

# **Korelacija i Regresija**

**Nedelja 9/10 - Vežbe**

**Dr Nikola N. Grubor**



**“Netačna pretpostavka da korelacija podrazumeva uzročnu vezu je verovatno jedna od dve ili tri najveće greške u čovekovom zaključivanju.”**

**Stephen Jay Gould**

# Korelacija ≠ Kauzalnost



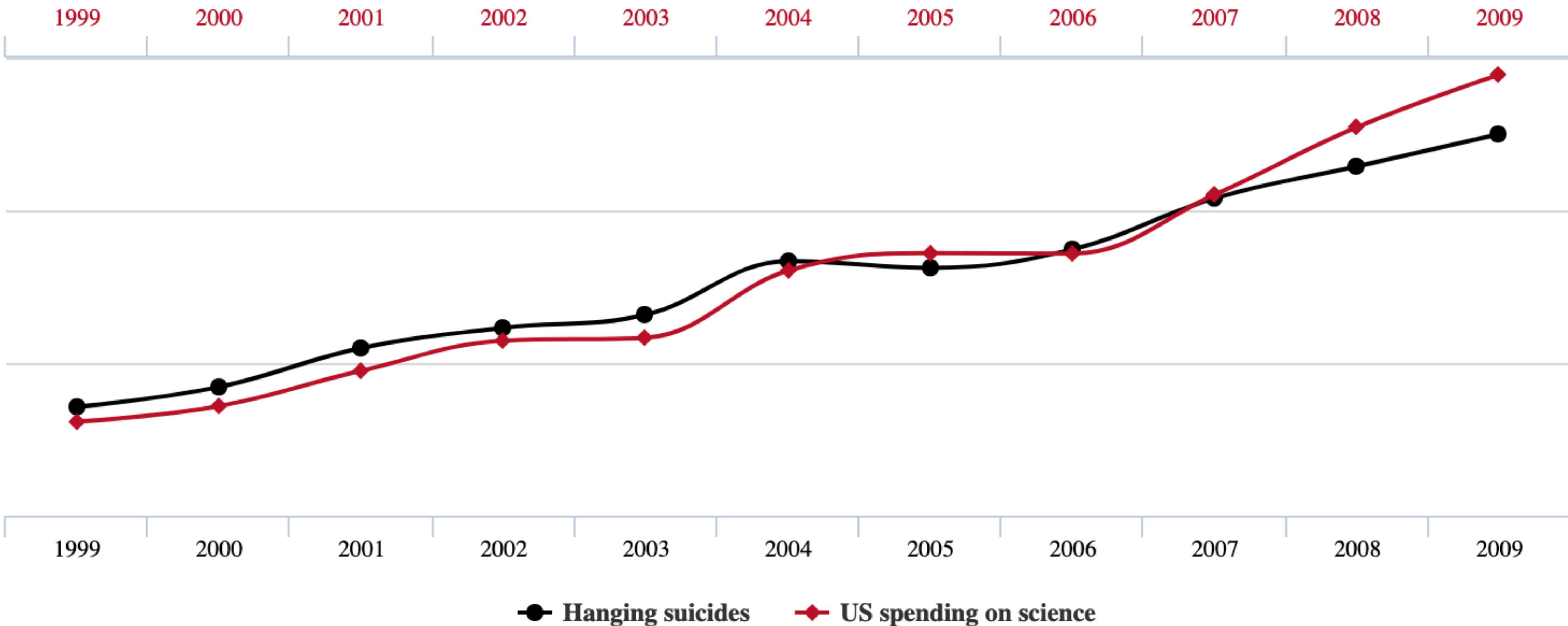
# Korelacija ≠ Kauzalnost

**US spending on science, space, and technology**

correlates with

**Suicides by hanging, strangulation and suffocation**

Correlation: 99.79% ( $r=0.99789126$ )



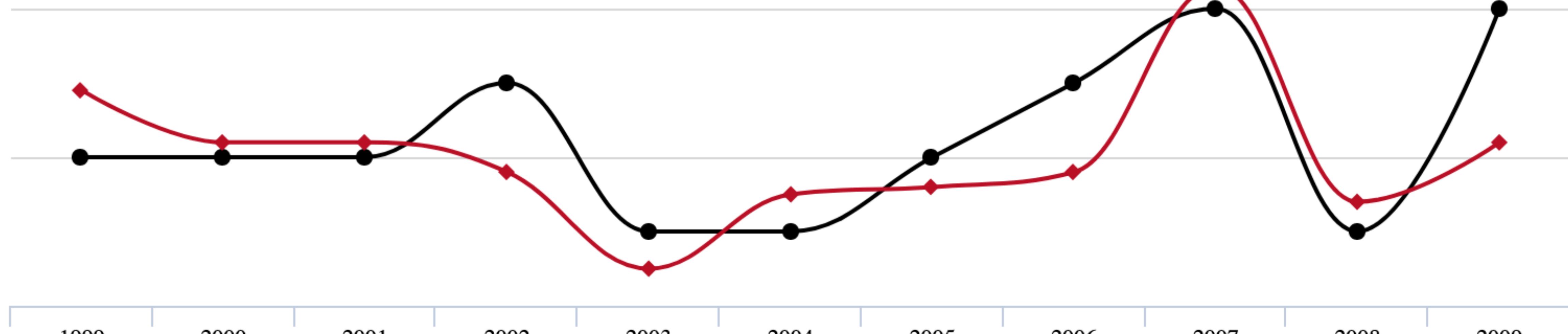
# Korelacija ≠ Kauzalnost

**Number of people who drowned by falling into a pool**

correlates with

**Films Nicolas Cage appeared in**

Correlation: 66.6% ( $r=0.666004$ )

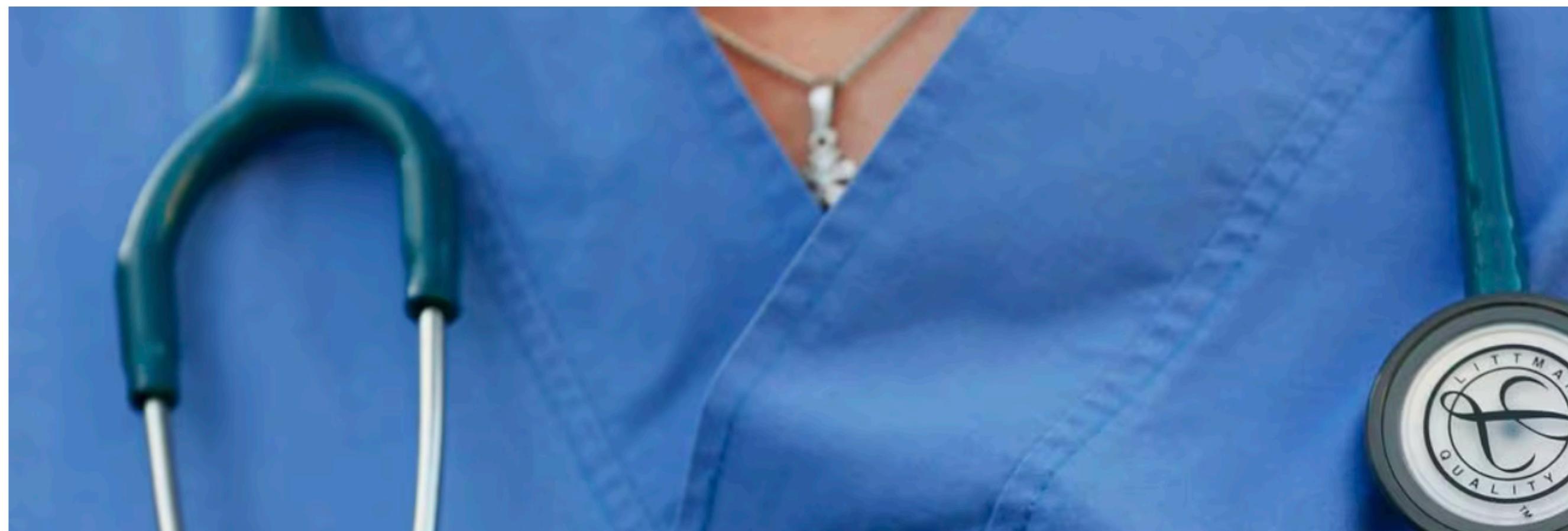


● Nicholas Cage

◆ Swimming pool drownings

# When doctors strike, fewer patients die

By Ryan Hoskins Globe Correspondent, February 9, 2016, 1:36 p.m.



## July effect

文 2 languages ▾

Article Talk

Read Edit View history Tools ▾

From Wikipedia, the free encyclopedia

"*Killing season*" redirects here. For other uses, see [Killing Season \(disambiguation\)](#).

The **July effect**, sometimes referred to as the **July phenomenon**, is a perceived but scientifically unfounded increase in the risk of medical errors and surgical complications that occurs in association with the time of year in which United States medical school graduates begin residencies.<sup>[1]</sup> A similar period in the United Kingdom is known as the **killing season** or, more specifically, **Black Wednesday**, referring to the first Wednesday in August when postgraduate trainees commence their rotations.

# Objašnjenje povezanosti



# Objašnjenje povezanosti



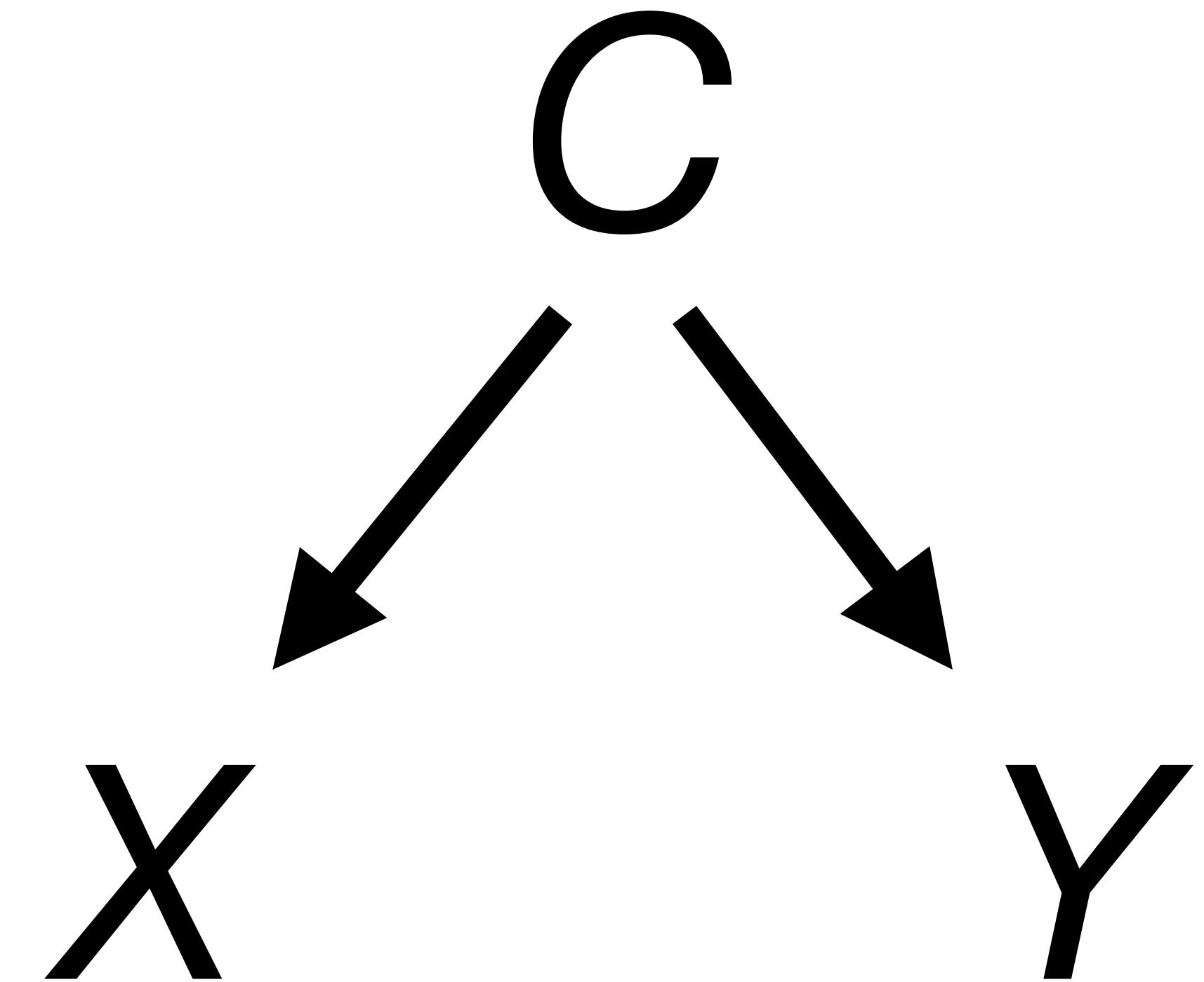
# Objašnjenje povezanosti

X

Y

Slučajno povezane

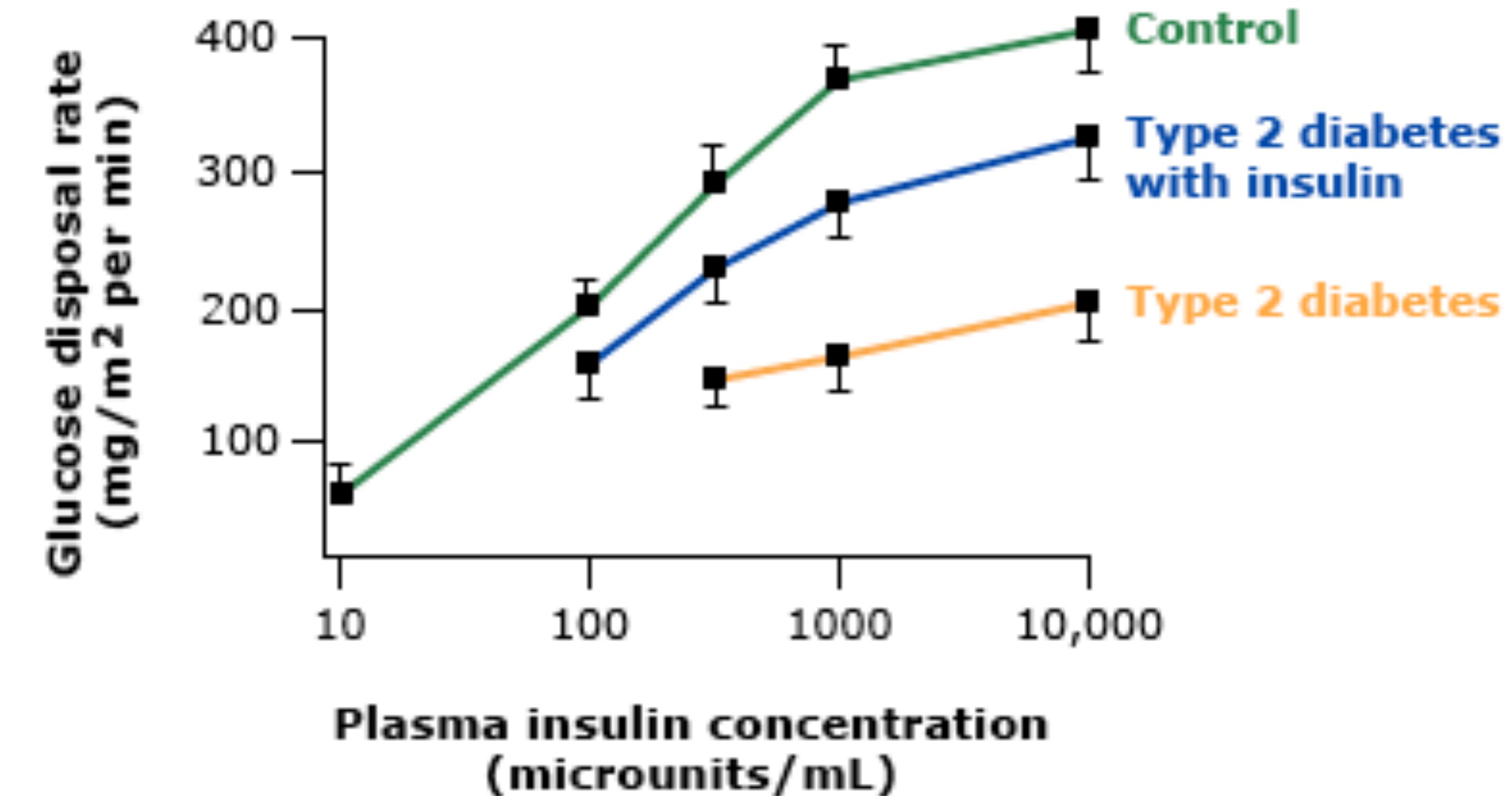
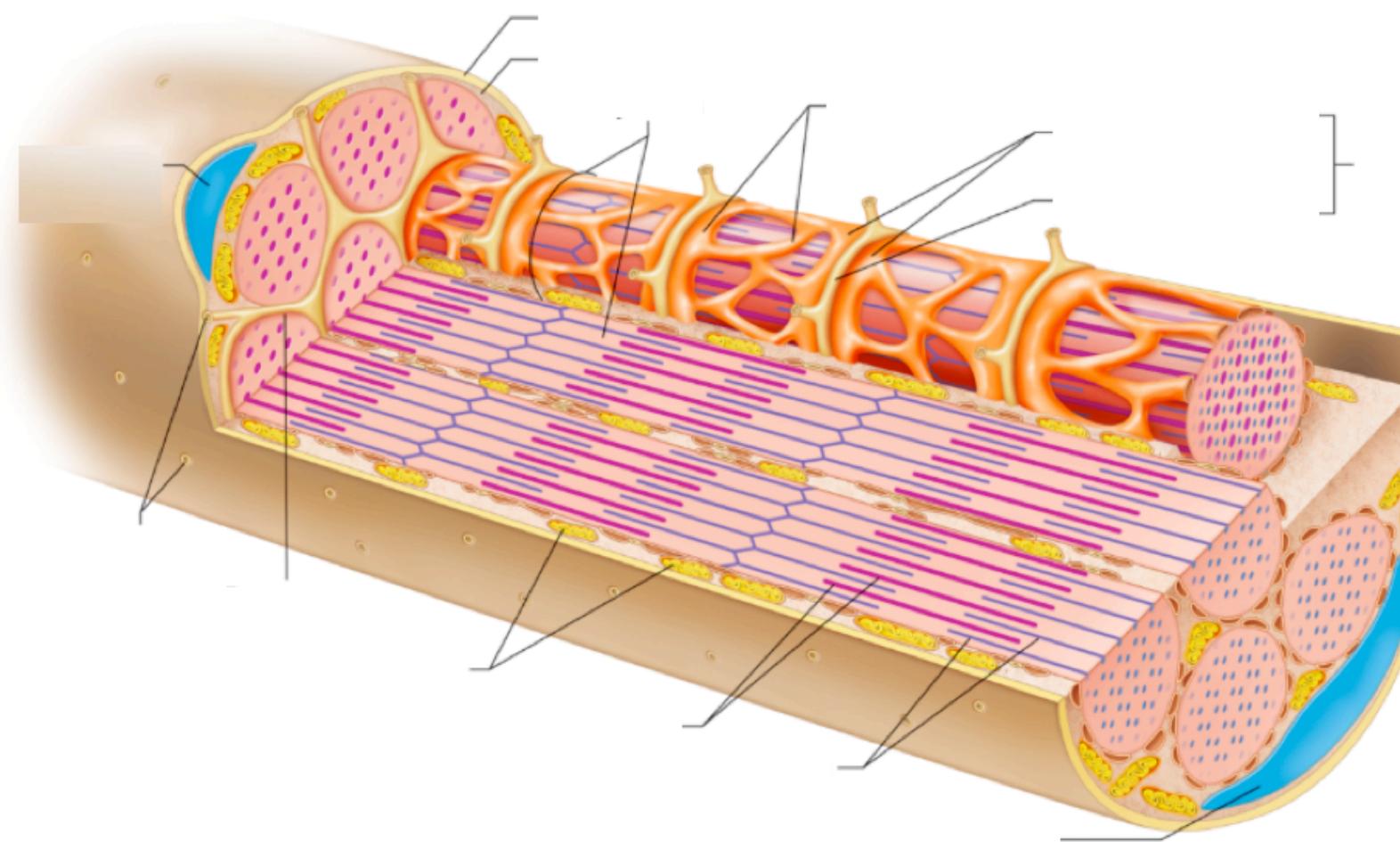
# Objašnjenje povezanosti



