

# Moderne biostatistične metode- multivariabilne metode

---

Univerza v Ljubljani  
Fakulteta *za farmacijo*



*doc. dr. Mitja Kos, mag. farm.*

Katedra za socialno farmacijo  
Univerza v Ljubljani- Fakulteta za farmacijo

# Analiza povezanosti

---

- Opazovani pojav= odvisna spremenljivka
  - Napovedni dejavnik= neodvisna spremenljivka
  - Statistični modeli:
    - Univariabilni:
      - en napovedni dejavnik
      - Povezava kot pomembna pokaže:
        - zaradi dejanske povezanosti napovednega dejavnika s pojavom
        - lahko tudi zaradi povezanosti z nekim drugim napovednim dejavnikom.
    - Multivariabilni:
      - več napovednih dejavnikov
-

# Linearna regresija

---

- Preprosta LR: matematični model = premica
  - za vsak posamezni  $y$

$$y = \beta_0 + \beta_1 x + \varepsilon$$

- $\beta_0$  in  $\beta_1$  sta parametera modela.
- $\varepsilon$  je napaka  $N(0, \sigma_e^2)$
- Pričakovana vrednost (povprečen  $y$ )

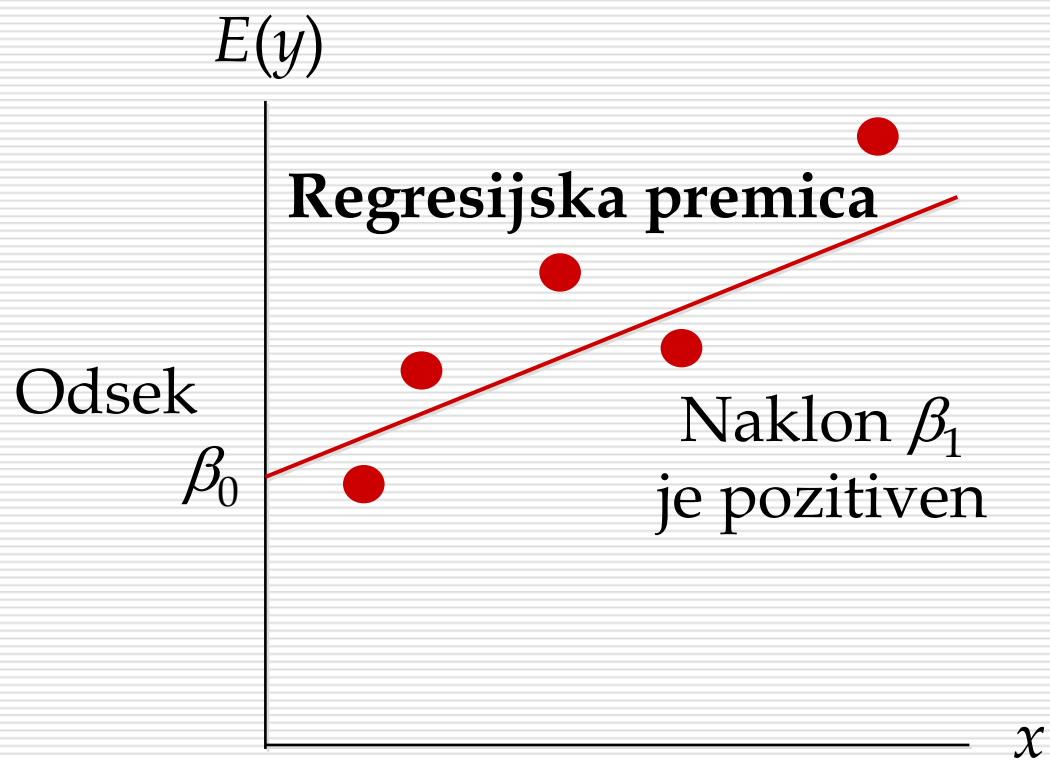
$$E(y) = \beta_0 + \beta_1 x$$

$$\beta_1 = \frac{\Delta y}{\Delta x}$$

---

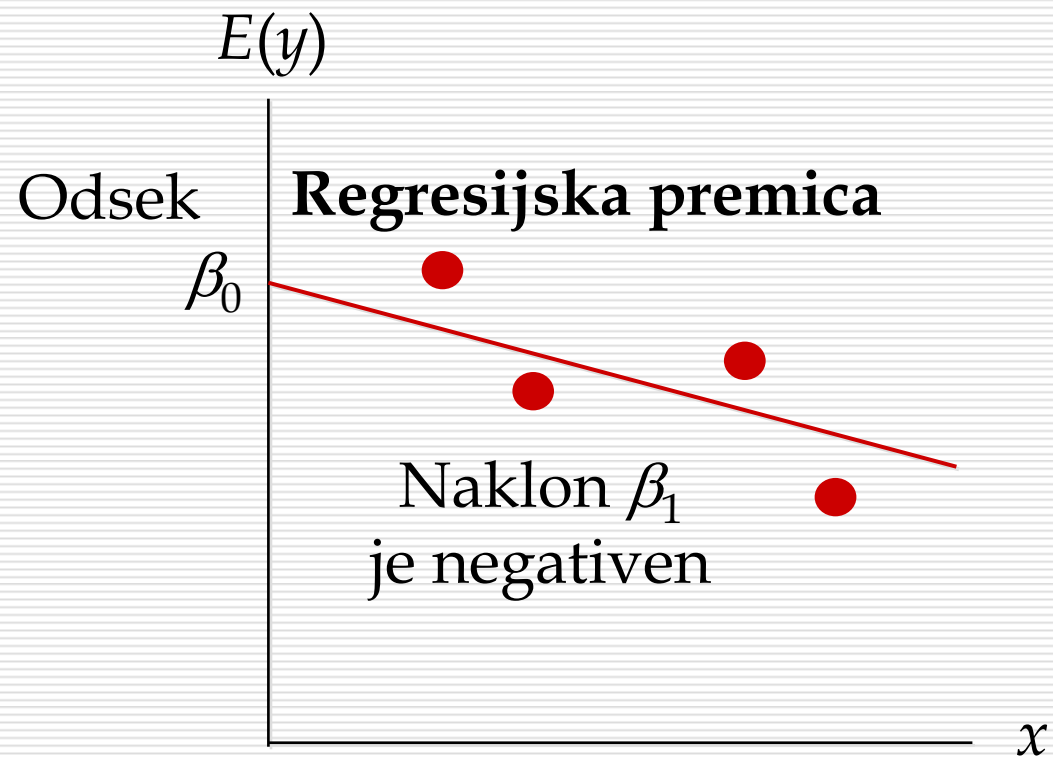
# Pozitiven linearni odnos

---



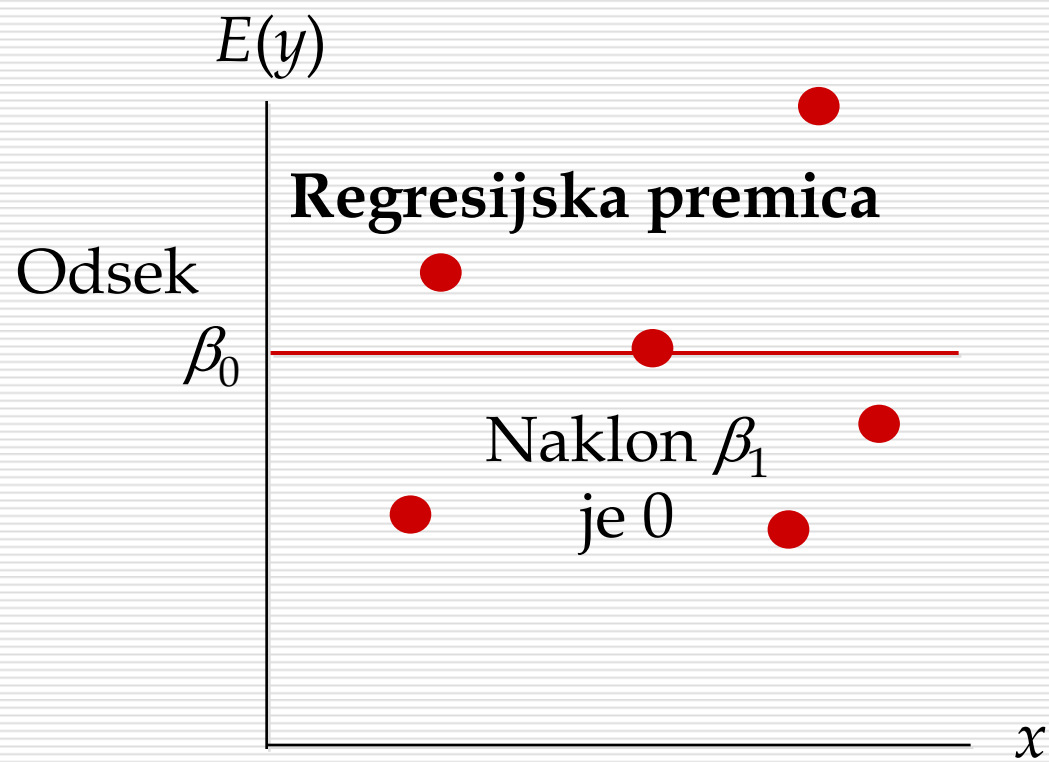
