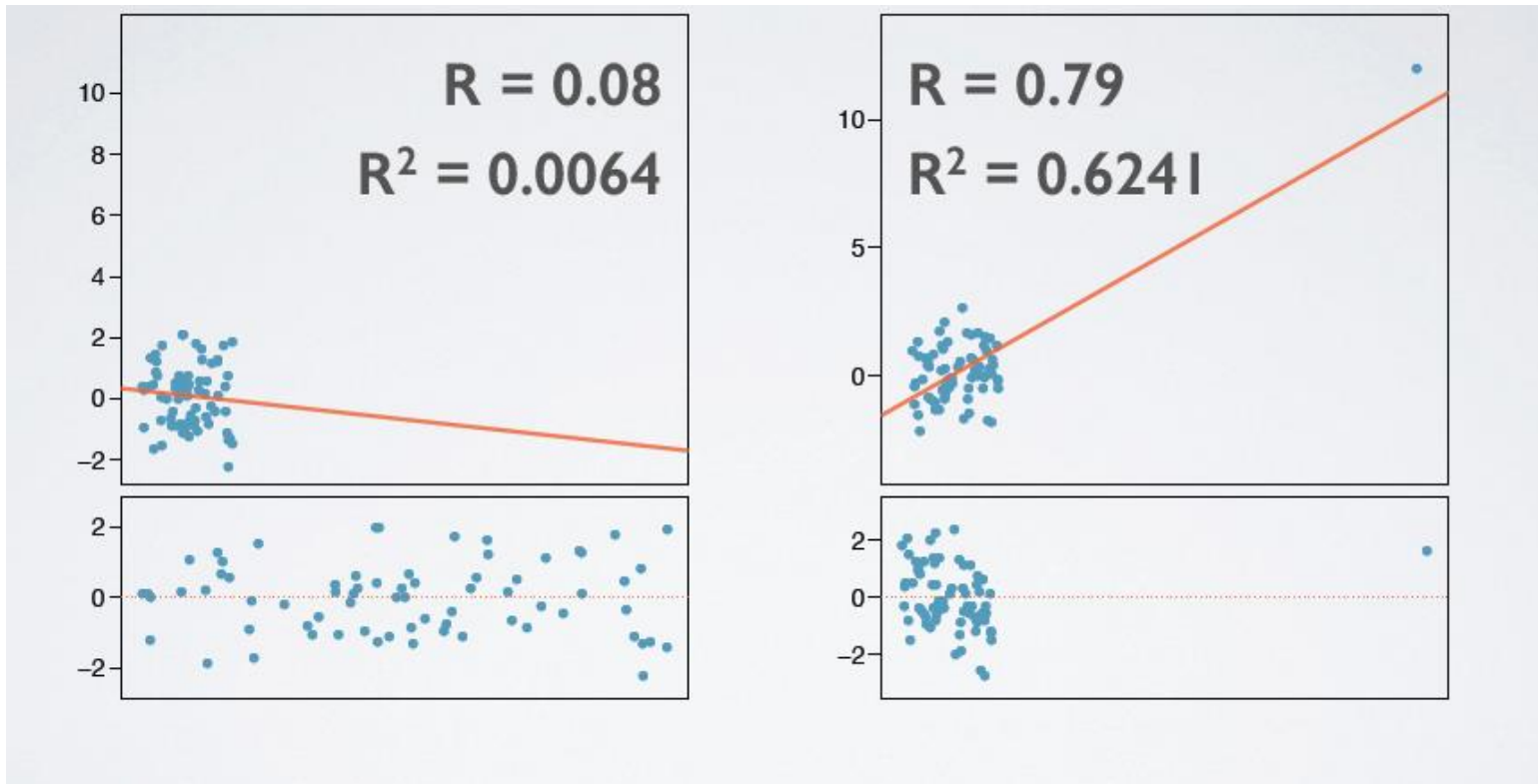
- Relacja pomiędzy intensywnością promieniowania i temperaturą powierzchni dla 47 gwiazd w clustrze CYG OBI