# Primer: lipidi i insulinska rezistencija

Borkman et al. (1993)

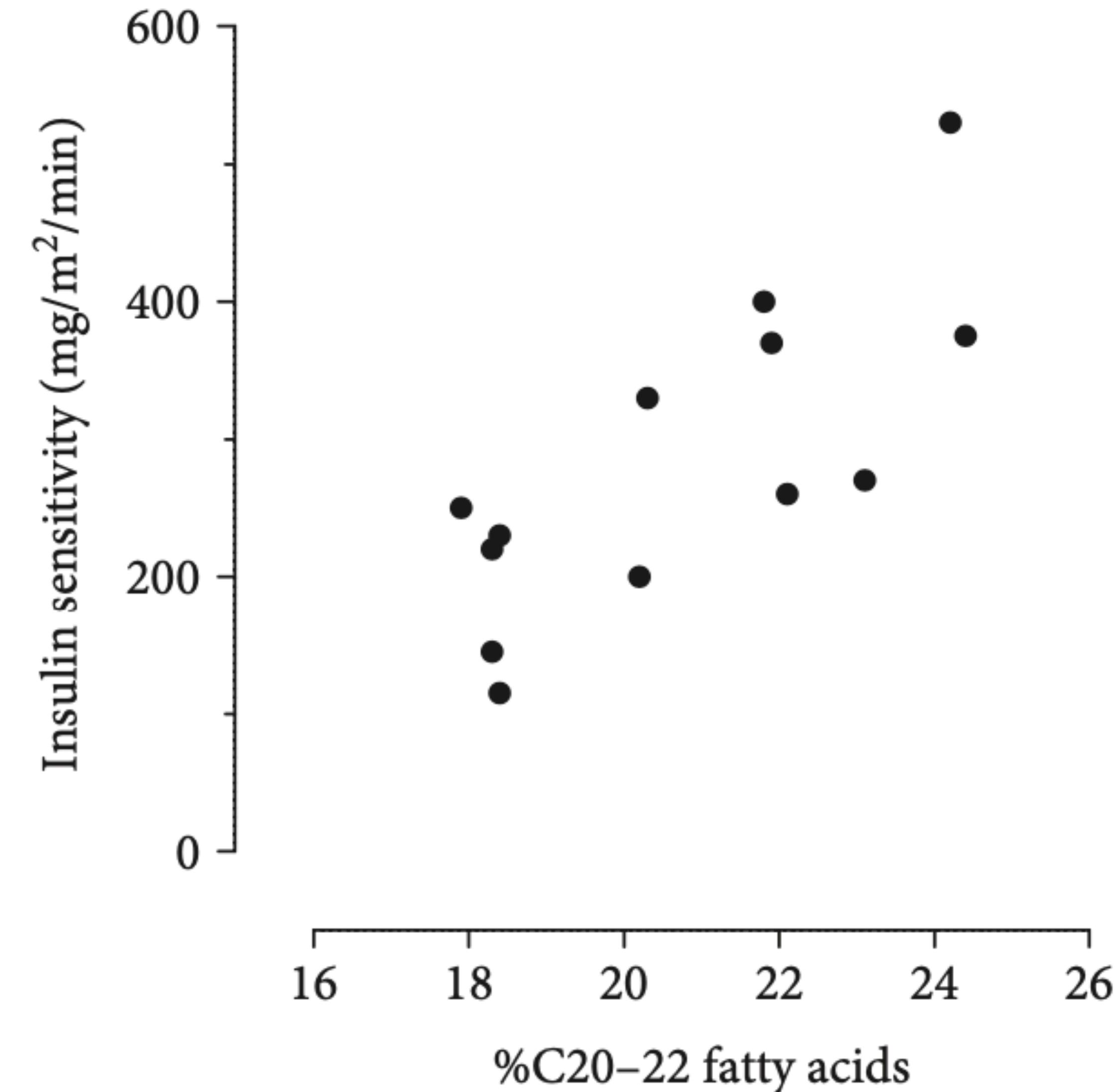
- **Hipoteza:** lipidni sastav membrane skeletnih mišićih ćelija utiče na insulinsku rezistenciju.



# Primer: lipidi i insulinska rezistencija

Borkman et al. (1993)

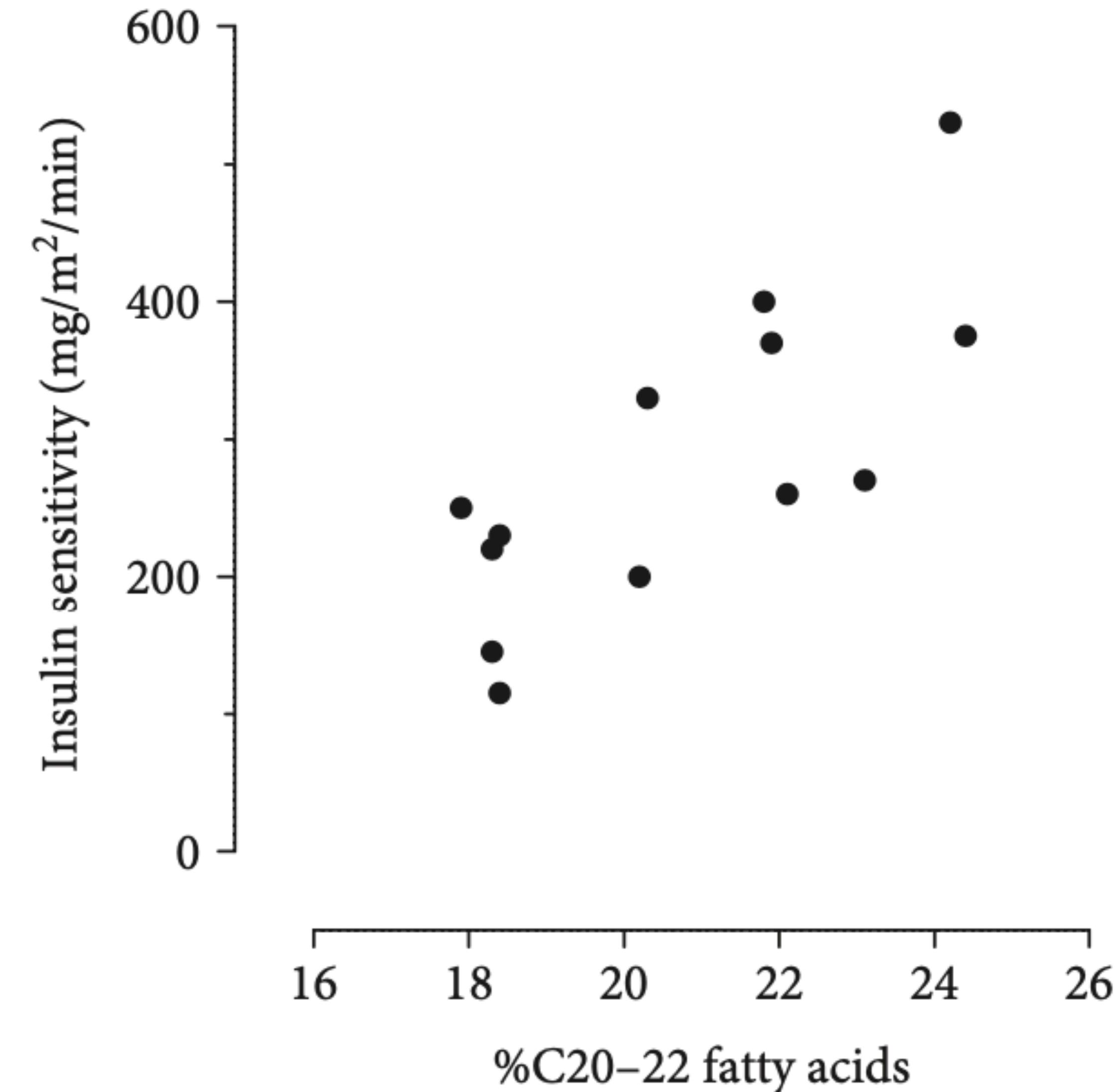
- Scatterplot (dijagram rasturanja)
- Kako da opišemo odnos ove dve varijable?

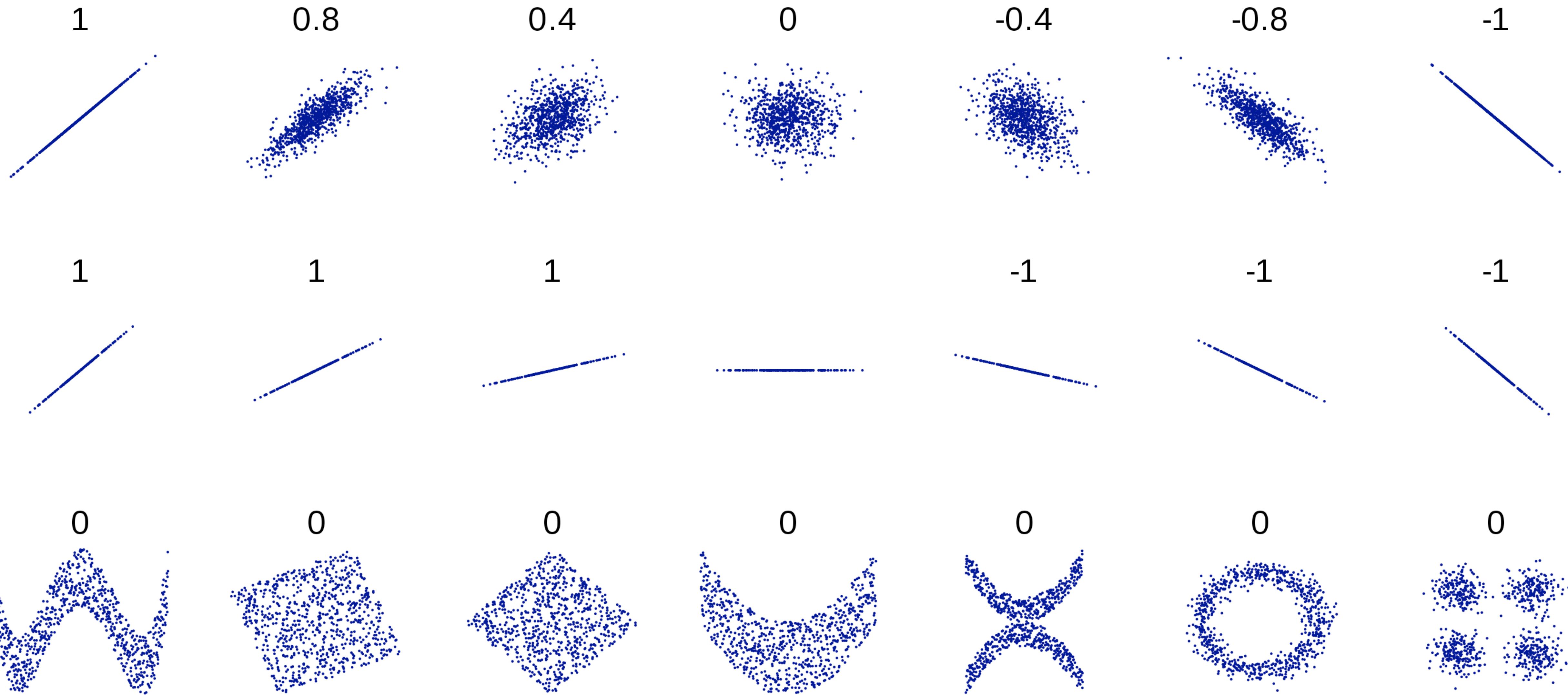


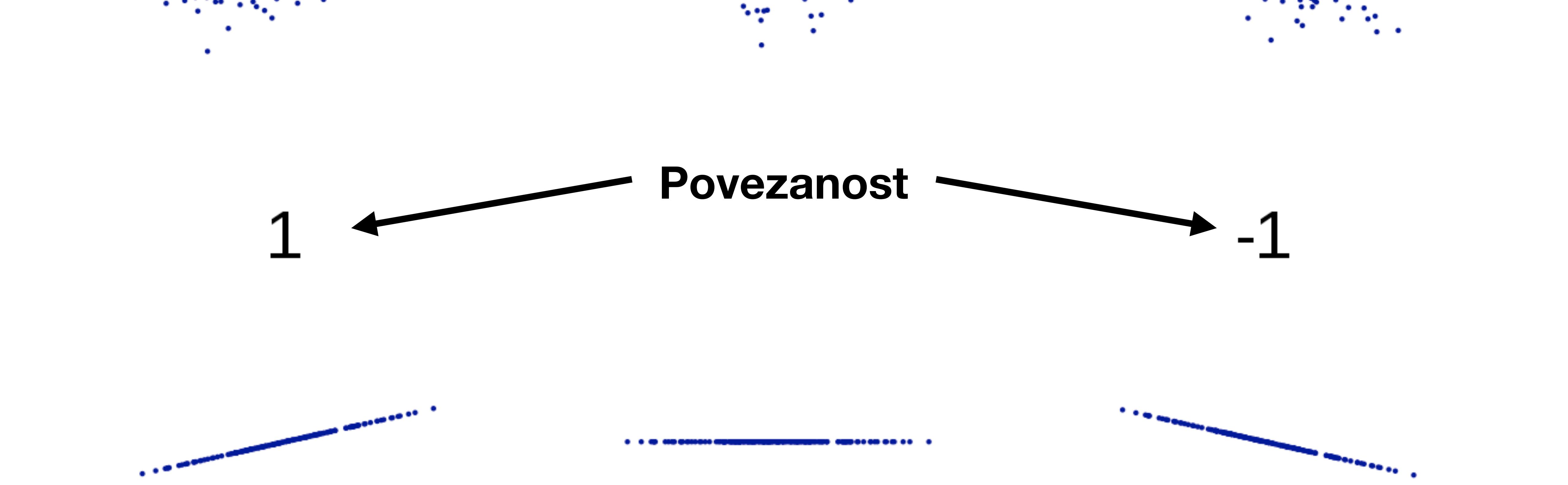
# Primer: lipidi i insulinska rezistencija

Borkman et al. (1993)

- Scatterplot (dijagram rasturanja)
- Kako da opišemo odnos ove dve varijable?
- **Korelacija**







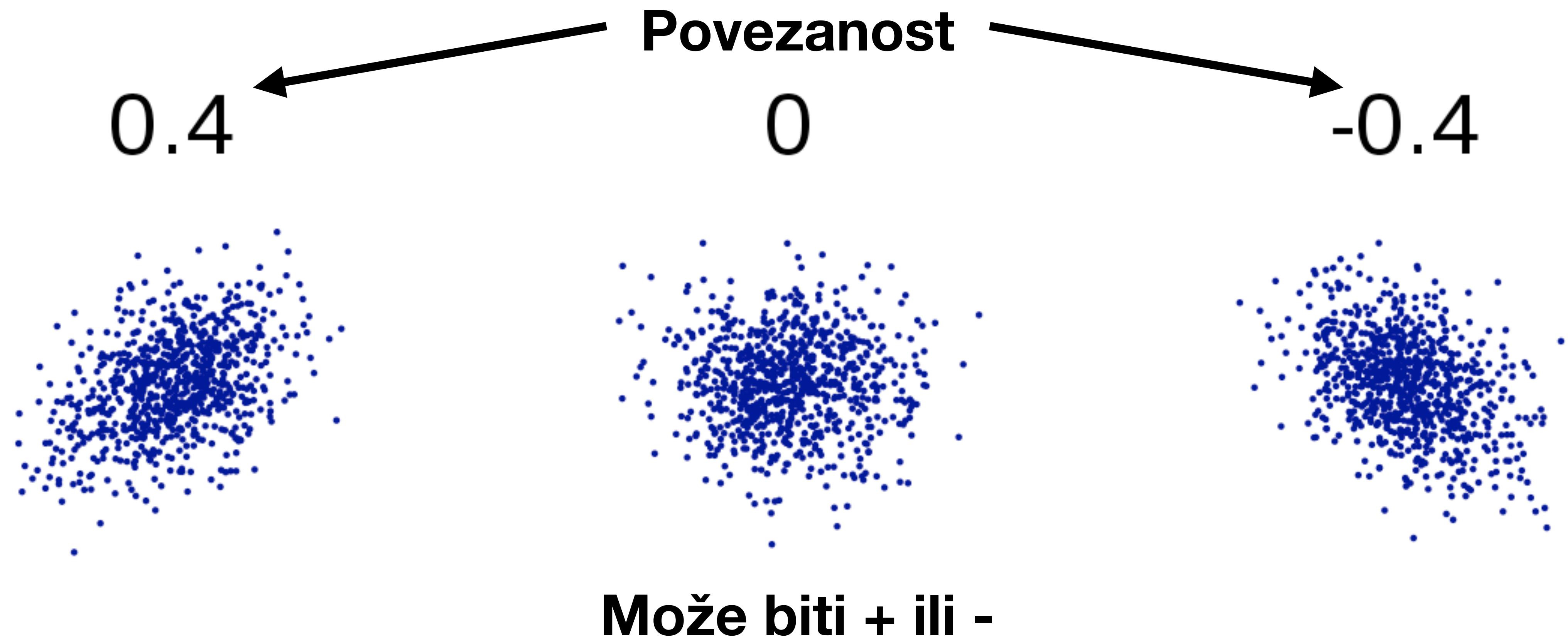
**Povezanost**

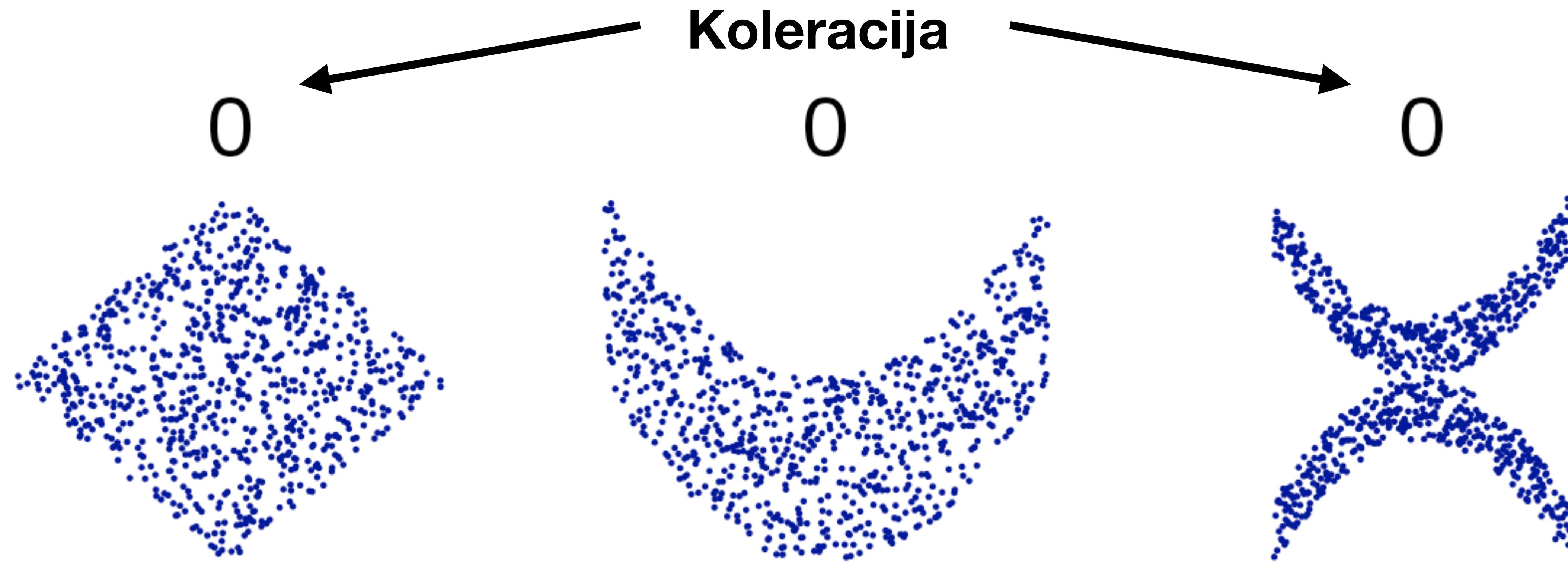
The background of the slide features a scatter plot with two distinct clusters of blue dots. A horizontal dashed blue line runs across the middle of the plot. Two thick black arrows point from the text labels '1' and '-1' towards this dashed line.