# Negativen linearni odnos

---

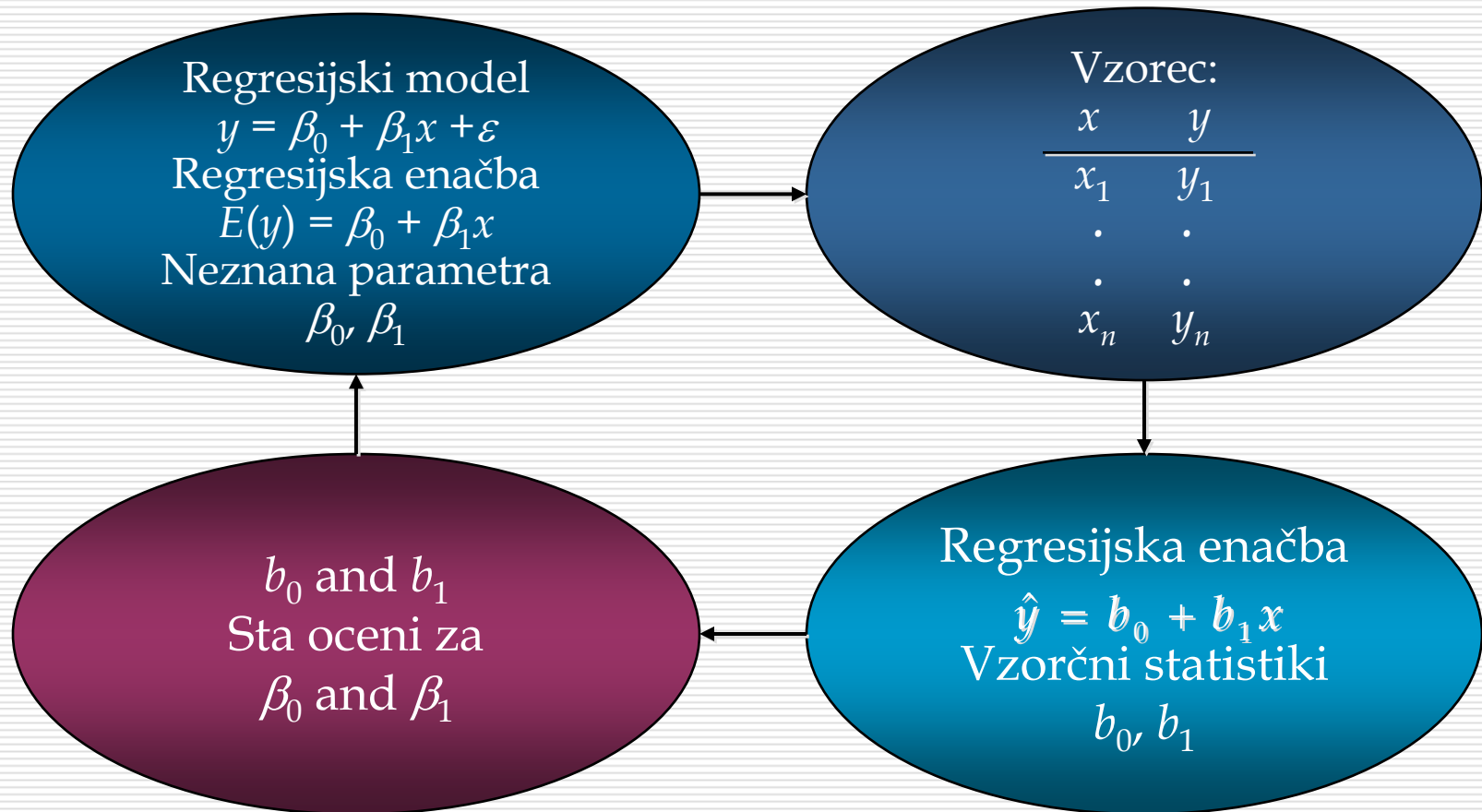


# Ni povezave

---



# Določanje parametrov modela



# Metoda najmanjših kvadratov

---

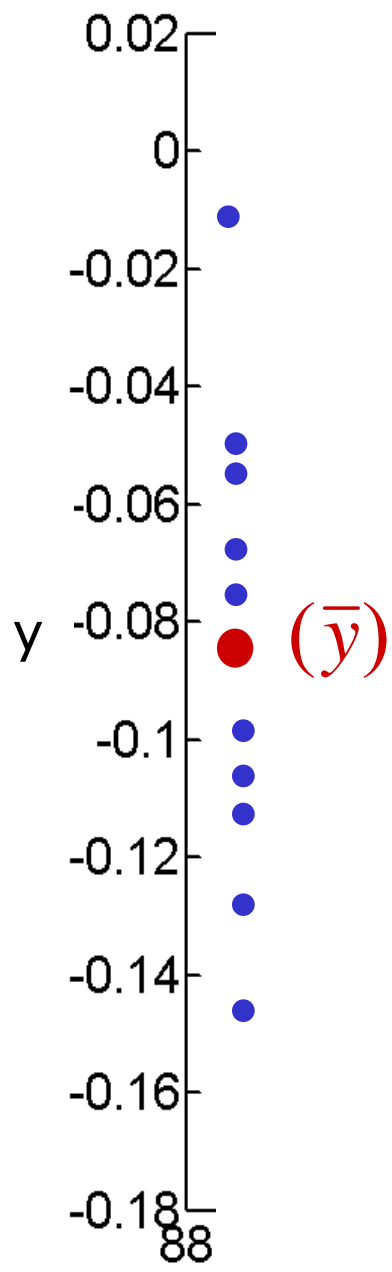
$$\min \sum (y_i - \hat{y}_i)^2$$

$\hat{y}_i$  = ocena i-te vrednosti odvisne spremenljivke

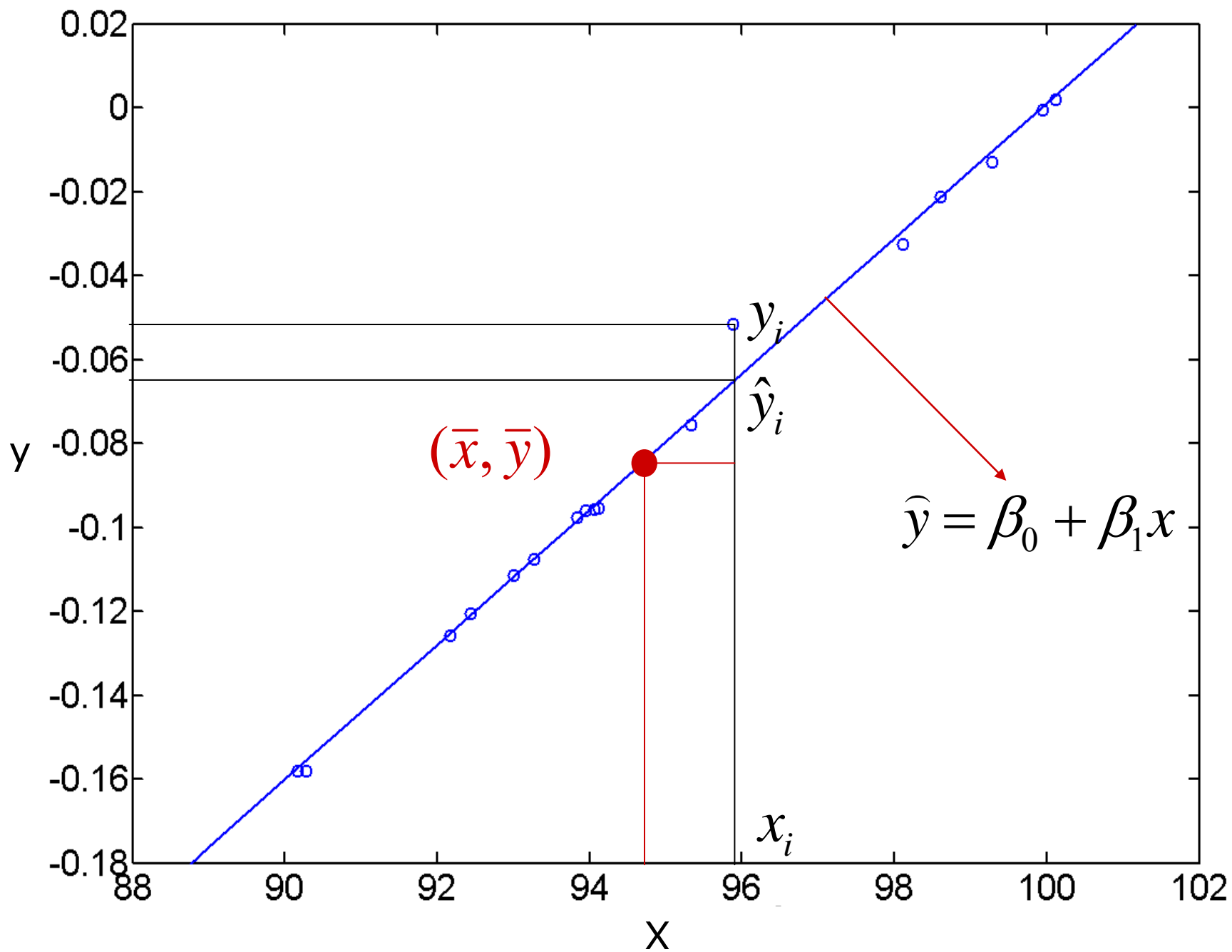
$y_i$  = opažena i-ta vrednosti odvisne spremenljivke