Influential outliers

31

- Nie muszą zmniejszać R^2



Wnioskowanie na podstawie regresji liniowej

32

□ Wyniki:

regression output:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

linear model:

$$\widehat{fosterIQ} = 9.2076 + 0.9014 \text{ bioIQ}$$

$$R^2: R^2 = 0.78$$

□ Wnioskowanie statystyczne: czy istnieje jakiś związek? Testujemy kierunek nachylenia krzywej

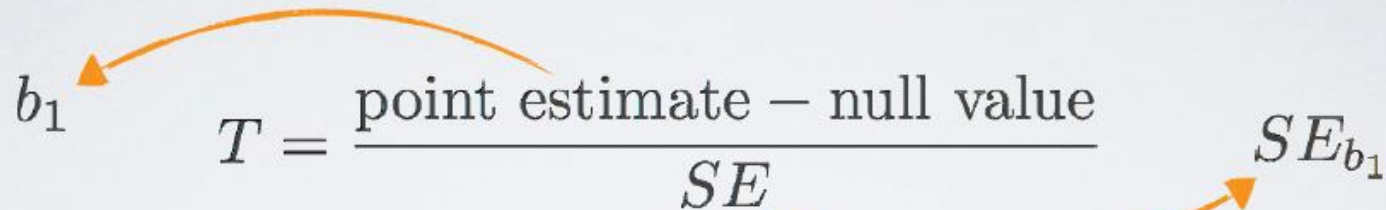
$$H_0 : \beta_1 = 0$$

$$H_A : \beta_1 \neq 0$$

Wnioskowanie statystyczne

33

- Używamy t-statystykę dla wnioskowania

$$T = \frac{\text{point estimate} - \text{null value}}{SE}$$


**t-statistic
for the slope:**

$$T = \frac{b_1 - 0}{SE_{b_1}} \quad df = n - 2$$


Odejmujemy 1 stopień swobody
dla każdego parametru: β_0 i β_1

Wnioskowanie statystyczne

34

- Używamy t-statystykę dla wnioskowania

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

$$T = \frac{0.9014 - 0}{0.0963} = 9.36$$

$$df = 27 - 2 = 25$$

$$p\text{-value} = P(|T| > 9.36) \approx 0$$

Wnioskowanie statystyczne

35

- Używamy t-statystykę dla przedziału ufności

point estimate \pm margin of error

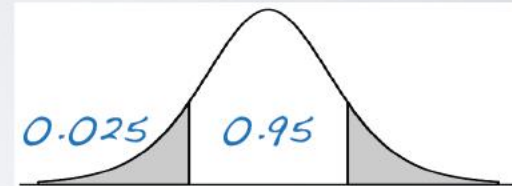
$$b_1 \pm t_{df}^* SE_{b_1}$$

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

$$df = 27 - 2 = 25$$

$$t_{25}^* = 2.06$$

$$0.9014 \pm 2.06 \times 0.0963 = (0.7, 1.1)$$



```
R
> qt(0.025, df = 25)
[1] -2.059539
```

Wnioskowanie statystyczne

36

hypothesis test:

$$T = \frac{b_1 - \text{null value}}{SE_{b_1}} \quad df = n - 2$$

confidence interval:

$$b_1 \pm t_{df}^* SE_{b_1}$$

- Hipoteza H_0 jest często $\beta_1 = 0$
- Wynik wnioskowania: b_1 , SE_{b_1} , oraz dwustronne przedział ufności for t-test.
- Wnioskowanie dla punktu przecięcia jest rzadko stosowane.

Analiza wariancji

37

- T-test pozwalał ocenić jak silna jest hipoteza o liniowej relacji pomiędzy x i y
- Alternatywy: sprawdzenie jak bardzo zmienność y jest wyjaśniona przez model: analiza wariancji.