1

-1

**Može biti + ili -**

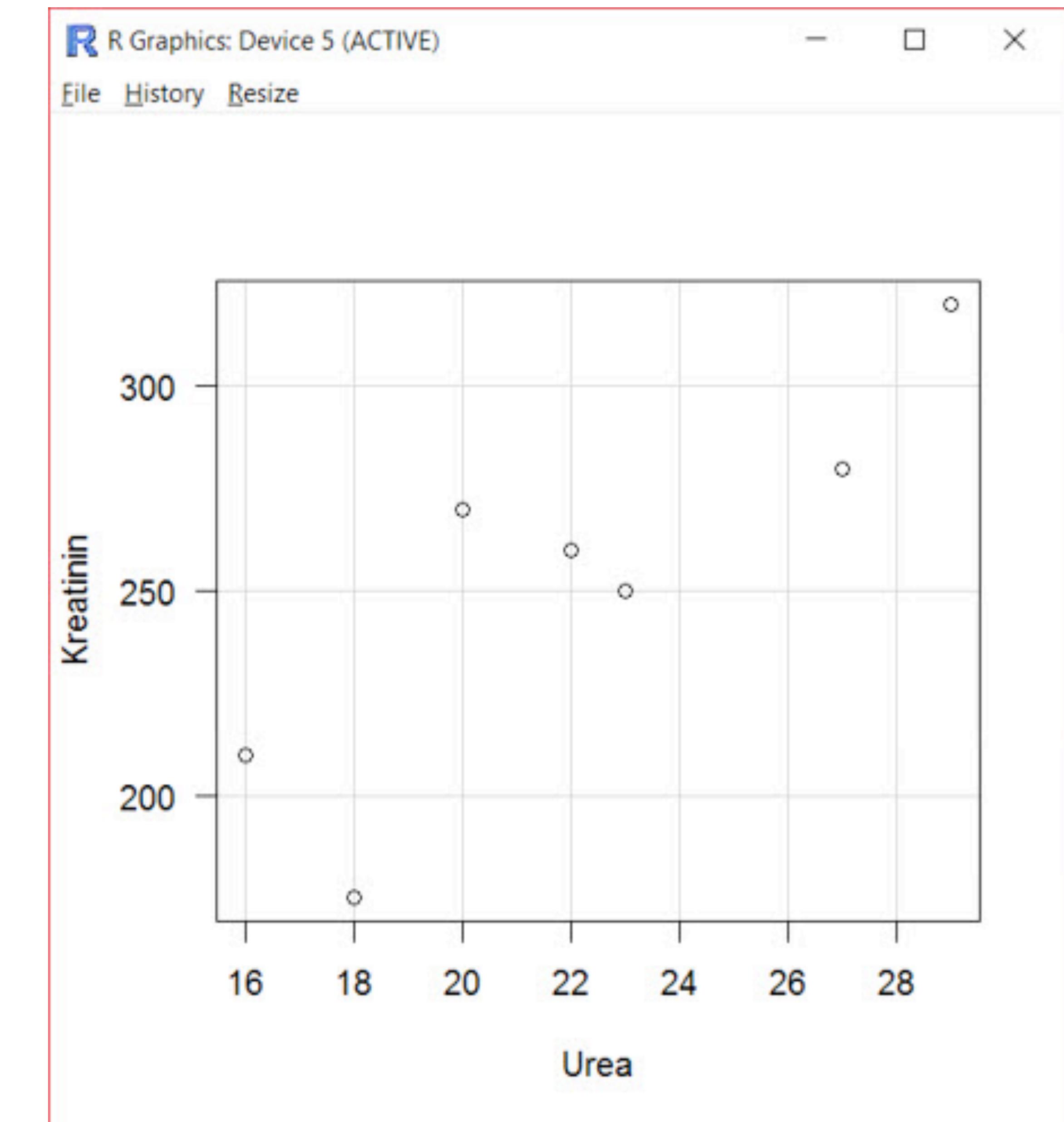




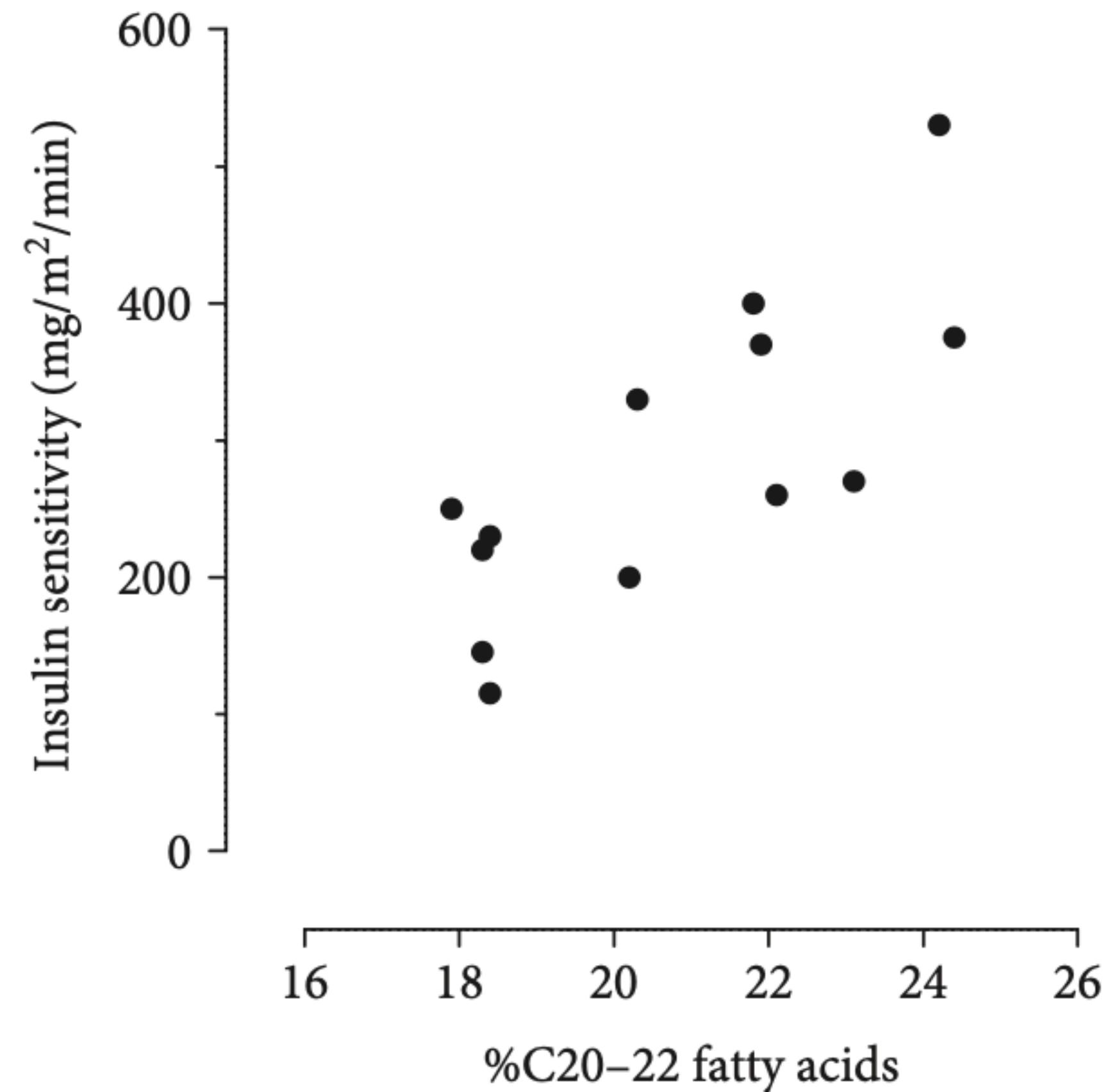
**Varijable mogu biti povezane i nelinearno**

# Povezanost Rezime

- **Smer** (pozitivan +, negativan -)
- **Stepen** (od -1 do 1 = “jačina”)
- **Oblik** (Linearan, nelinearan)
- Scatterplot (dijagram rasturanja)



**Napraviti dijagram rasturanja u kome će se vrednosti uree naneti na x osu, a vrednosti kreatinina na y osu.**



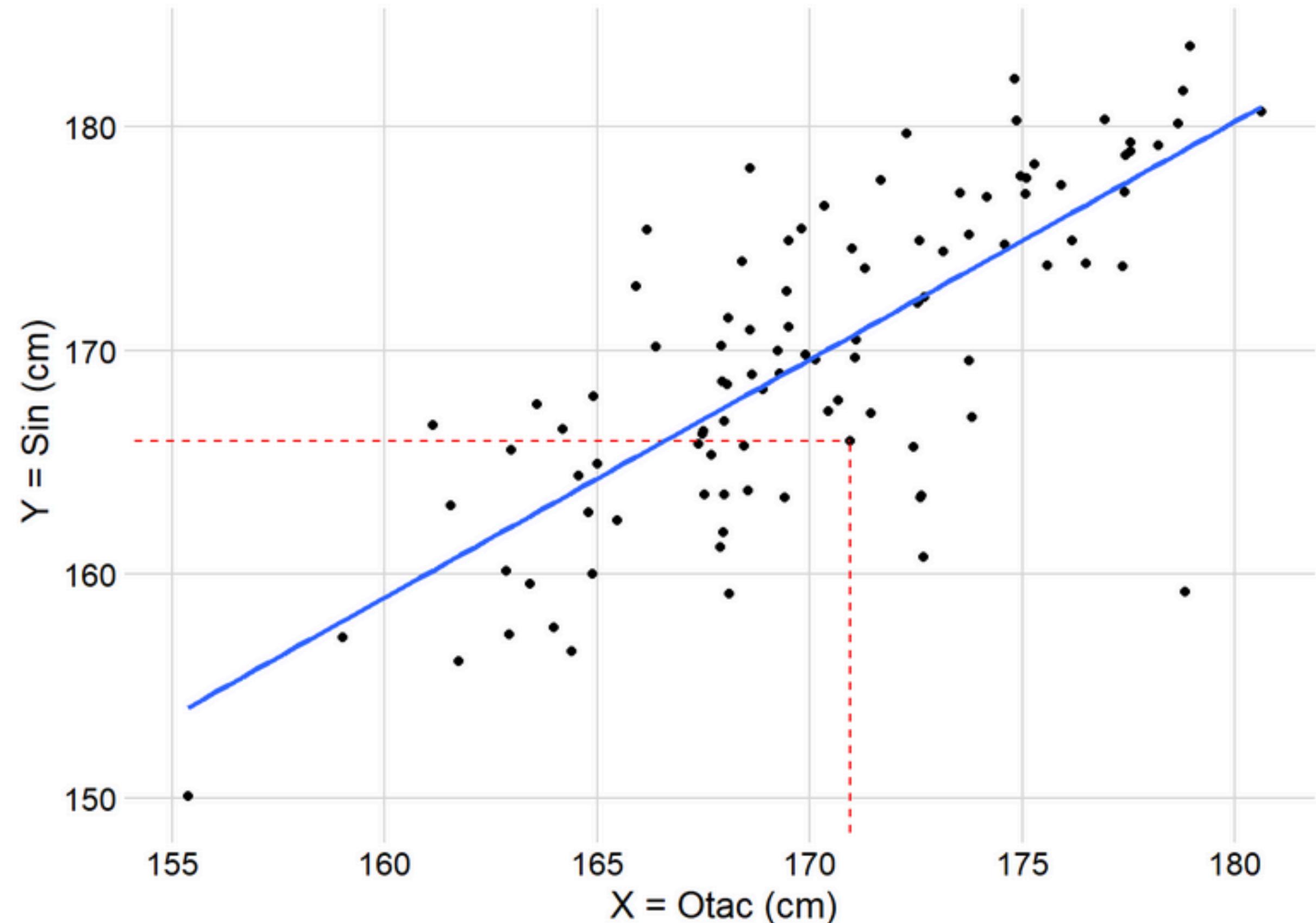
rb	Urea_mmol_L	Cr_umol_L
1	27	280
2	29	320
3	16	210
4	18	175
5	20	270
6	23	250
7	22	260

# Korelacioni koeficijent

**Koeficijent korelacije je statistika koja kvantificuje jačinu povezanosti varijabli.**

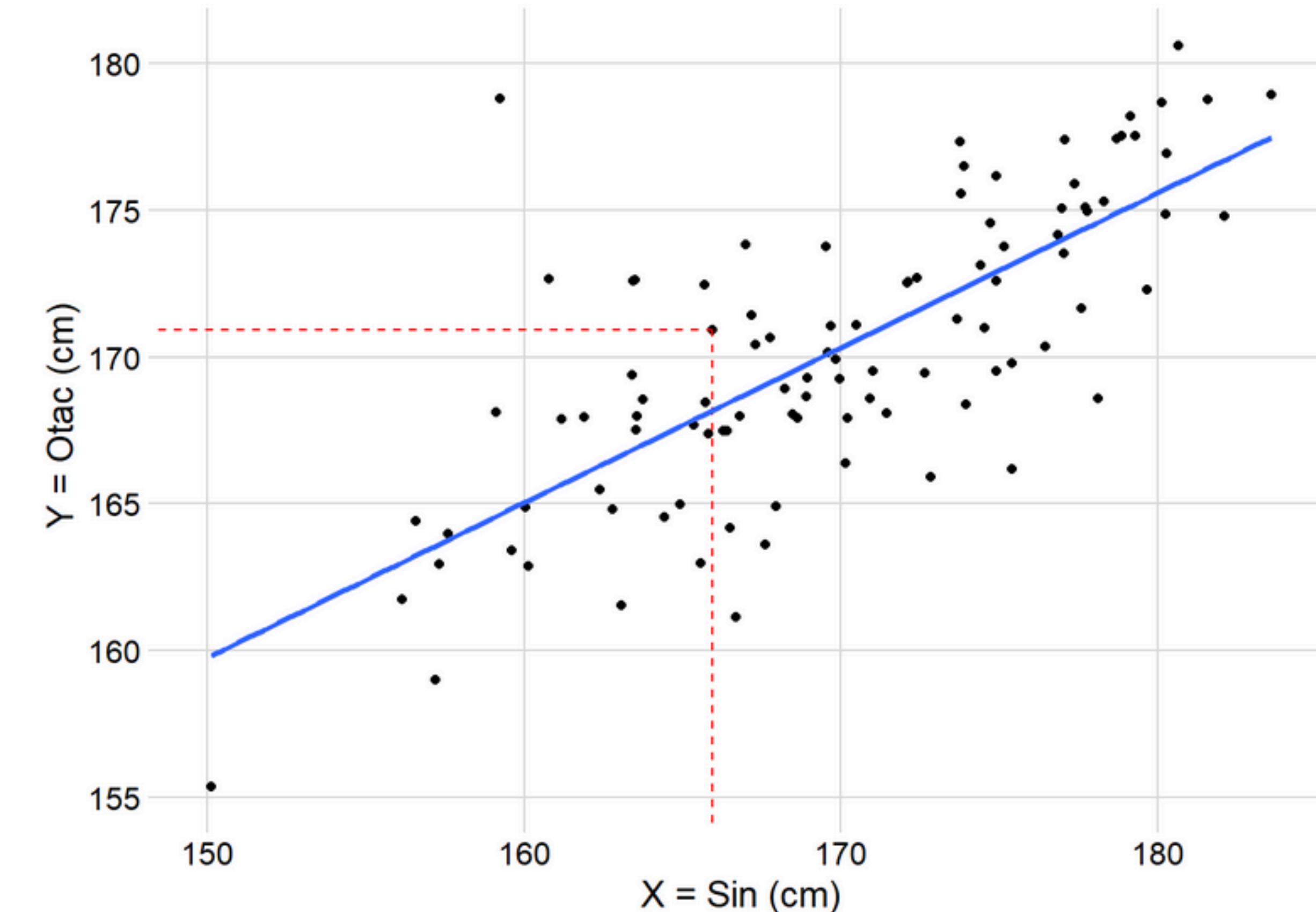
- Pearsonov koeficijent korelacije ( $r$ )
- Spearmanov koeficijent korelacije ranga ( $\rho$ )