---



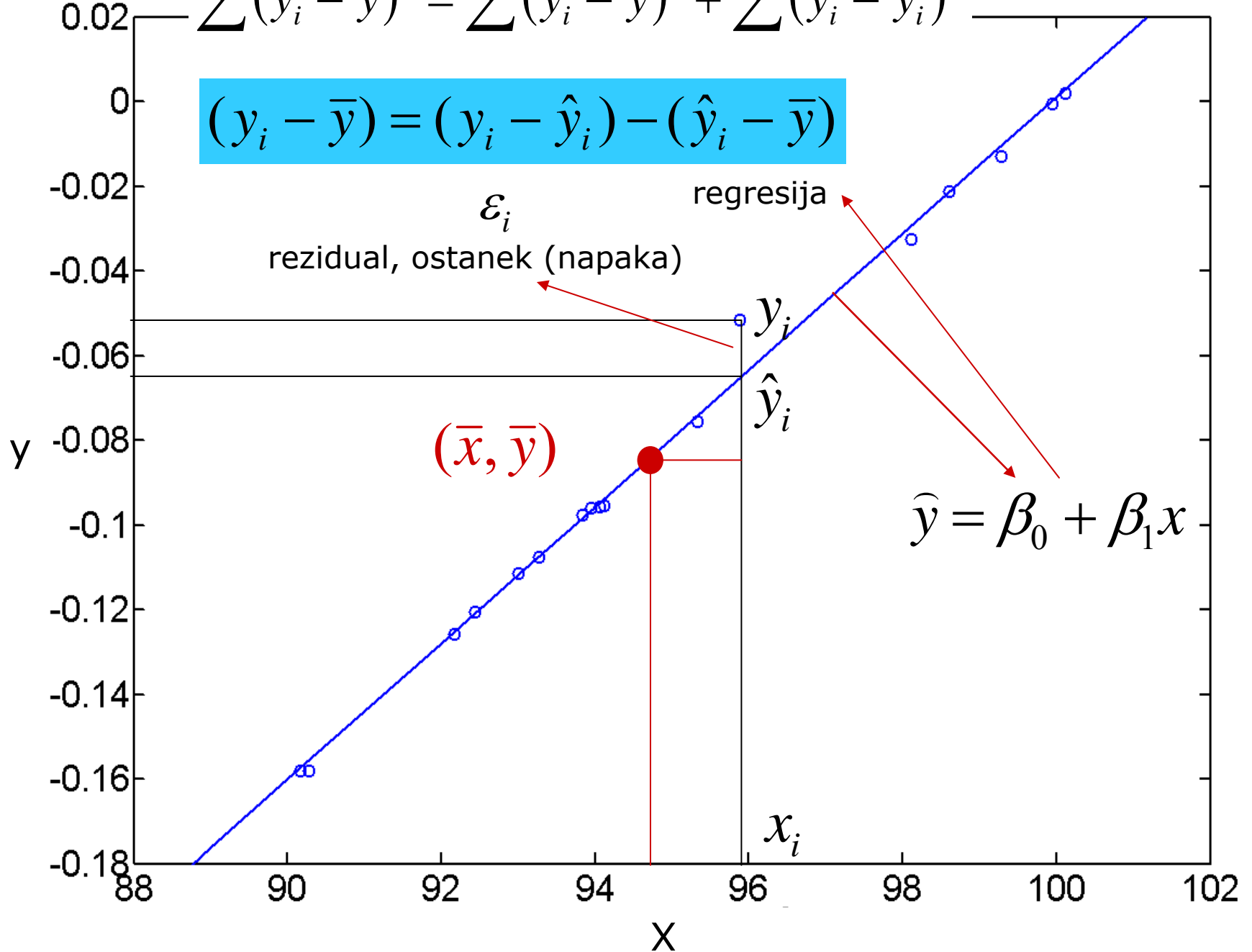


88



$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

$$(y_i - \bar{y}) = (y_i - \hat{y}_i) - (\hat{y}_i - \bar{y})$$



# Multipli regresijski model

---

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

- $\beta_0, \beta_1, \beta_2, \dots, \beta_p$  so parametri.
  - Standardizirani parametri ( $x_i \rightarrow z_i$ )
  - $\varepsilon$  je napaka, ki je slučajna spremenljivka.
-

# Metoda najmanjših kvadratov

---

□ Kriterij

$$\min \sum (y_i - \hat{y}_i)^2$$

□ Določitev koeficientov