Analiza wariancji

38

anova
output

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
bioIQ	1	5231.13	5231.13	87.56	0.0000
Residuals	25	1493.53	59.74		
Total	26	6724.66			

sum of squares

total variability in y :

$$SS_{Tot} = \sum (y - \bar{y})^2 = 6724.66$$

unexplained variability in y (residuals):

$$SS_{Res} = \sum (y - \hat{y})^2 = \sum e_i^2 = 1493.53$$

explained variability in y :

$$SS_{Reg} = 6724.66 - 1493.53 = 5231.13$$

Analiza wariancji

39

anova
output

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
bioIQ	1	5231.13	5231.13	87.56	0.0000
Residuals	25	1493.53	59.74		
Total	26	6724.66			

degrees of freedom

total degrees of freedom: $df_{Tot} = 27 - 1 = 26$

regression degrees of freedom: $df_{Reg} = 1$ *only 1 predictor*

residual degrees of freedom: $df_{Res} = 26 - 1 = 25$

Analiza wariancji

40

anova
output

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
bioIQ	1	5231.13	5231.13	87.56	0.0000
Residuals	25	1493.53	59.74		
Total	26	6724.66			

mean squares

MS regression:

$$MS_{Reg} = \frac{SS_{Reg}}{df_{Reg}} = \frac{5231.13}{1} = 5231.13$$

MS residual:

$$MS_{Res} = \frac{SS_{Res}}{df_{Res}} = \frac{1493.53}{25} = 59.74$$

F statistic

*ratio of explained to
unexplained variability*

$$F_{(1,25)} = \frac{MS_{Reg}}{MS_{Res}} = 87.56$$

Analiza wariancji

41

anova

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
bioIQ	1	5231.13	5231.13	87.56	0.0000
Residuals	25	1493.53	59.74		
Total	26	6724.66			



small p-value → reject H_0

$$H_0 : \beta_1 = 0$$

$$H_A : \beta_1 \neq 0$$

The data provide convincing evidence that the slope is significantly different than 0, i.e. the explanatory variable is a significant predictor of the response variable.

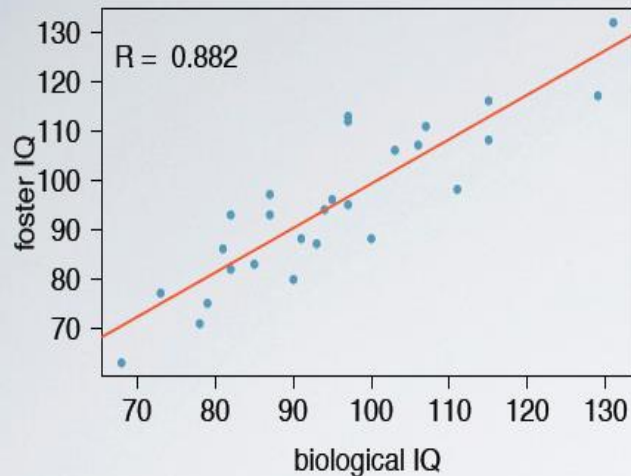
Jeszcze raz na temat R^2

42

- R^2 – proporcja zmienności wyjaśniona przez model
 - Duże: istnieje liniowy związek pomiędzy x i y
 - Małe: nie można wyciągnąć takiego wniosku
- Dwa sposoby aby to policzyć
 - Na podstawie korelacji: policzyć współczynnik korelacji
 - Z definicji: policzyć procent wyjaśnionej zmienności

Analiza wariancji

43



	Df	Sum Sq	Mean Sq	F value	Pr(>F)
bioIQ	1	5231.13	5231.13	87.56	0.0000
Residuals	25	1493.53	59.74		
Total	26	6724.66			

$$(1) \quad R^2 = \text{square of correlation coefficient} = 0.882^2 \approx 0.78$$

$$(2) \quad R^2 = \frac{\text{explained variability}}{\text{total variability}} = \frac{SS_{Reg}}{SS_{Tot}} = \frac{5231.13}{6724.66} \approx 0.78$$