# Simetričnost korelacija



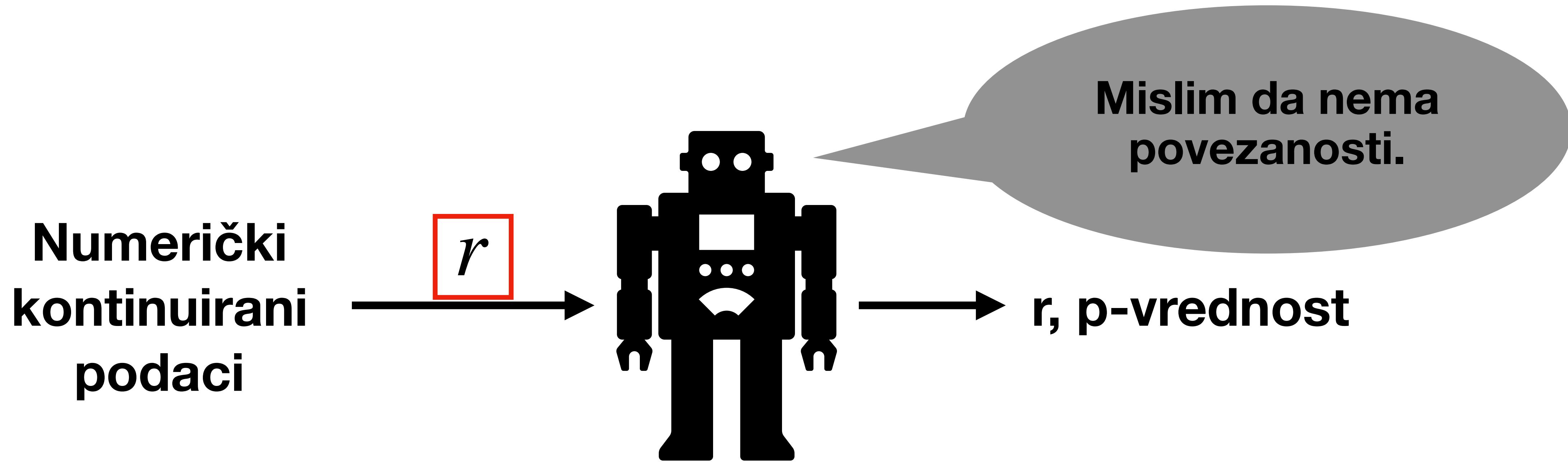
(a) Očevi predviđaju sinove

$$V_d = \frac{V_{otac} + V_{majka}}{2} \pm 13 \text{ cm}$$

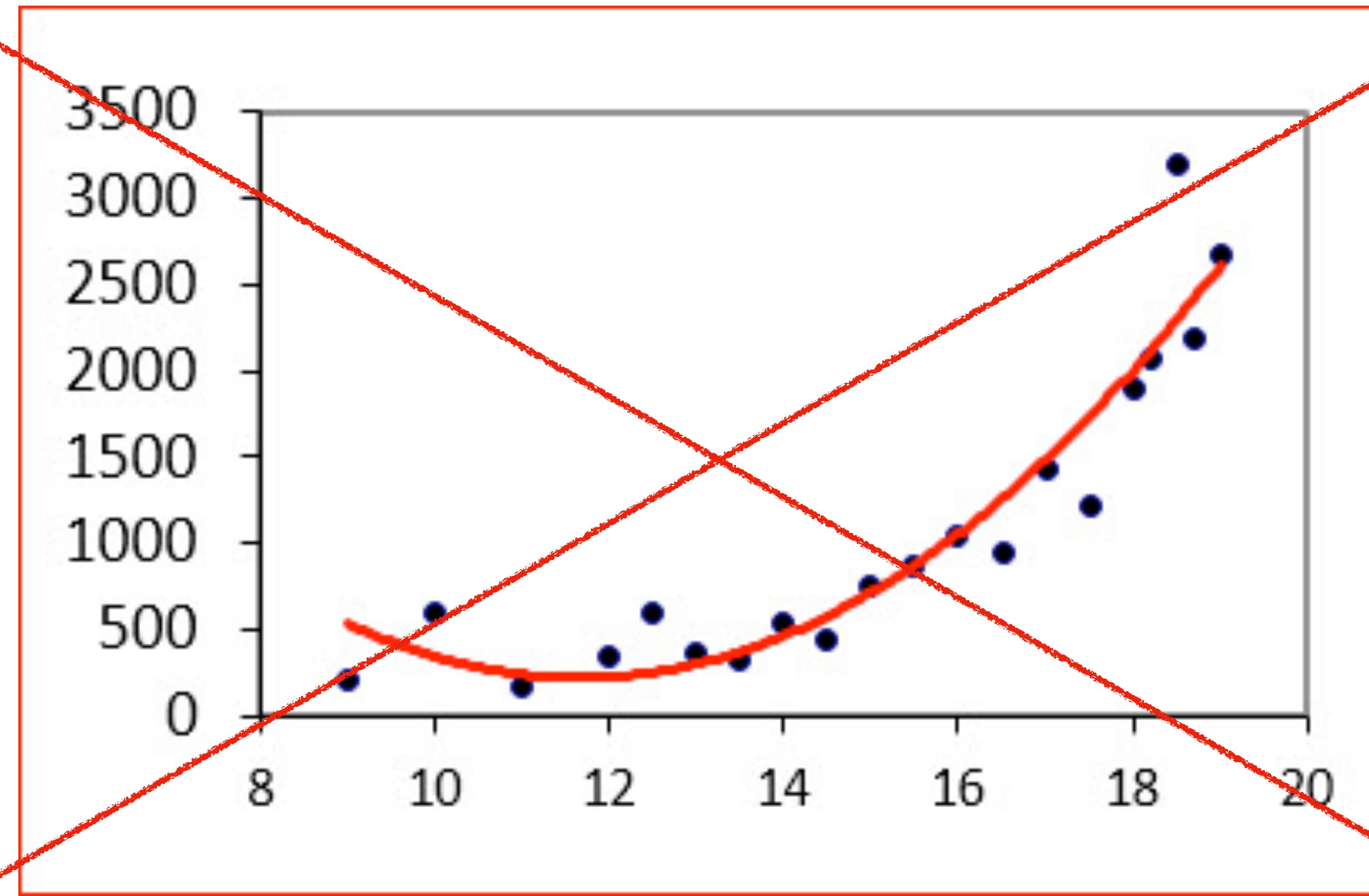


(b) Sinovi predviđaju očeve

# Pearsonov koeficijent korelaciјe

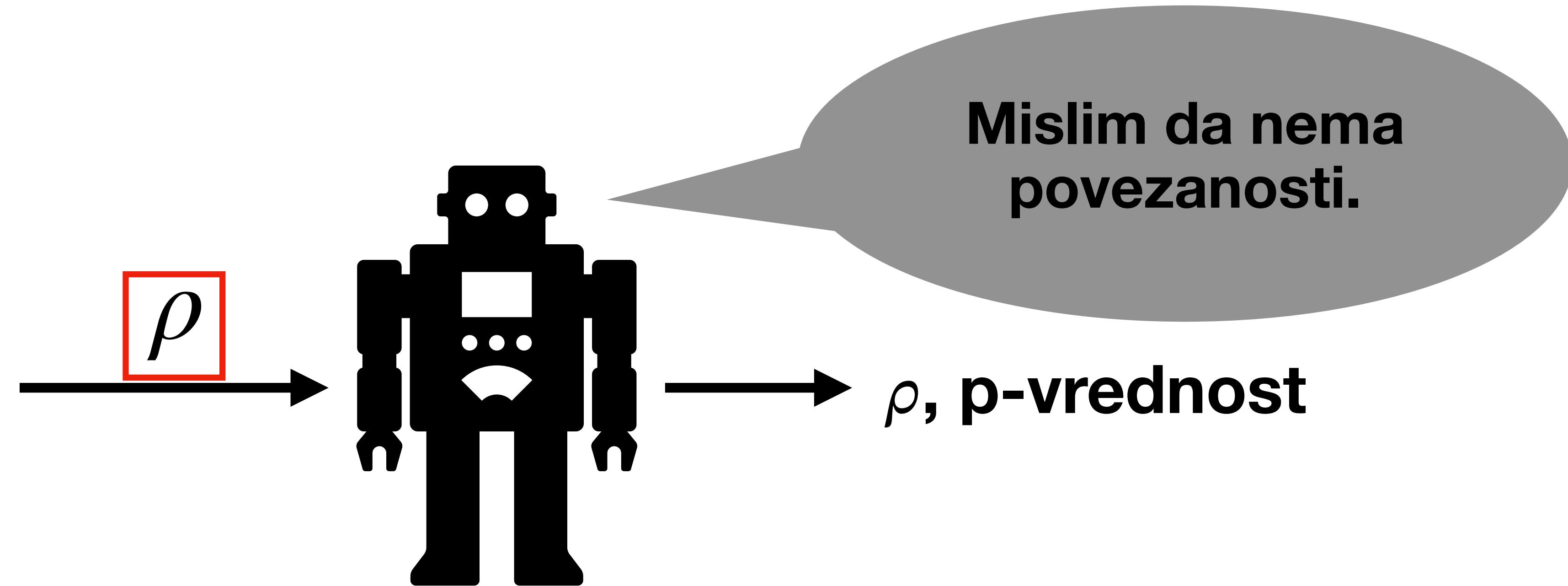


1. Nezavisne opservacije
2. Normalna raspodela u populaciji
3. Svaka varijabla mora da ima svog para
4. Linearan odnos varijabli (scatterplot)



# Spearmanov koeficijent korelaciјe rangova

Ordinalni  
podaci



1. **Varijабla nema normalnu raspodelu**
2. **Podaci su ordinalni ili rangovi**
3. **Svaka varijabla mora da ima svog para**
4. **Nelinearan odnos varijabli (scatterplot)**
5. **Monotonija\***

# Interpretacija koeficijenta korelacije

Koeficijent	Jačina povezanosti <i>(Interpretacija jednako važi i za vrednosti sa negativnim predznakom)</i>
0.70 – 1.00	Jaka povezanost
0.30 – 0.69	Osrednja povezanost
< 0.30	Slaba povezanost
oko 0.00	Nema povezanosti

Ako je  $p \leq 0.05$ ? Ako je  $p > 0.05$ ?

**Odrediti koeficijent korelacije  
između uree i kreatinina.**

# Odrediti koeficijent korelacija između uree i kreatinina.

- 1. Proveriti normalnost (numerical summaries)**
- 2. Pearson / Spearman ?**

```
> numSummary(Dataset[,c("Kreatinin", "Urea")], statistics=c("mean", "sd",
+ "quantiles", "cv"), quantiles=c(0,.25,.5,.75,1))
      mean        sd        cv   0% 25% 50% 75% 100% n
Kreatinin 252.14286 47.421615 0.1880744 175 230 260 275 320 7
Urea       22.14286  4.670067 0.2109062   16   19   22   25   29 7
```

Statistical analysis Graphs and tables Tools Help Original menu

Discrete variables

Continuous variables

Nonparametric tests

Survival analysis

Accuracy of diagnostic test

Matched-pair analysis

Metaanalysis and metaregression

Calculate sample size

-16)

```
, c("Kreatinin", "Urea")],  
) , quantiles=c(0,.25,.5,.
```

```
(0,10), main="Sample")
```

```
####
```

Model: <No active model>

Numerical summaries

Smirnov-Grubbs test for outliers

Kolmogorov-Smirnov test for normal distribution

Confidence interval for a mean

Single-sample t-test

Two-variances F-test

Two-sample t-test

Paired t-test

Bartlett's test

One-way ANOVA

Repeated-measures ANOVA

Multi-way ANOVA

ANCOVA

Test for Pearson's correlation

Linear regression

## Test for Pearson's correlation

Click pressing Ctrl key to select multiple variables

Variables (pick two)

Kreatinin

rb

Urea

Alternative Hypothesis

Two-sided

Correlation < 0

Correlation > 0

Condition to limit samples for analysis. Ex1. age>50

<all valid cases>

<

Help

Reset

OK

**Da li postoji povezanost  
depresije i sistolne arterijske  
tenzije? Testirati za nivo  
značajnosti 0.05.**

rb	skor	ta
1	23	139
2	19	109
3	26	113
4	23	128
5	19	124
6	17	105
7	23	116
8	26	135
9	20	120
10	19	124

Discrete variables

Continuous variables

Nonparametric tests

Survival analysis

Accuracy of diagnostic test

Matched-pair analysis

Metaanalysis and metaregression

Calculate sample size

Model:  <No active model>

Mann-Whitney U test

Wilcoxon's signed rank test

Kruskal-Wallis test

Friedman test

Jonckheere-Terpstra test

Spearman's rank correlation test

 Spearman's rank correlation test

Click pressing Ctrl key to select multiple variables

Variables (pick two)

Depresivnost

rb

TA

Alternative Hypothesis

 Two-sided Correlation < 0 Correlation > 0

Method

 Spearman Kendall

Condition to limit samples for analysis. Ex1. age&gt;50 &amp;

&lt;all valid cases&gt;

 Help Reset OK

## Output

```
> (res <- cor.test(Dataset$Depresivnost, Dataset$TA, alternative="two.sided",
+   method="spearman"))
```

Spearman's rank correlation rho

data: Dataset\$Depresivnost and Dataset\$TA

S = 89.695, p-value = 0.1849

alternative hypothesis: true rho is not equal to 0

sample estimates:

rho
0.4563926

vrednost Spearmanovog  
koeficijenta korelacije

statistička značajnost  
koeficijenta korelacije

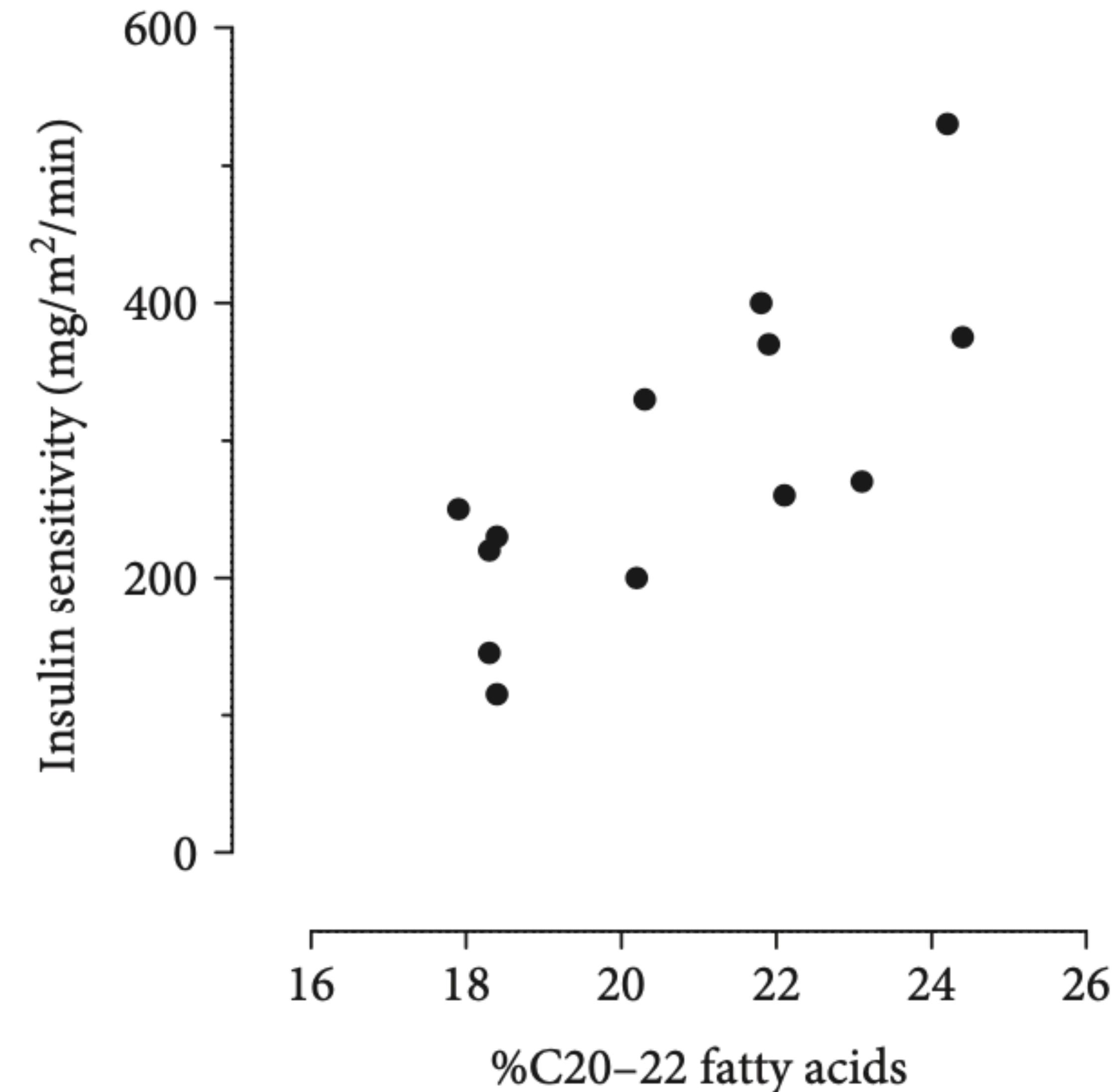
Preuzimanjem **baza DZ** odgovori na sledeća istraživačka pitanja:

1. Da li postoji povezanost između **starosti i ukupnog holesterola?**
2. Da li postoji povezanost između **triglicerida i stepena uhranjenosti?**

# Primer: lipidi i insulinska rezistencija (2)

Borkman et al. (1993)

CORRELATION	
r	0.77
95% CI	0.3803 to 0.9275
$r^2$	0.5929
P (two-tailed)	0.0021
Number of XY pairs	13



Opisali smo povezanost

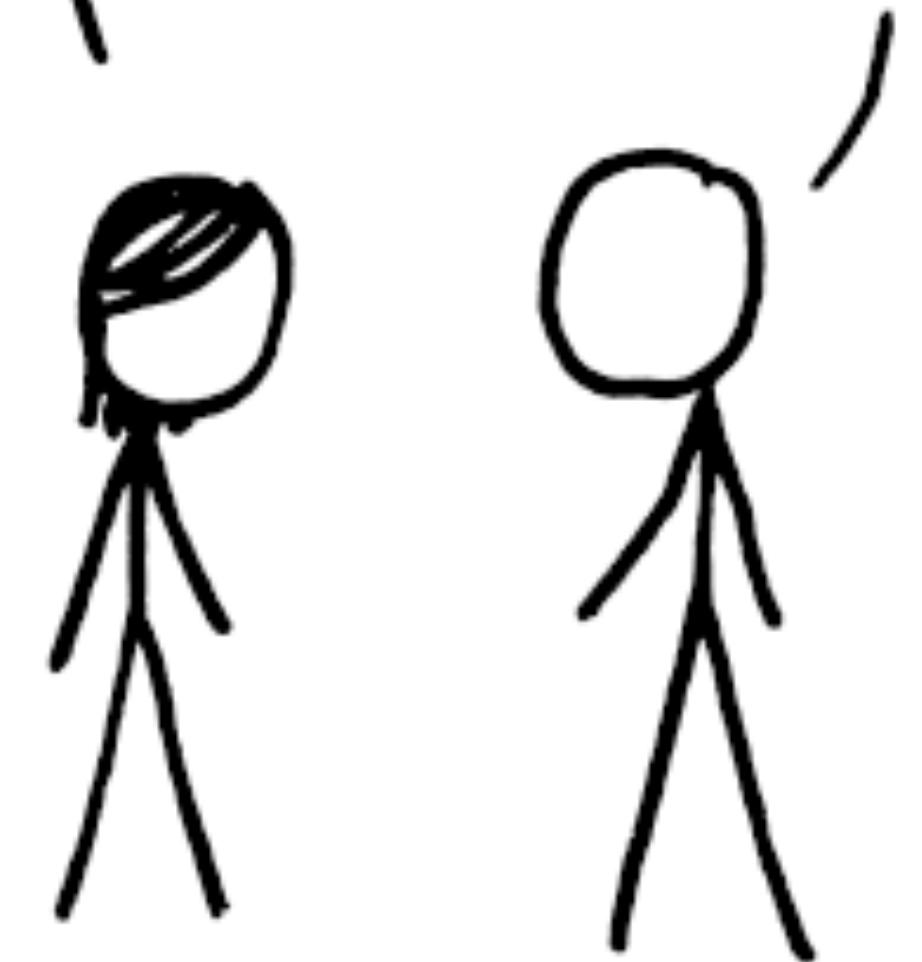
I USED TO THINK  
CORRELATION IMPLIED  
CAUSATION.



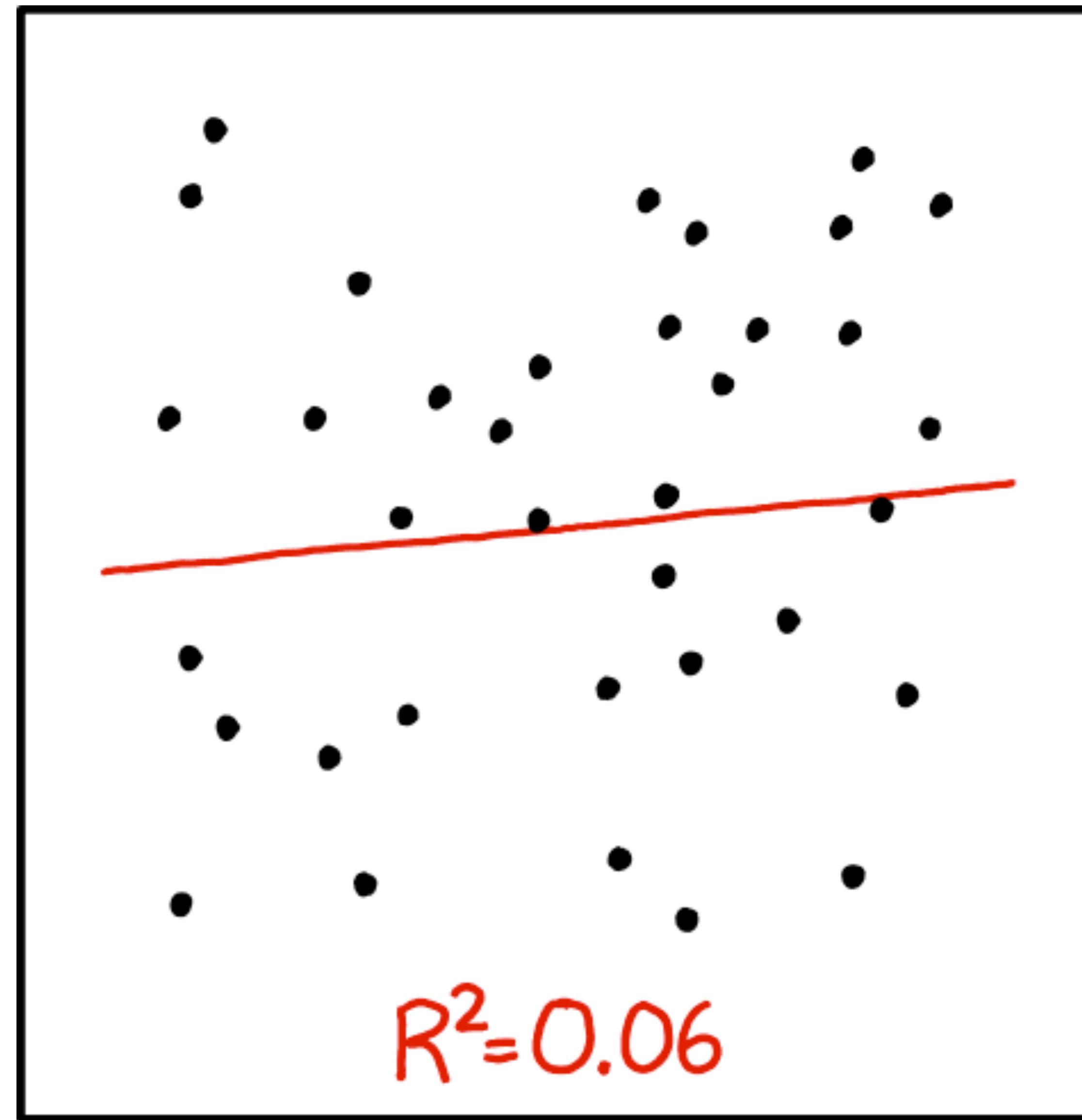
THEN I TOOK A  
STATISTICS CLASS.  
NOW I DON'T.