Kompleksna algebra. Uporaba statističnih programskih paketov.

□ Interpretacija koeficientov

$$E(y) = a + b_1x_1 + b_2x_2 + \dots + b_kx_k$$

$$b_k = \frac{\Delta y}{\Delta x_k}$$

---

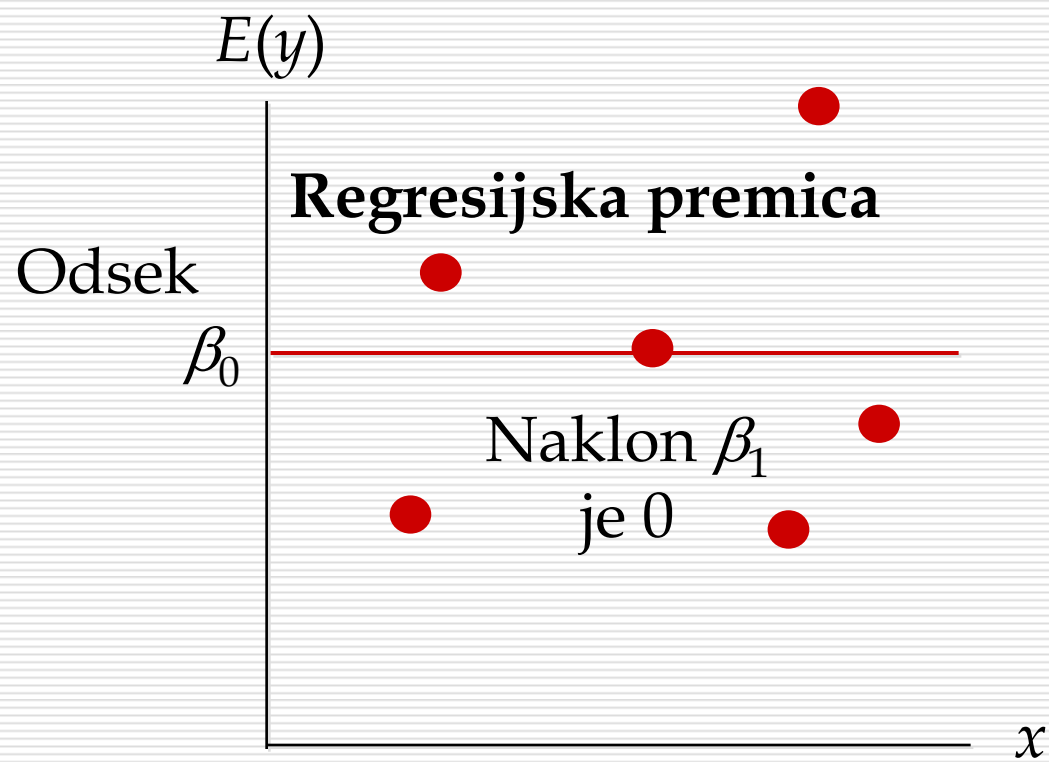
# Statistično sklepanje

---

- Pri enostavni linearni regresiji nas  $F$  and  $t$  test vodita k istim sklepom.
  - V multipli regresiji uporabimo  $F$  in  $t$  test za različna namena.
-

# Ni povezave

---



# F test

---

- Pomembnost modela kot celote.
- Test for overall significance.
- Hipotezi:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$$