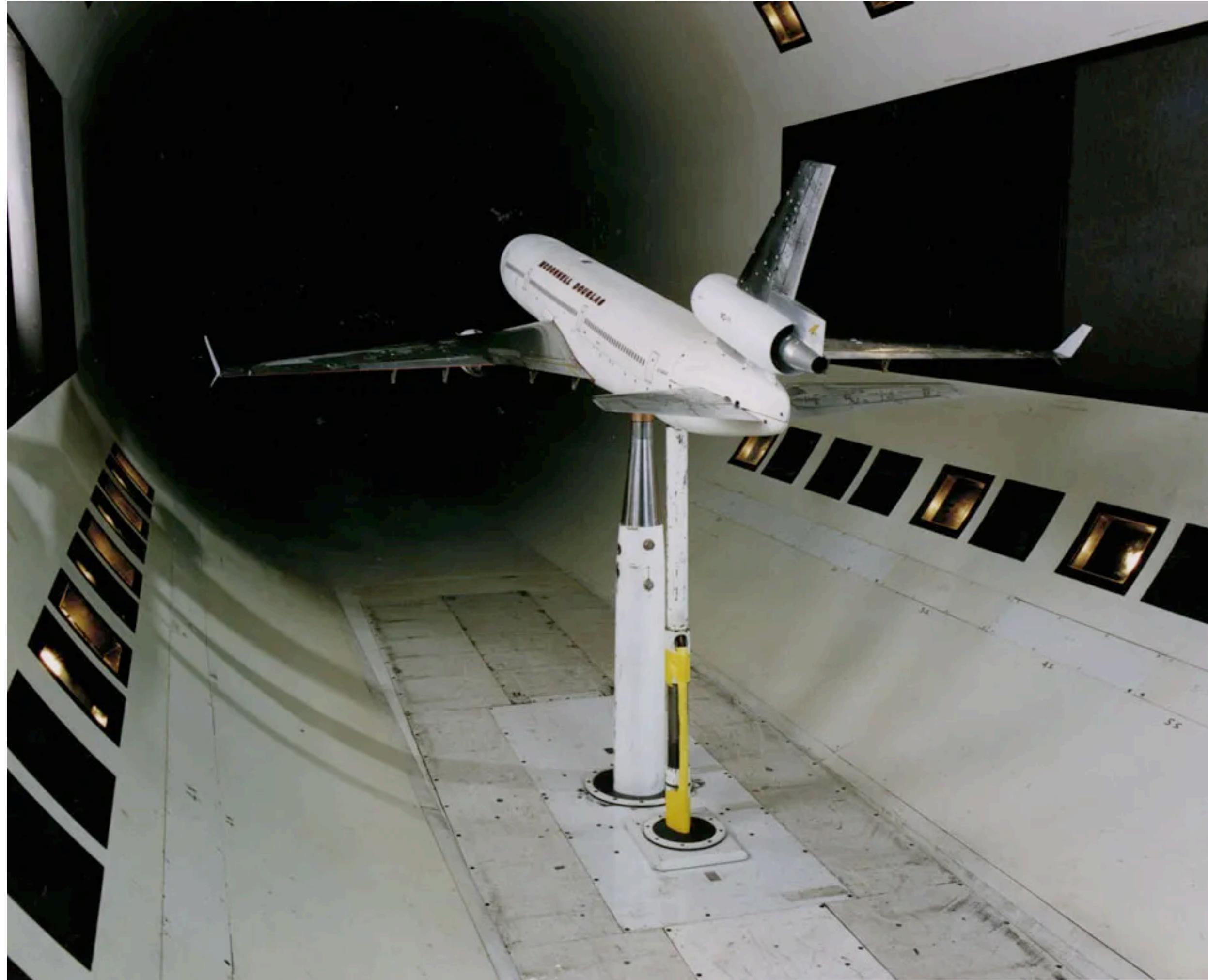


SOUNDS LIKE THE  
CLASS HELPED.  
WELL, MAYBE.



# Regresija





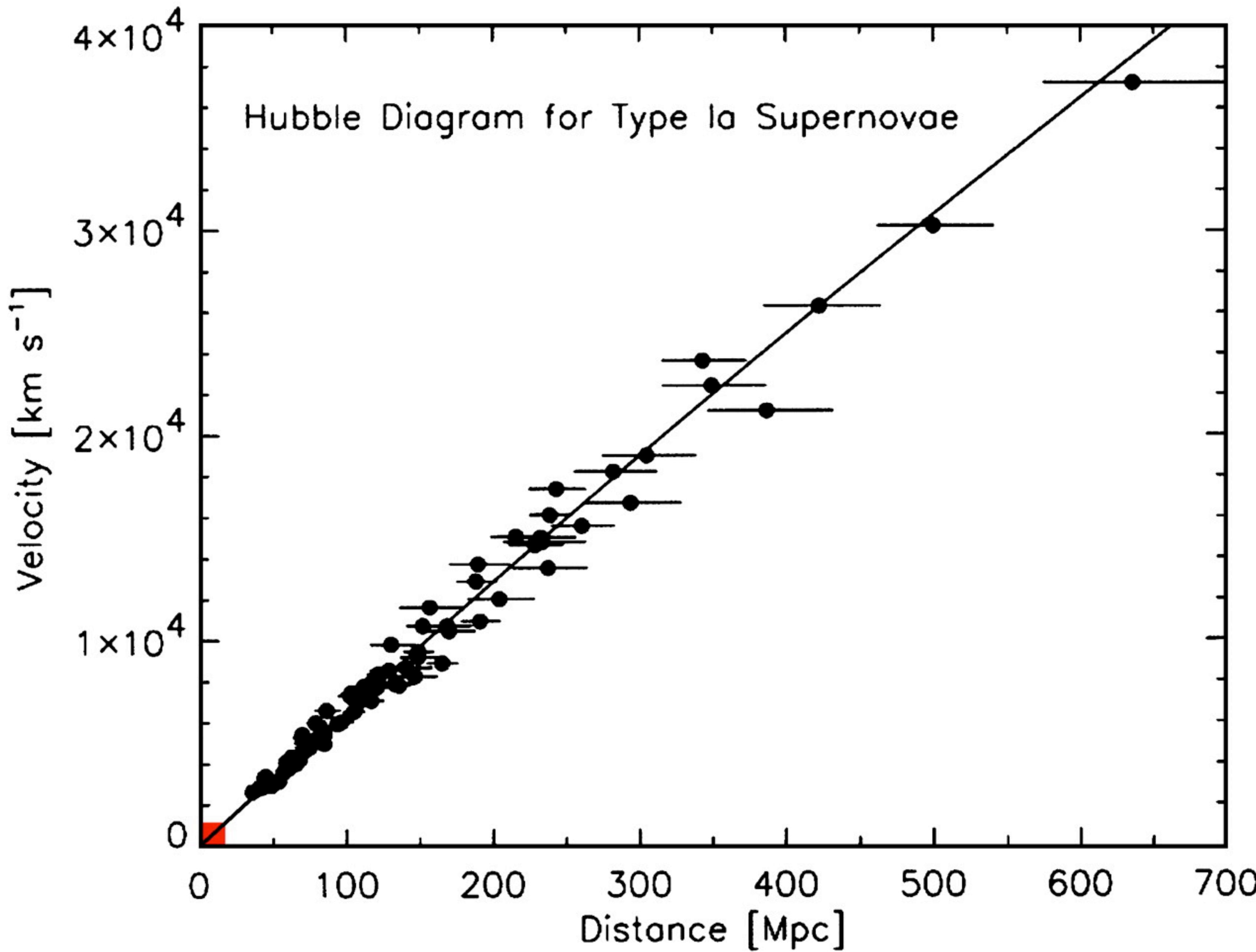
Matematički model

$$y = \alpha + \beta x$$

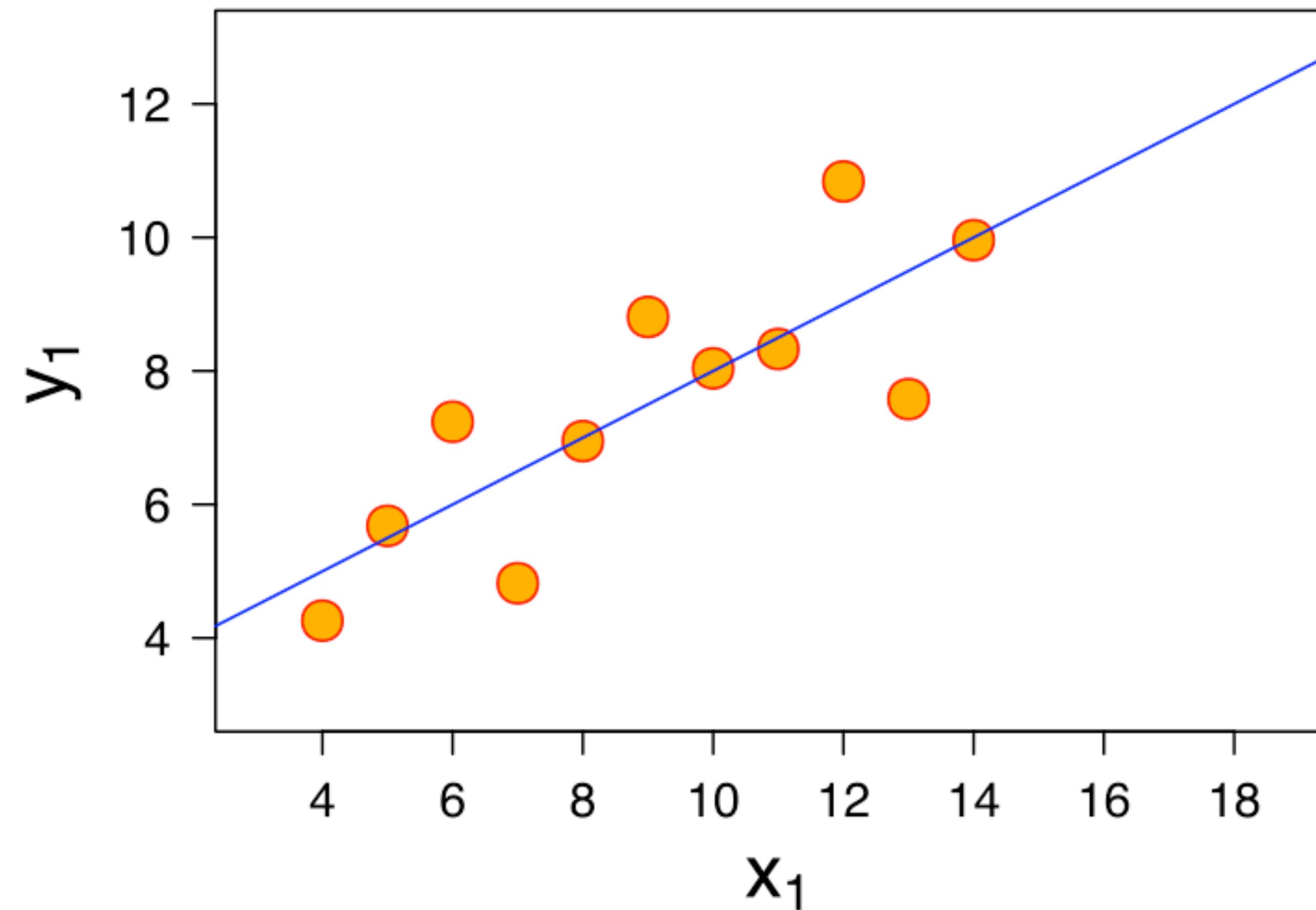


Stvarnost

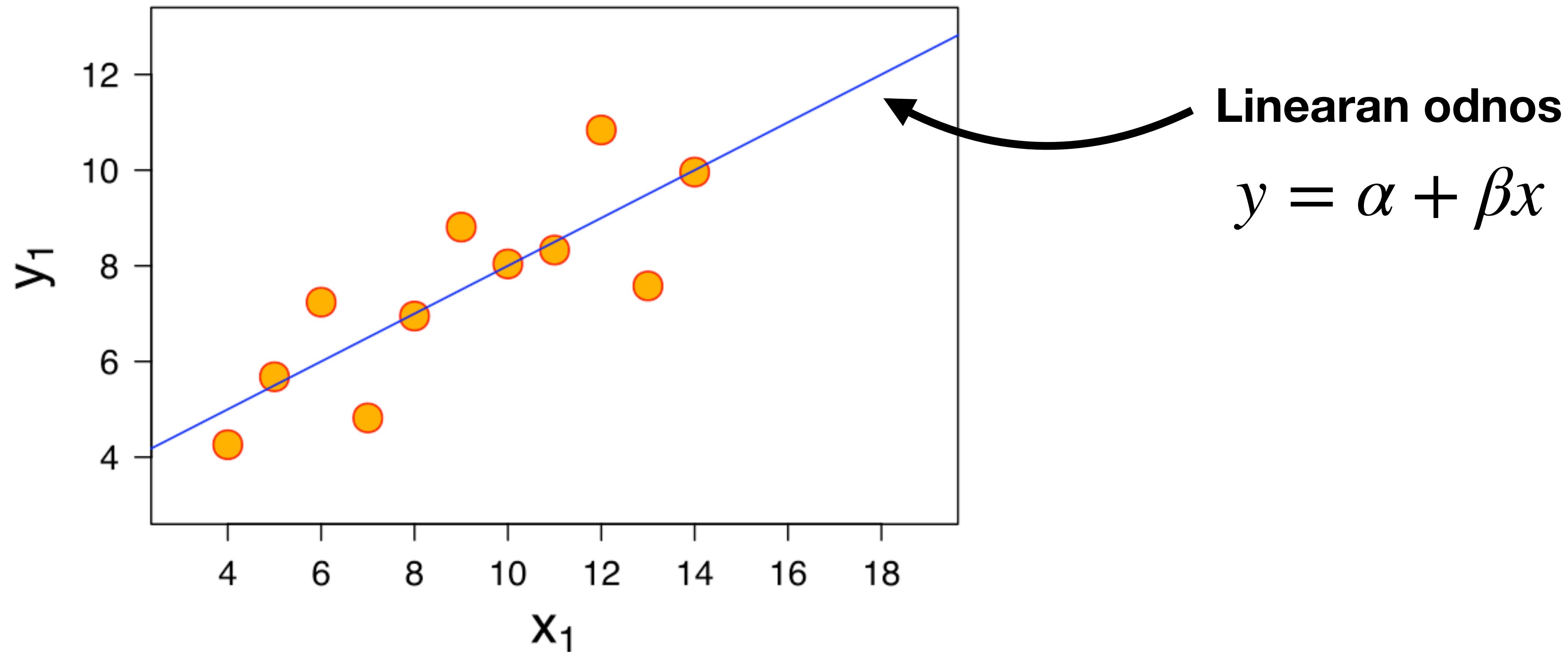
$y = \text{komplikovano...}$



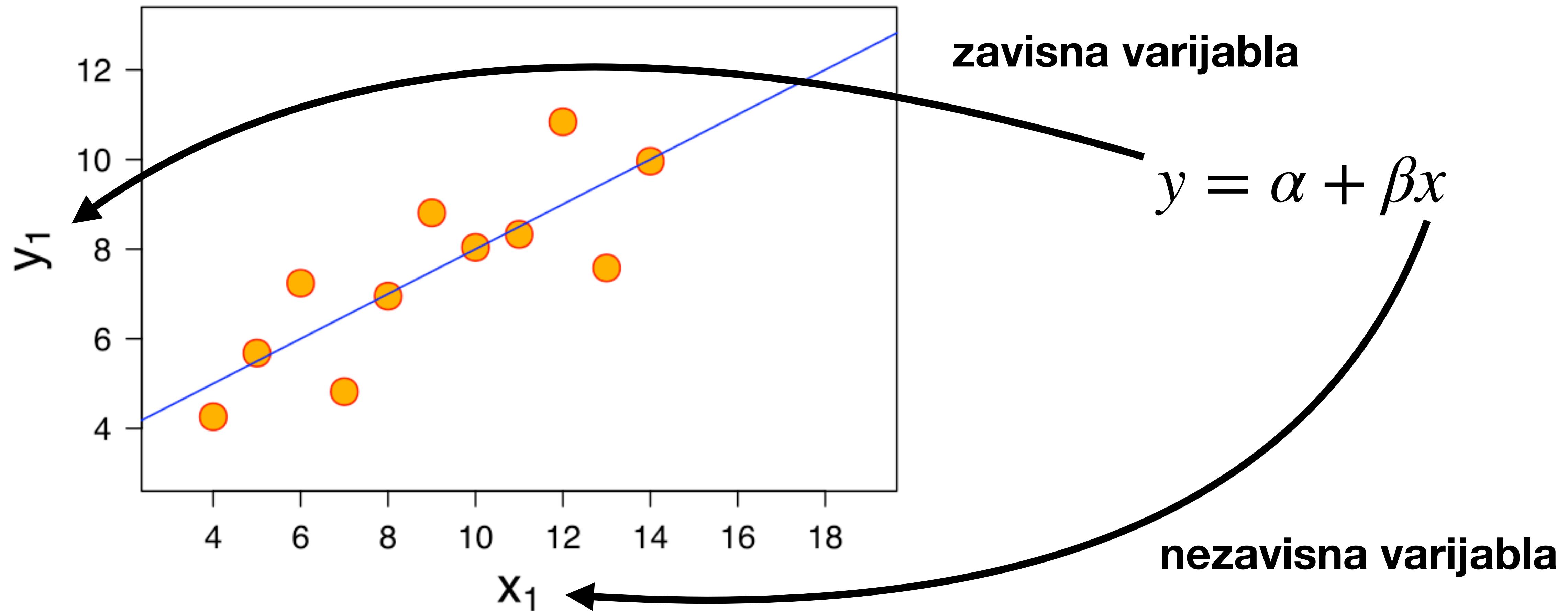
# Anatomija regresione jednačine



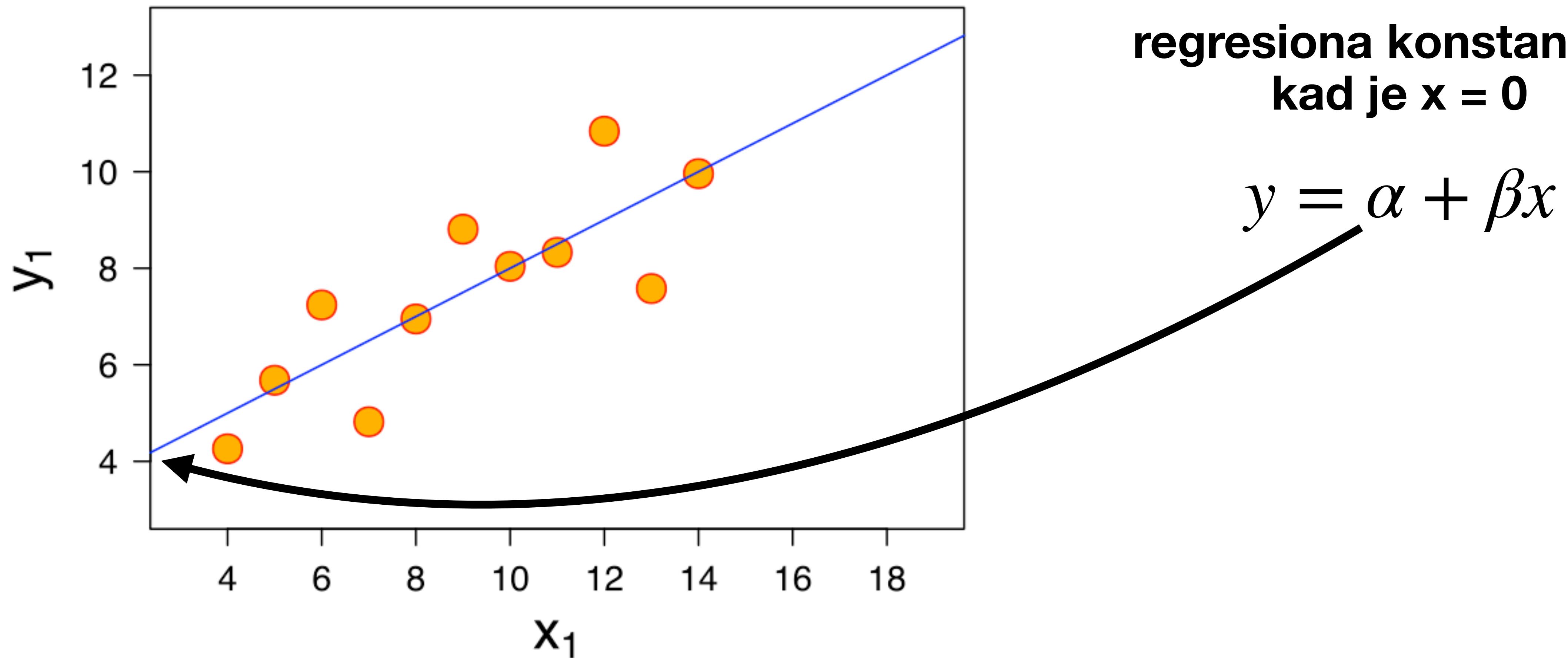
# Anatomija regresione jednačine



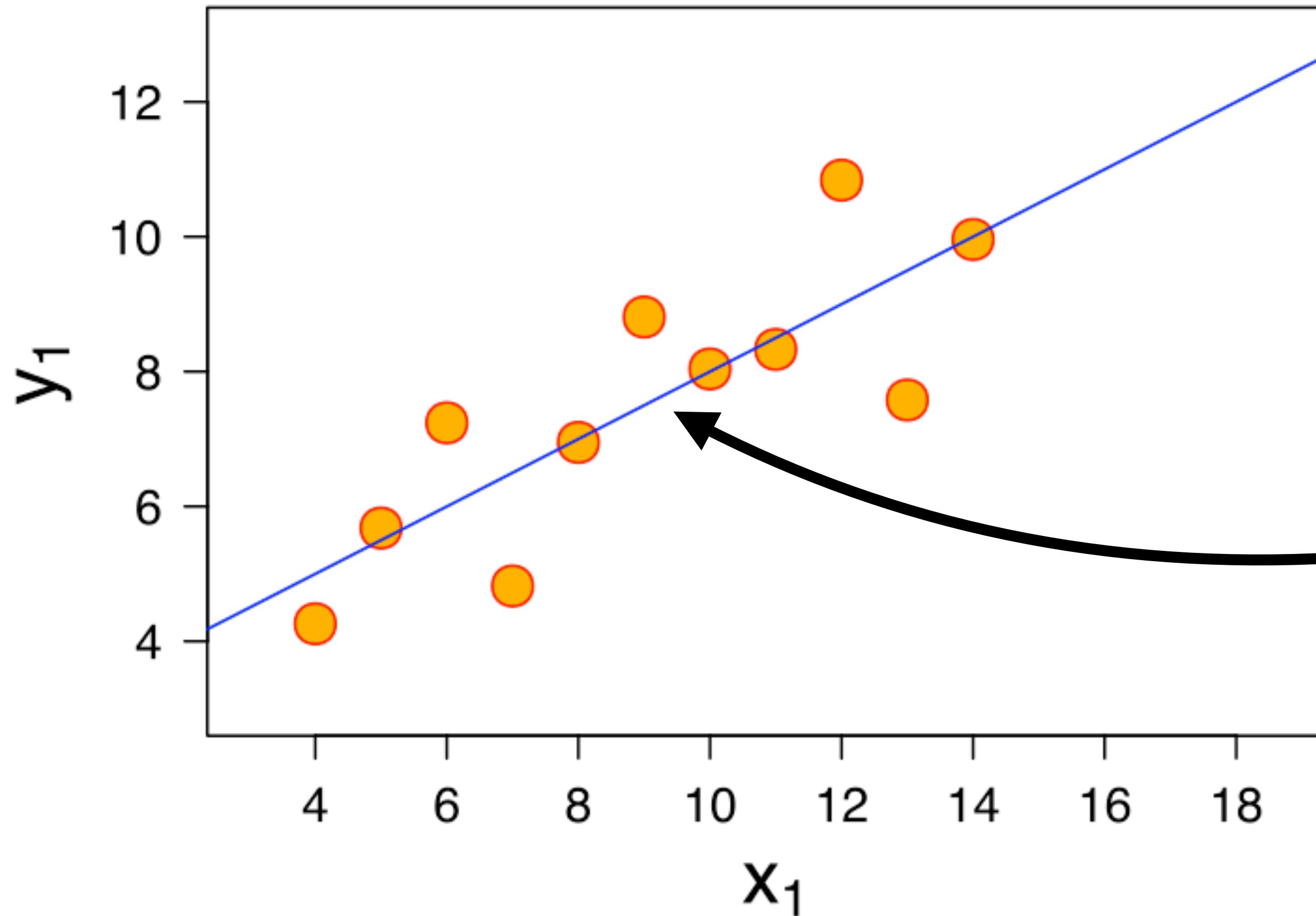
# Anatomija regresione jednačine



# Anatomija regresione jednačine



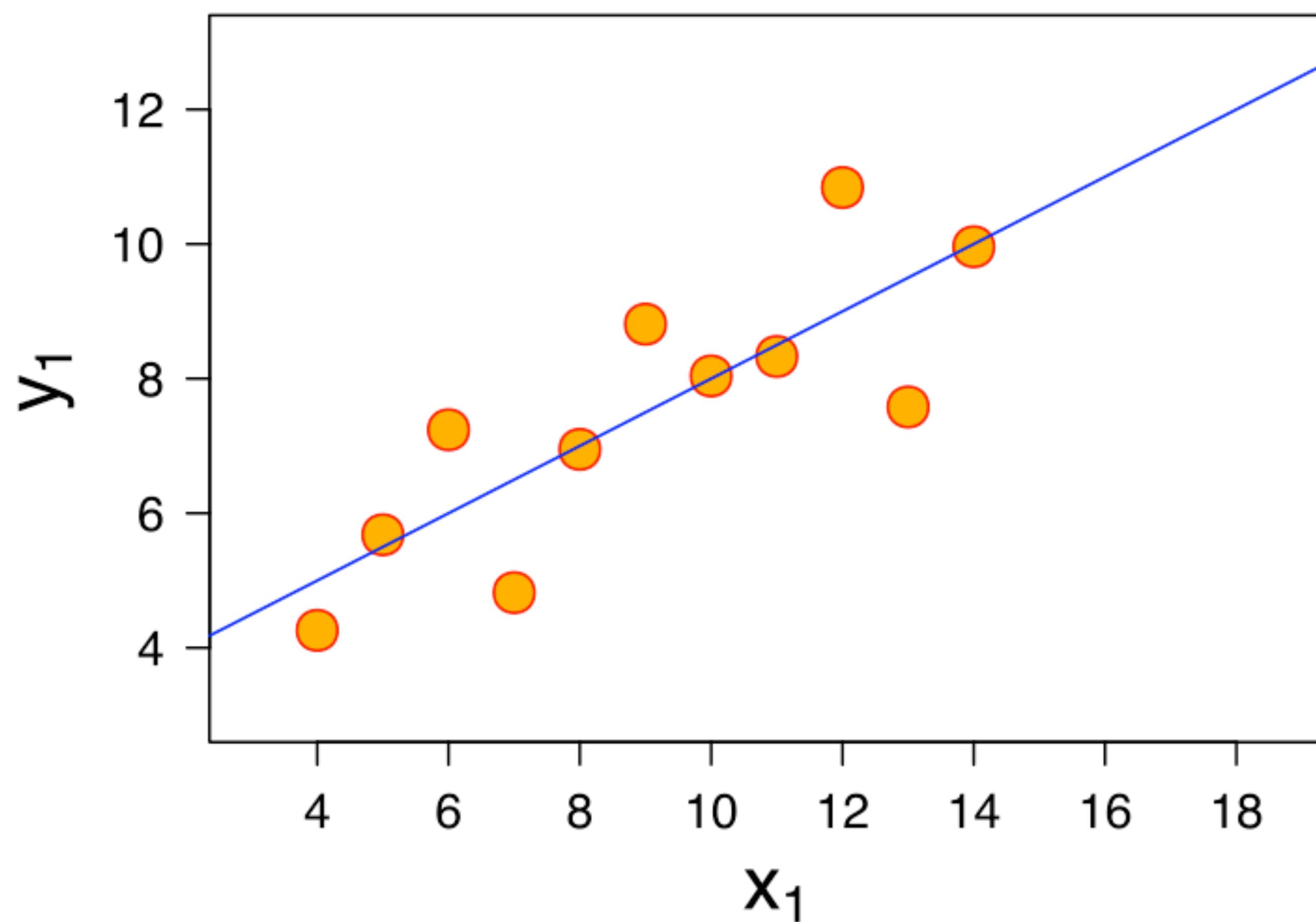
# Anatomija regresione jednačine



$$y = \alpha + \beta x$$

**regresioni koeficijent b:**  
nagib / vrednost za  
koju se Y promeni ako  
povećamo X za 1.

# Anatomija regresione jednačine



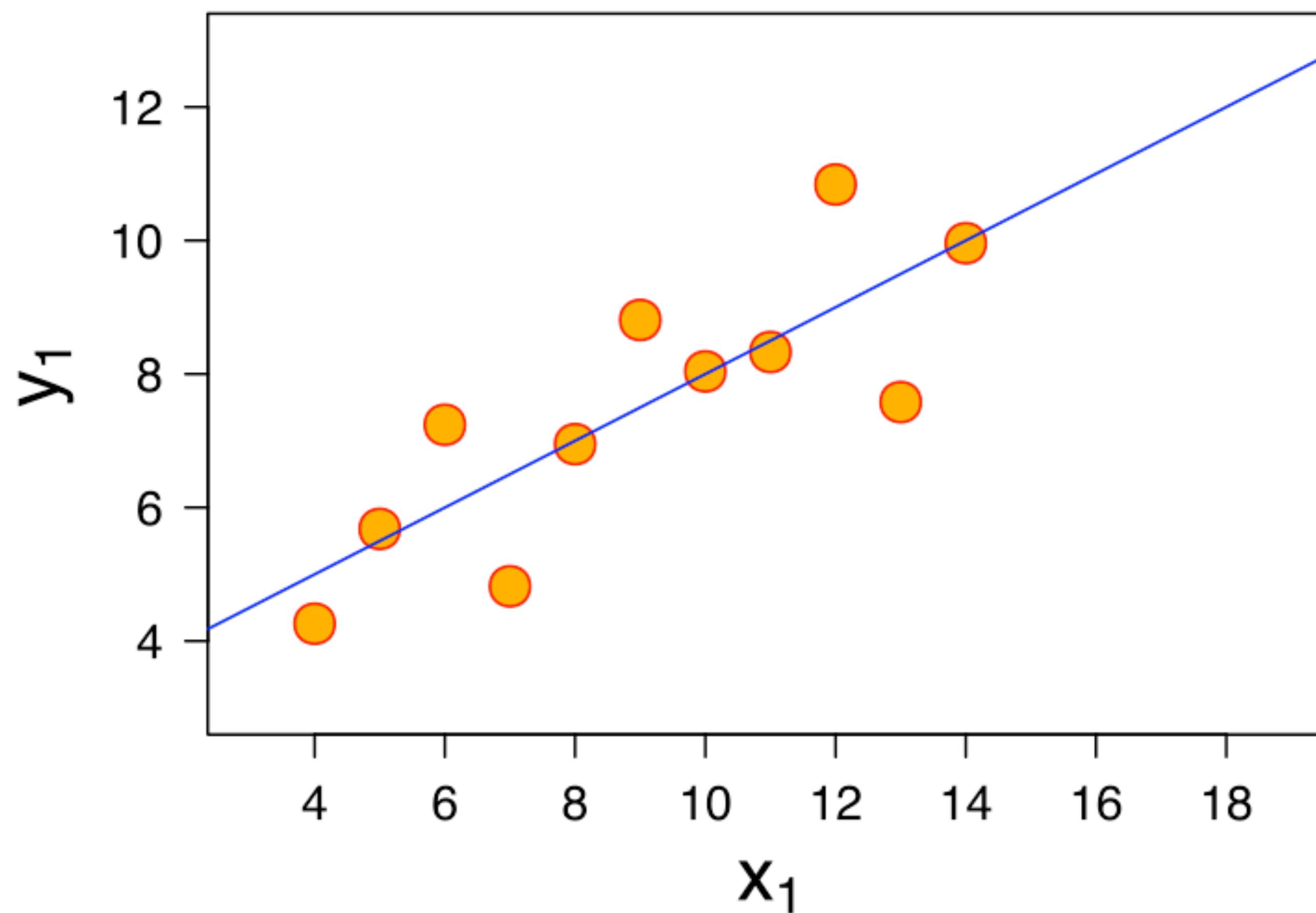
$$y = \alpha + \beta x$$

**Koeficijent determinacije**

$$r^2 = \text{cor}(x, y)^2, \text{ od } 0 \text{ do } 1$$

Udeo objasnjene varijacije

# Anatomija regresione jednačine



$$y = \alpha + \beta x$$

$$y = -4.16 + 0.13x$$

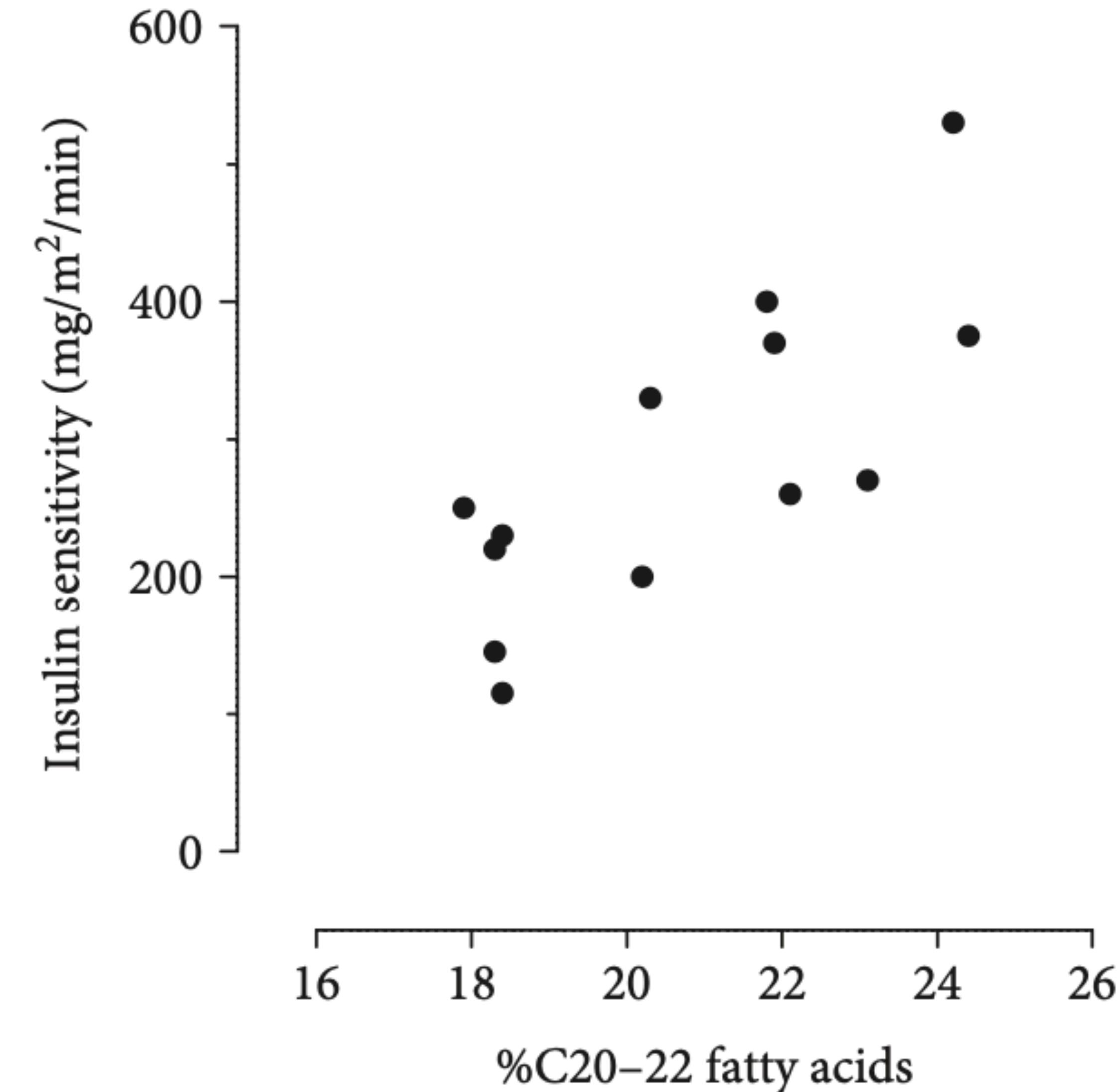
$$r^2 = 0.5 \times 100 = 50\%$$

# Primer: lipidi i insulinska rezistencija (2)

Borkman et al. (1993)

CORRELATION	
r	0.77
95% CI	0.3803 to 0.9275
$r^2$	0.5929
P (two-tailed)	0.0021
Number of XY pairs	13

Opisali smo povezanost



# Primer: inteligencija i koncentracije spermatozoida

Arden et al. (2008)

	r	$r^2$	P VALUE
Log (sperm count per ml)	0.15	0.023	0.0019
Log (sperm count per ejaculate)	0.19	0.036	< 0.0001
Fraction of sperm that are motile (%)	0.14	0.020	0.0038

# Primer: inteligencija i koncentracije spermatozoida

Arden et al. (2008)

	r	$r^2$	P VALUE
Log (sperm count per ml)	0.15	0.023	0.0019
Log (sperm count per ejaculate)	0.19	0.036	< 0.0001
Fraction of sperm that are motile (%)	0.14	0.020	0.0038