$H_a$ : Najmanj eden izmed parametrov ni enak 0

---



# t test

---

- Če je izid  $F$  testa za model kot celota značilen, uporabimo  $t$  teste za ugotavljanje značilnosti vplivov posameznih neodvisnih spremenljivk.
- Za vsako neodvisno spremenljivko izvedemo en  $t$  test.
- Vsak  $t$  test imenujemo tudi test posamične značilnosti (test for individual significance).
- Hipotezi:

$$H_0: \beta_i = 0$$

$$H_a: \beta_i \neq 0$$

---

# Koeficient determinacije

---

- Multipli koeficient determinacije  $R^2$
  - Kako dobro model opisuje podatke?
-

# Primer:

## Vpliv kajenja na pljučno funkcijo

---

Z raziskavo želimo opredeliti vpliv kajenja na pljučno funkcijo. Naključno smo izbrali 1000 ljudi obeh spolov in spremljali pljučno funkcijo (parameter FEV), poleg tega smo beležili še njihovo starost in telesno višino. Za analizo rezultatov raziskave smo uporabili metodo multiple regresije. Neodvisne spremenljivke v analizi so bile:

starost [leta]

telesna višina [cm]

spol (0 = ženski, 1 = moški)

kajenje (0 = nekadilec, 1 = kadilec)

Odvisna spremenljivka pa je bila FEV [liter].

---

# Primer:

## Vpliv kajenja na pljučno funkcijo

**Model Summary**

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.497 <sup>a</sup>	.247	.243	.25269

a. Predictors: (Constant), kajenje, t. višina (cm), starost (leta), spol

**ANOVA<sup>b</sup>**

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	20.786	4	5.196	81.383	.000 <sup>a</sup>
	Residual	63.533	995	.064		
	Total	84.319	999			

a. Predictors: (Constant), kajenje, t. višina (cm), starost (leta), spol

b. Dependent Variable: FEV (l)

**Coefficients<sup>a</sup>**

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	2.666	.118		22.568	.000
	starost (leta)	-1.9E-03	.001	-.063	-2.296	.022
	t. višina (cm)	-5.6E-04	.001	-.024	-.844	.399
	spol	.219	.017	.376	13.208	.000
	kajenje	-.275	.023	-.327	-11.857	.000

a. Dependent Variable: FEV (l)



# Tveganje, obeti ter razmerja

---

Relativno tveganje  
ang. Relative Risk

$$RR = \frac{p_1}{p_2}$$

Razmerje obetov  
ang. Odds Ratio

$$\theta = \frac{\frac{p_1}{1 - p_1}}{\frac{p_2}{1 - p_2}}$$

Povezava med razmerjem obetov ter relativnim tveganjem:

$$\theta = \frac{\frac{p_1}{1 - p_1}}{\frac{p_2}{1 - p_2}} = \frac{p_1}{p_2} \times \frac{1 - p_2}{1 - p_1} = RR \times \frac{1 - p_2}{1 - p_1}$$

---

# Preprosta logistična regresija

---

- Opazovani pojav = binarna spremenljivka 0/1
  - Kako uporabiti linearni model za opisovanje odnosa?
-

# Preprosta logistična regresija

---

Vrednost binarne spremenljivke zapisati kot:

□ Verjetnost, da zavzame vrednost 1 pri danem  $X$ :

□  $\pi(x)$ - populacija

□  $p(x)$ - vzorec

■ Zavzame le vrednosti med 0 in 1- preozek razpon

□ Razmerje verjetnosti, da dogodek zgodi in da se ne zgodi:

$$\frac{p(x)}{1 - p(x)}$$

■ Zavzame vrednosti med 0 in  $+\infty$  neskončnostjo



# Preprosta logistična regresija

---

- Logaritmiranje  $\Rightarrow$  med  $-$  in  $+$  neskončnostjo

$$\ln \left[ \frac{p(x)}{1-p(x)} \right] = \text{logit}$$

Populacija:

$$\ln \left[ \frac{\pi(x)}{1-\pi(x)} \right] = a + bx = \text{logit}$$

Ocena parametrov iz vzorca:

$$\ln \left[ \frac{p(x)}{1-p(x)} \right] = a + bx = \text{logit}$$

---

$$p(x) = \frac{e^{a+bx}}{1+e^{a+bx}}$$

Verjetnost, da dogodek  
zgodil pri danem X

# Pomen regresijskega koef. b

---

□ Lin. regresija:  $b = y_{(x+1)} - y_{(x)}$   $\ln\left[\frac{p(x)}{1-p(x)}\right] = a + bx = \text{logit}$