# Primer: inteligencija i koncentracije spermatozoida

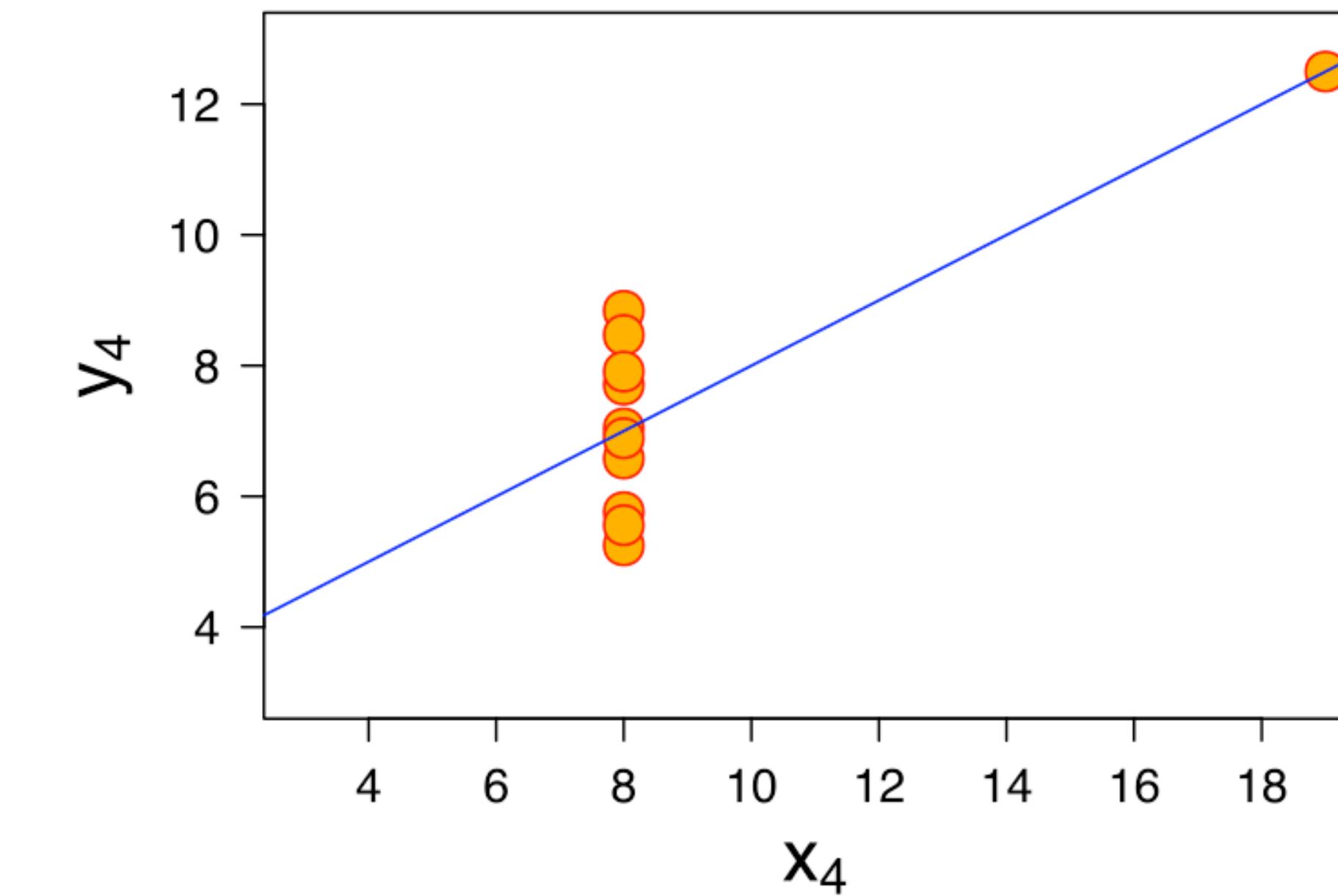
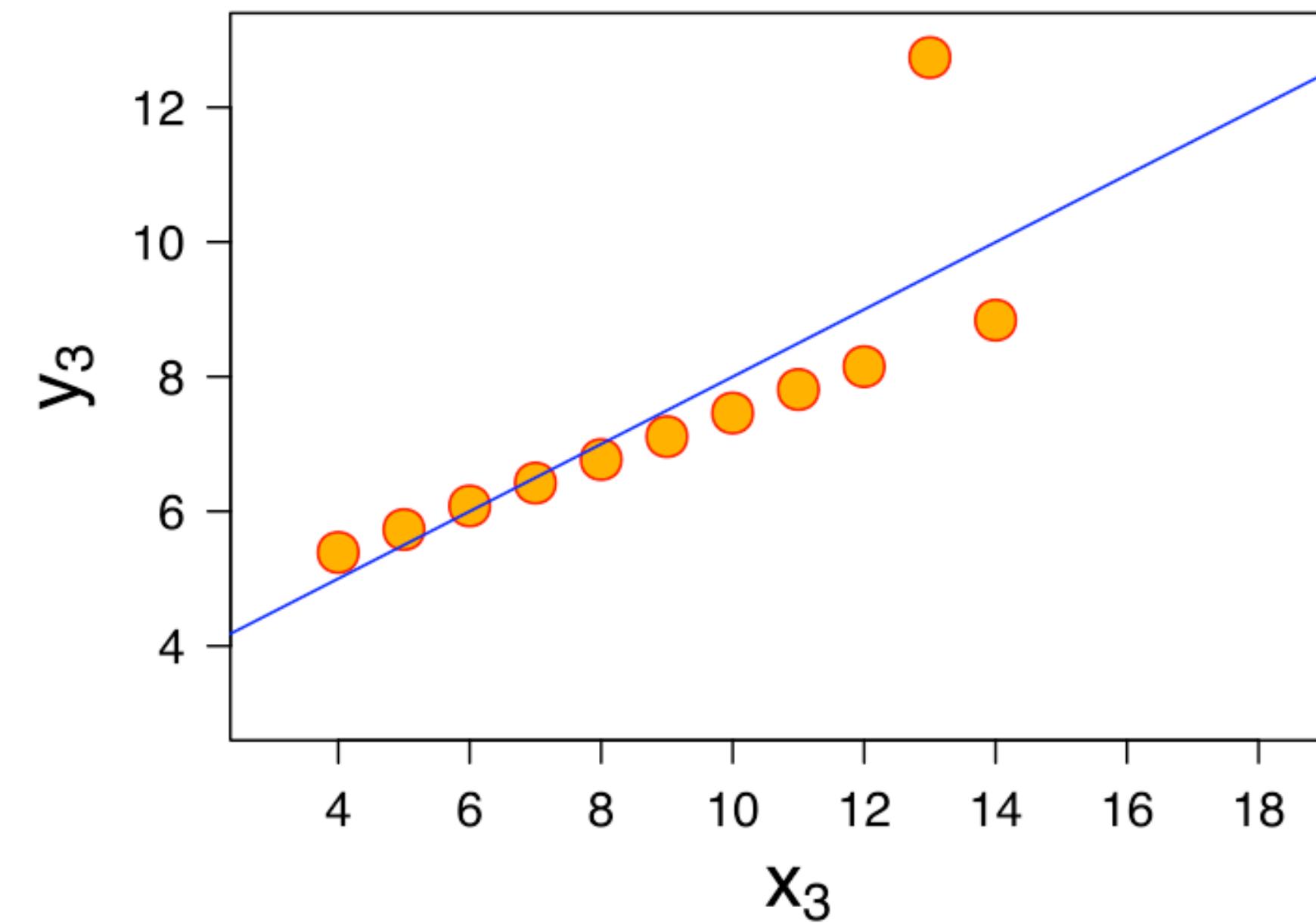
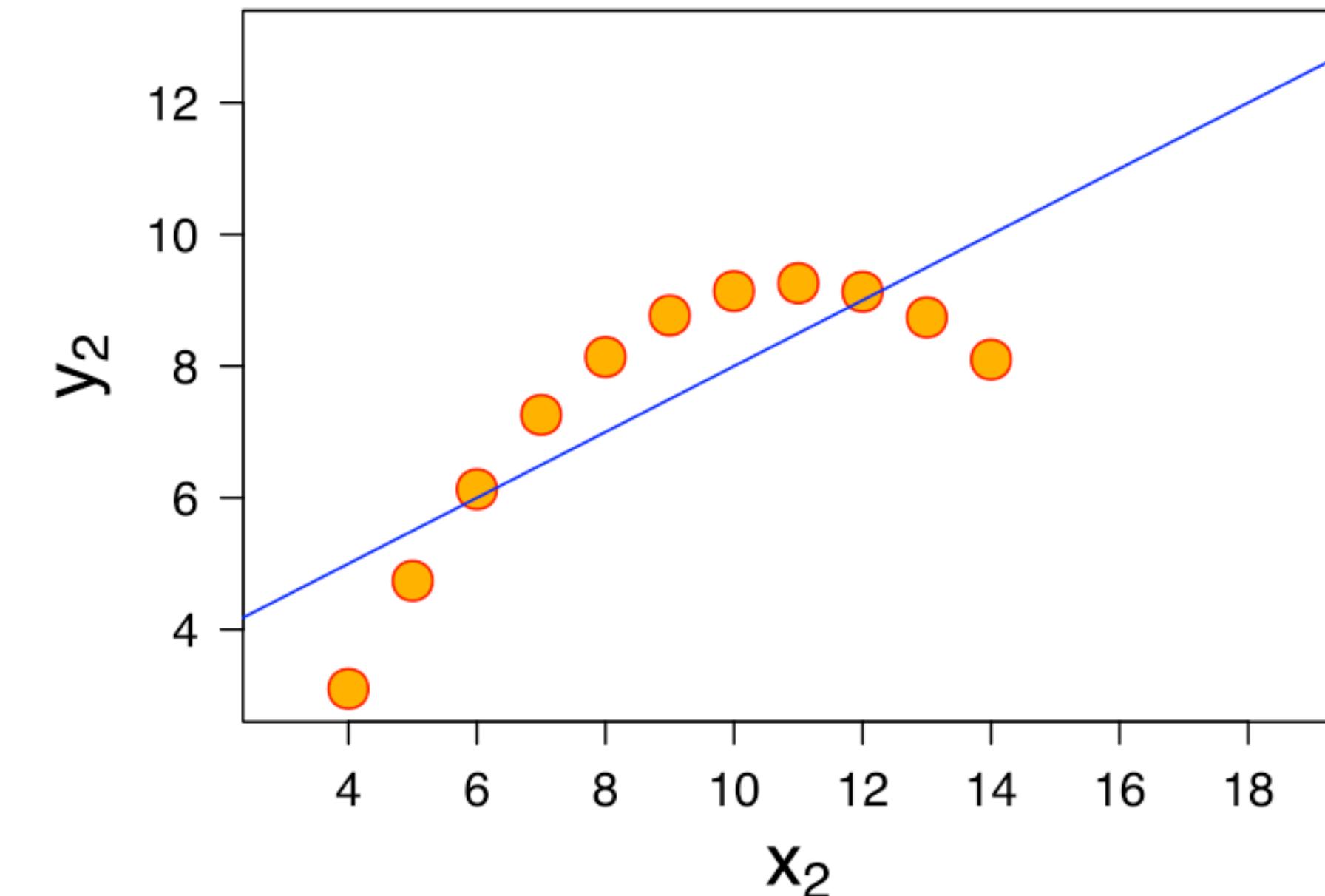
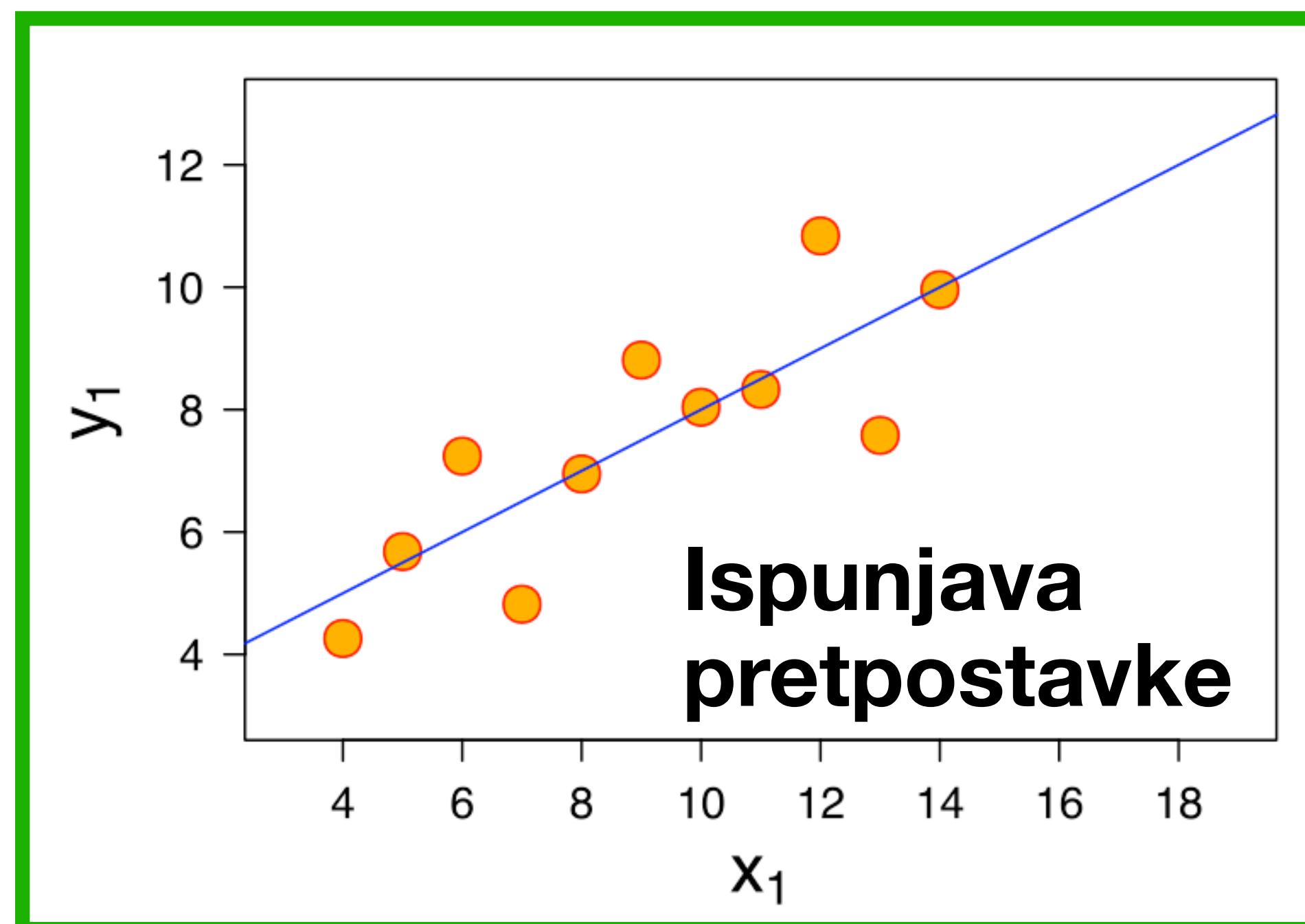
Arden et al. (2008)

% objašnjene varijanse

	r	r <sup>2</sup>	P VALUE
Log (sperm count per ml)	0.15	0.023	0.0019
Log (sperm count per ejaculate)	0.19	0.036	< 0.0001
Fraction of sperm that are motile (%)	0.14	0.020	0.0038

# Pretpostavke linearna regresije

1. Normalnost
2. Linearost
3. Nezavisnost (merenja nisu ponovljena)
4. Jednakost varijansi

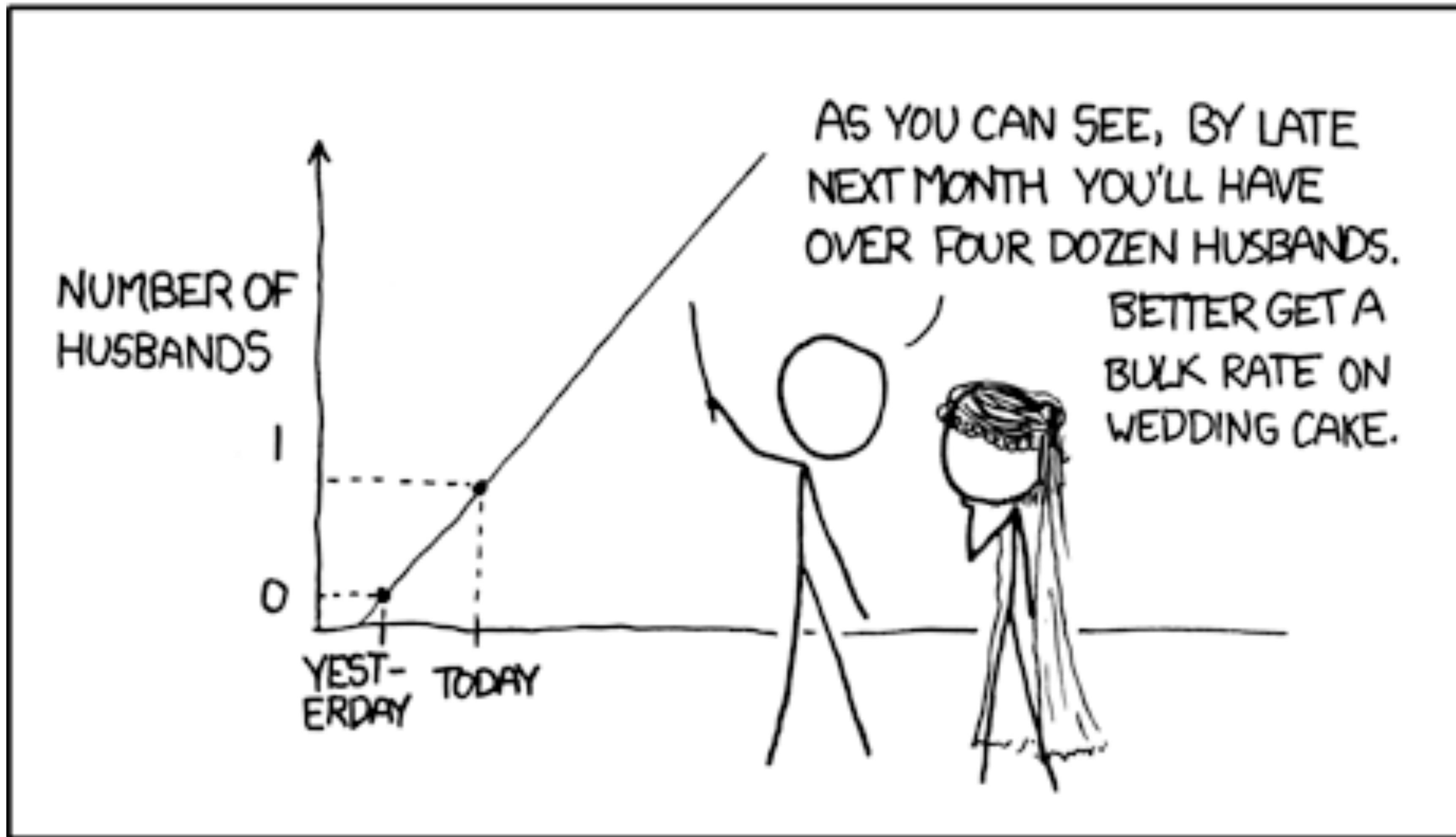


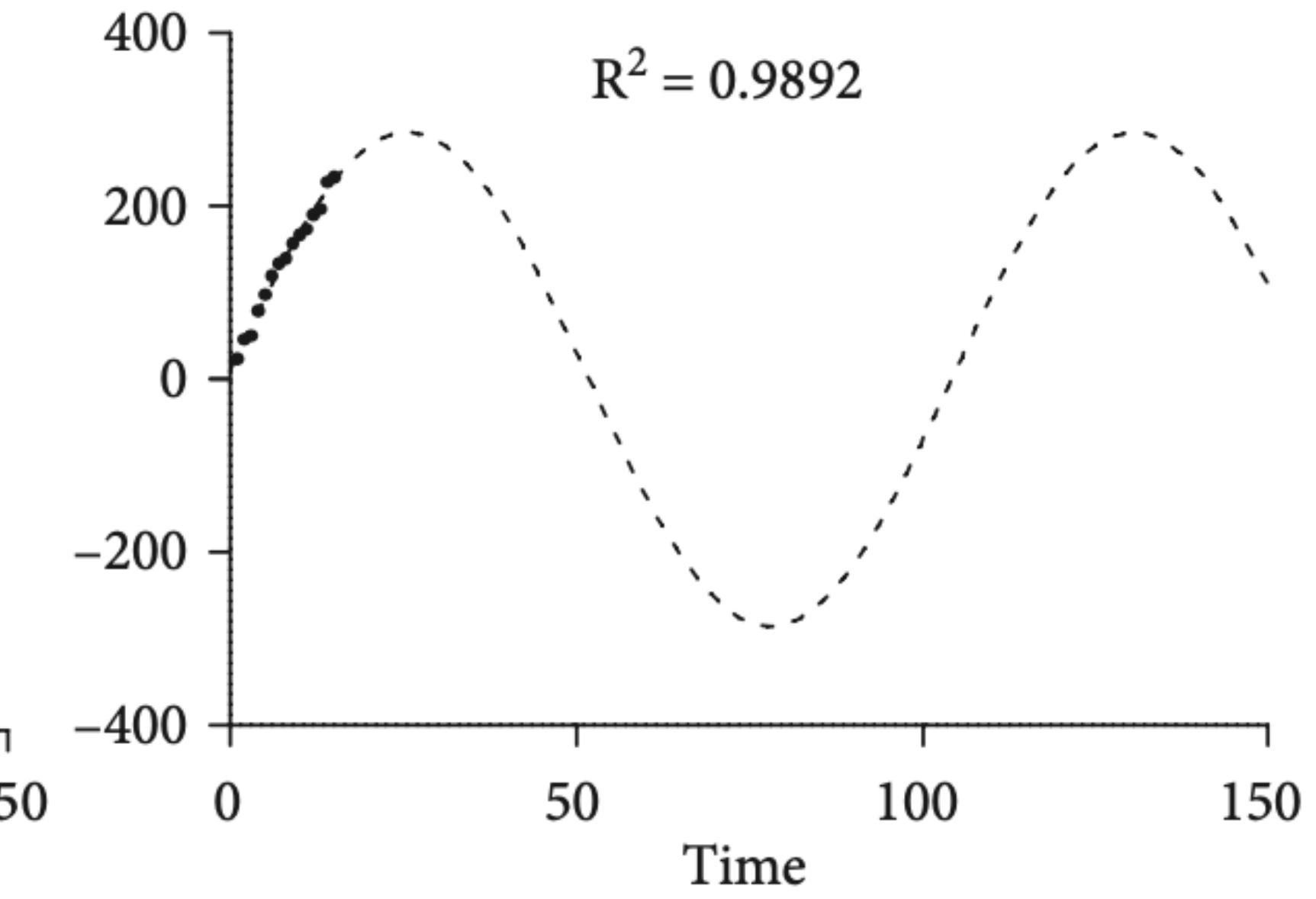
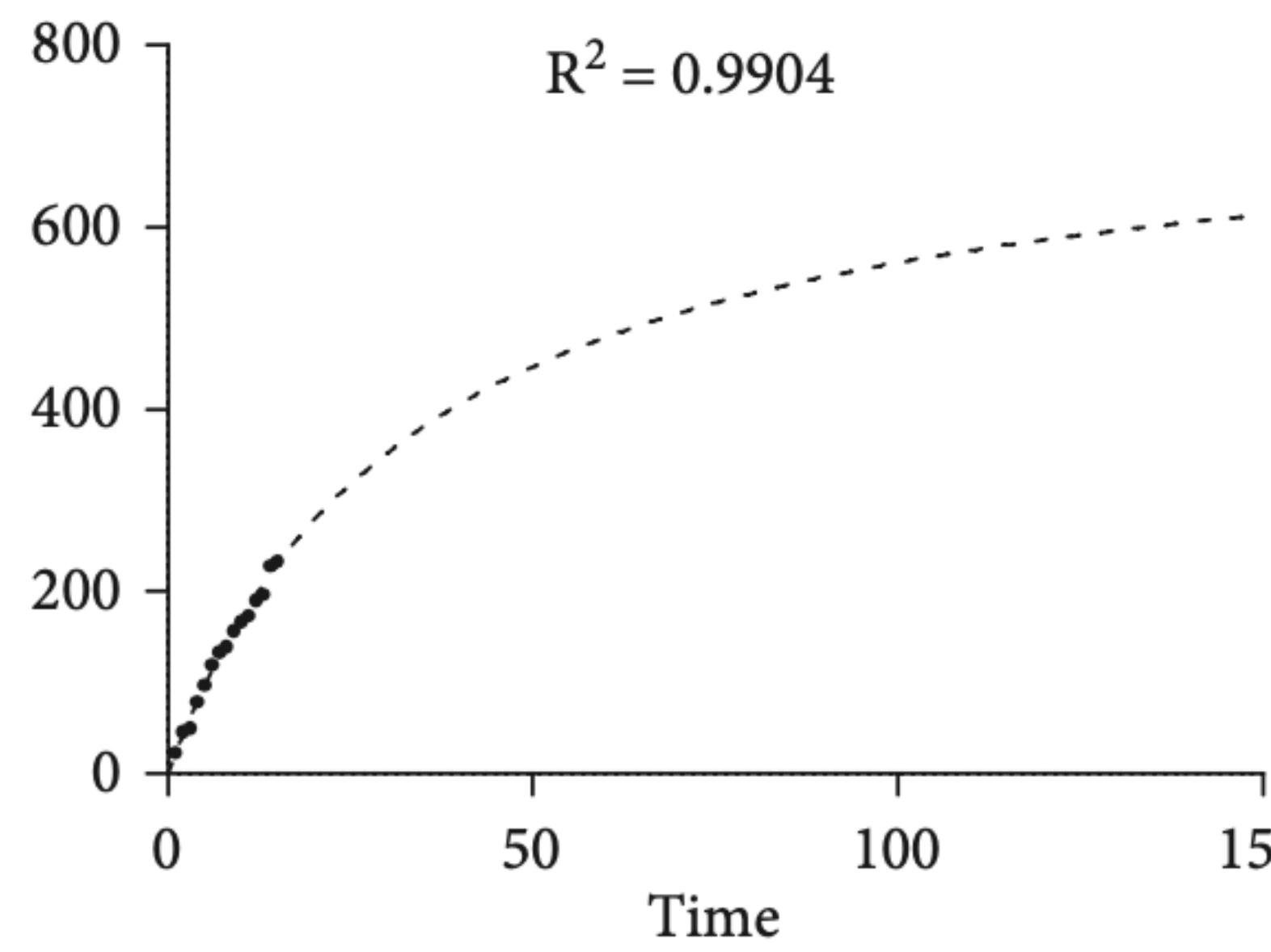
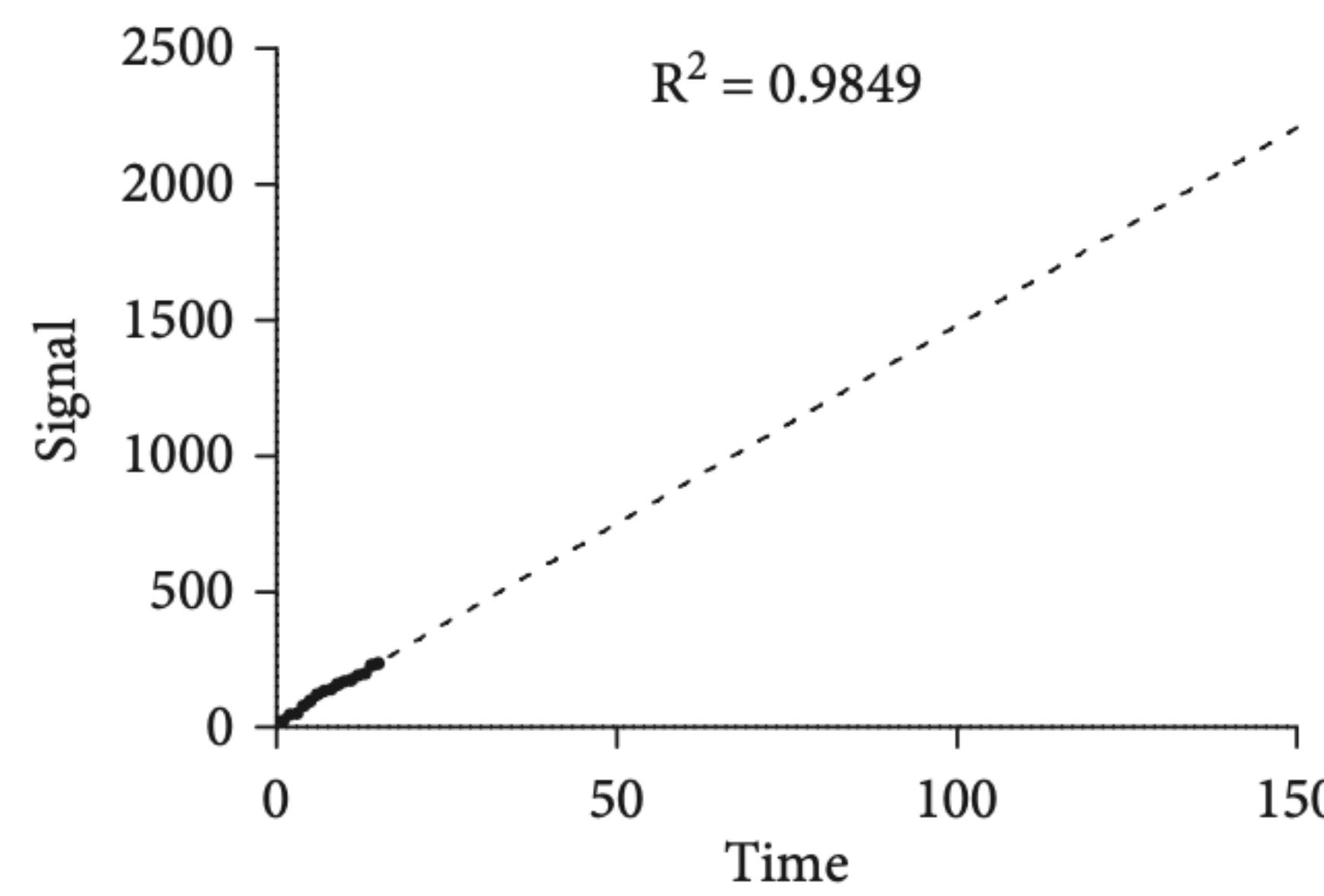
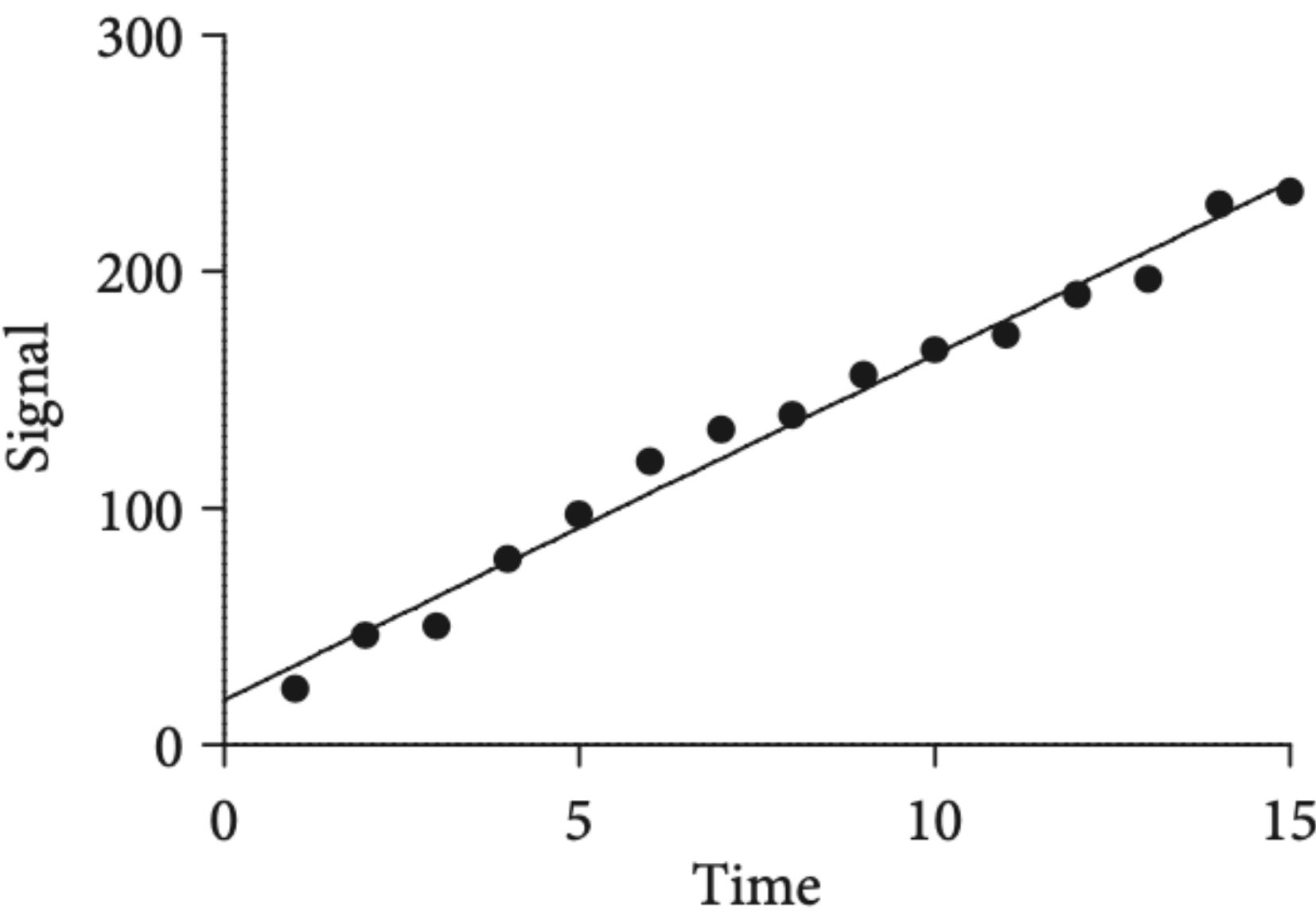
unutar

# Ekstrapolacija i interpolacija

Izvan

MY HOBBY: EXTRAPOLATING





Discrete variables

Continuous variables

Nonparametric tests

Survival analysis

Accuracy of diagnostic test

Matched-pair analysis

Metaanalysis and metaregression

Calculate sample size

, signif(res\$p.value, dig:

(Dataset\$Hemoglobin, Dataset  
a"))

Numerical summaries

Smirnov-Grubbs test for outliers

Kolmogorov-Smirnov test for normal distribution

Confidence interval for a mean

Single-sample t-test

Two-variances F-test

Two-sample t-test

Paired t-test

Bartlett's test

One-way ANOVA

Repeated-measures ANOVA

Multi-way ANOVA

ANCOVA

Test for Pearson's correlation

Linear regression

## Linear regression

Enter name for model: RegModel.1

Click pressing Ctrl key to select multiple variables.

Response variable (pick one)

Hemoglobin

rb

SaO2

 Wald test for overall p-value for factors with >2 levels Keep results as active model for further analyses Show basic diagnostic plots Stepwise selection based on AIC Stepwise selection based on BIC Stepwise selection based on p-value

Condition to limit samples for analysis. Ex1. age&gt;50 &amp; Sex==0 Ex2. age&lt;50 | Sex==1

&lt;all valid cases&gt;

Help

Reset

OK

Cancel

Apply

Coefficients:	regresiona konstanta $\alpha$				
	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	339.390	56.536	6.003	0.00184	**
SaO2	-1.984	0.673	-2.948	0.03195	*
---					
Signif. codes:	0 '***'	1	0.05	'.'	0.1 ' '

koeficijent nagiba  $\beta$

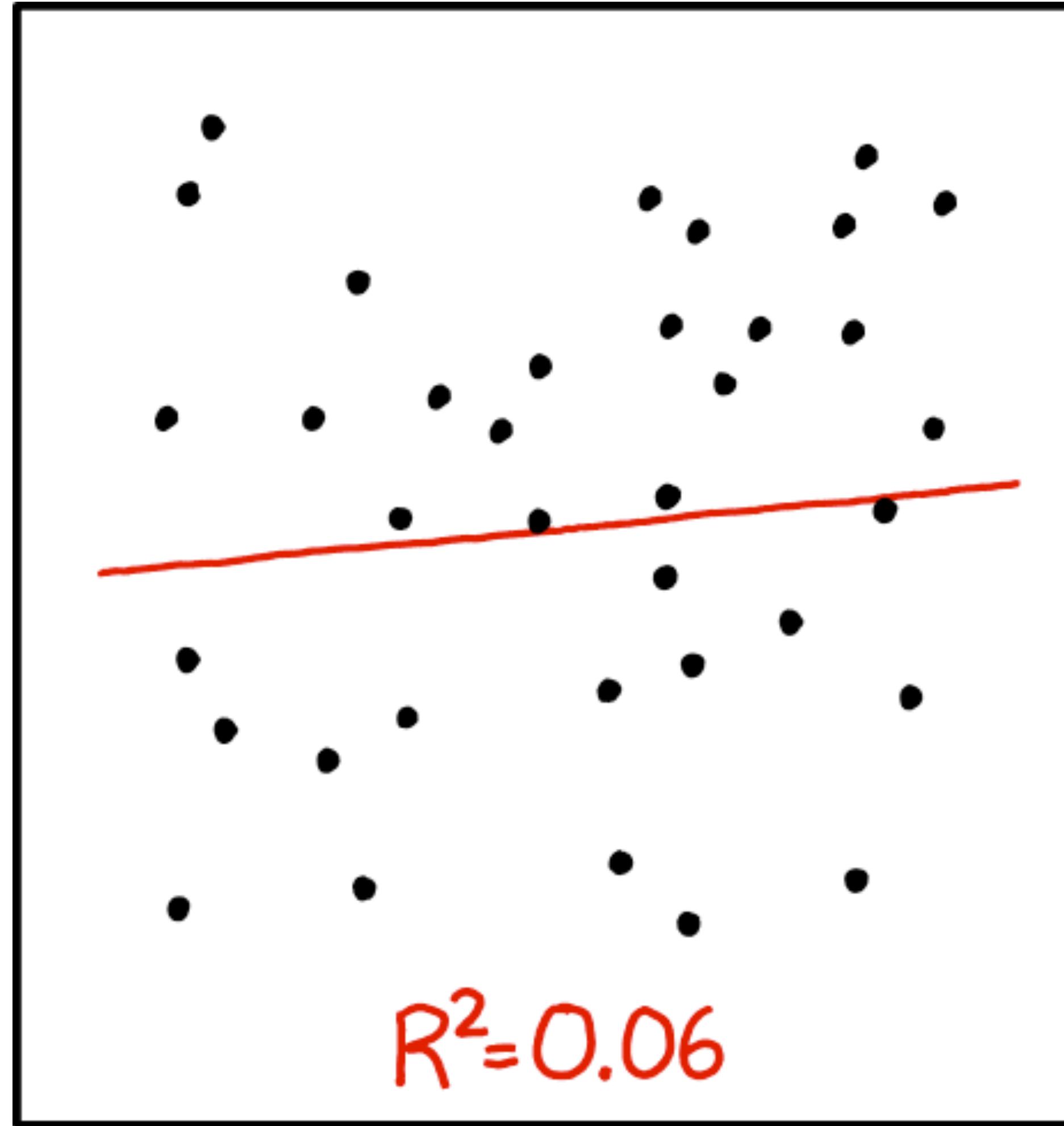
$y = 338.39 - 1.984x$

Koristeći dopunjenu **bazu DZ** odgovori na sledeća istraživačka pitanja:

1. Da li postoji povezanost između **sistolne tenzije i ukupnog holesterola?**
2. Za ispitanika sa ukupnim holesterolom od **6.2 mmol/l** izračunaj vrednost **sistolne tenzije (mmHg).**

# Korelacija vs. Regresija

- Korelacija opisuje povezanost varijabli, ne određuje najbolju liniju
- Kod korelacija nije potrebno razmišljati o uzrocima (simetrija)
- Kod regresije  $Y \leftarrow X \neq X \leftarrow Y$  (nema simetrije)



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER  
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE  
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.