□ Log. regresija:  $b = \text{logit}_{(x+1)} - \text{logit}_{(x)}$

□ Binarni X (0/1)  $\Rightarrow x=0, x+1=1.$

□ Obet da bo  $y=1$ :

$$O_{x=1} = \frac{p(1)}{1-p(1)}$$

$$O_{x=0} = \frac{p(0)}{1-p(0)}$$

□ Logaritmiranje:

$$\text{logit}_{x=1} = \ln\left[\frac{p(1)}{1-p(1)}\right]$$

$$\text{logit}_{x=0} = \ln\left[\frac{p(0)}{1-p(0)}\right]$$

---

# Pomen regresijskega koef. $b$

---

- Log. razmerja obetov:

$$RO = \frac{O_{x=1}}{O_{x=0}} = \frac{\left[ \frac{p(1)}{1-p(1)} \right]}{\left[ \frac{p(0)}{1-p(0)} \right]}$$

$\leftarrow O_{x=1} = \frac{p(1)}{1-p(1)}$   
 $\leftarrow O_{x=0} = \frac{p(0)}{1-p(0)}$

$$\ln RO = \ln \frac{\left[ \frac{p(1)}{1-p(1)} \right]}{\left[ \frac{p(0)}{1-p(0)} \right]} = \log it(1) - \log it(0)$$

$$\ln RO = b = \log it_{(1)} - \log it_{(0)}$$

$$RO = e^b$$

# Ocenjevanje b v vzorcu

---

□  $p(x)? \Rightarrow b$  in  $a$  ?

$$\ln \left[ \frac{p(x)}{1-p(x)} \right] = a + bx = \text{logit}$$

- Lin. regresija: metoda najmanjših kvadratov ostankov.
- Log. regresija: metoda največjega verjetja (maximum likelihood method).
  - Funkcija največjega verjetja
  - Oz. logaritem funkcije verjetja ("log likelihood"):
    - Nelinearna funkcija parametrov modela  $a$  in  $b$ .
    - Iteracijska metoda, več ocen parametrov. Nove ocene, dokler še zveča funkcijo največjega verjetja. Lokalna/globalna točka največjega verjetja. Start?

# Ocenjevanje b v populaciji

---

- Vzorčna porazdelitev b- ja: normalna porazdelitev = > interval zaupanja:

$$\beta = b \pm z * SE_{(b)}$$

- Vzorčna ocena razmerja obetov: nenormalna, nesimetrična porazdelitev
- Interval zaupanja za RO iz b (spodnjo/zgornjo mejo):

$$e^{b \pm z * SE_b}$$

---

# Multipla logistična regresija

---

- Matematični model:

$$p(x) = \frac{e^{a+b_1x_1+b_2x_2+\dots+b_kx_k}}{1+e^{a+b_1x_1+b_2x_2+\dots+b_kx_k}}$$

- $b(k)$ : sprememba v logitu, ki spremlja spremembo  $X(k)$  za 1 enoto, medtem ko se ostale sprem. ne spreminjajo.
  - Regresijski koef.  $b$  = logaritem razmerja obojev
  - Razmerje obojev = antilogaritem  $b$ :  $e^b$
  - Metoda največjega verjetja
-

# Testiranje pomembnosti modela kot celote

---

- Ali se log verjetja modela s spremenljivkami  $X$  statistično značilno poveča v primerjavi z log verjetja modela brez njih.
- G- statistika- test razmerja dveh verjetij:

$$G = 2 \ln \left[ \frac{\text{verjetje}_{SPREM.X}}{\text{verjetje}_0} \right]$$

- Porazdeljuje po hi- kvadrat, df:  $k-1$
  - $k$ : št. vseh spremenljivk v modelu
  - V SPSS-u: "Omnibus test"
-

# Vrednotenje prileganja modela

---

- Nagelkerkejev  $R^2$
  
  - Test hi- kvadrat
  - Hosmer- Lemeshov test:
    - Oba primerjata opazovano število enot, pri katerem se je opazovani dogodek zgodil, s pričakovanim številom, ki temelji na enačbi logistične regresije.
-



# Testiranje pomembnosti b

---

□  $H_0: b=0, H_a: b \neq 0$

□ Testi:

■ Z: 
$$z = \frac{b-0}{SE_b}$$

■ Waldova statistika:

$$\chi^2 = \left( \frac{b}{SE_b} \right)^2$$

■ G- statistika:

$$G = 2 \ln \left[ \frac{\text{verjetje}_x}{\text{verjetje}_0} \right]$$

G- statistika:  
Model, ki  
vključuje X, pove  
več o Y.

---