



# Teorija pogrešaka

---

## Linearna Regresija



# Procjena pogreška

Pretpostavka: mjerimo neku veličinu relativno nepreciznim instrumentom, te uvijek dobivamo istu vrijednost  $\bar{x} = x_1 = x_2 = \dots = x_n$ .

Znači li to da nema pogreške?

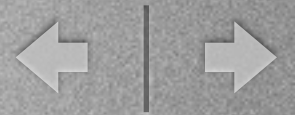
Odgovor: nikako!

To samo znači da instrument nije dovoljno osjetljiv.

U tom slučaju **procijenjujemo pogrešku** i pišemo:

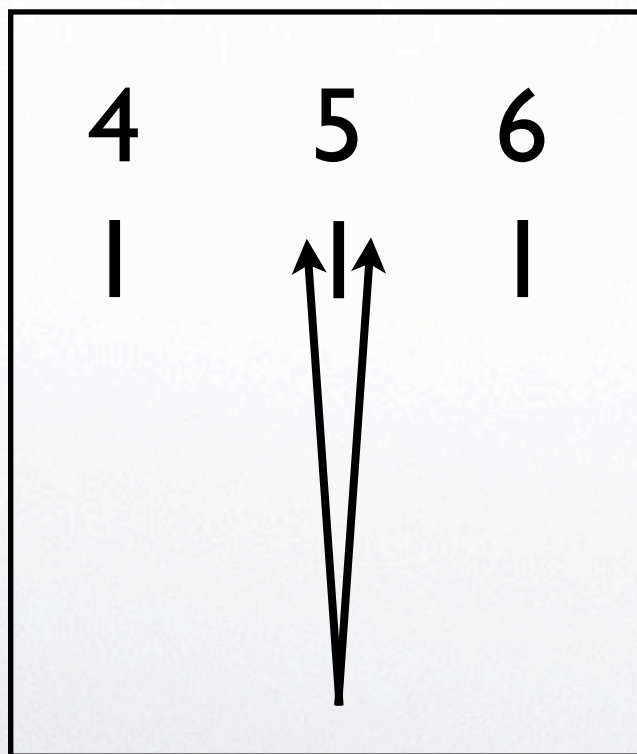
$$X = \bar{x} \pm \Delta x, \quad R = \frac{\Delta x}{\bar{x}} \cdot 100(\%)$$





# Procjena pogreška

Kako procijeniti  $\Delta x$ ?



$\Delta x$  je najmanji razmak u kojem na određenom instrumentu možemo uočiti promjenu

$\Delta x$  ima značenje *maksimalne pogreške*, koju možemo definirati i na skupu mjerenih rezultata  $x_1, x_2, \dots, x_n$

$$\Delta x = |\bar{x} - x_i|_{\max}$$



# Pogreška ovisnih veličina

Pretpostavimo da je fizikalna veličina  $F$  funkcija neposredno mjenjenih veličina  $y_1, y_2, \dots, y_n$

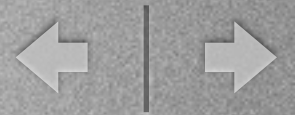
Primjeri:

1. mjerimo duljinu, širinu i visinu kvadra, a zanima nas njegov obujam ( $V = a \cdot b \cdot c$ )
2. mjerimo mase tri predmeta,  $m_1, m_2$  i  $m_3$ , a zanima nas ukupna masa ( $m = m_1 + m_2 + m_3$ )
3. mjerimo vrijeme  $t$  padanja kuglice s vrha nebodera, a želimo odrediti visinu nebodera ( $h = gt^2/2$ )

Najvjerojatnija vrijednost veličine  $\bar{F}$  je:

$$\bar{F} = F(\bar{y}_1, \bar{y}_2, \dots, \bar{y}_n)$$





# Pogreška ovisnih veličina

Primjer 1.  $\bar{a} = 2.31 \text{ cm}$ ,  $\bar{b} = 1.43 \text{ cm}$ ,  $\bar{c} = 4.11 \text{ cm}$

$$\bar{V} = \bar{a} \cdot \bar{b} \cdot \bar{c} = 13.58 \text{ cm}^3$$

No, što s pogreškom? Znamo pogreške osnovnih veličina, ali u kakvoj su one vezi s pogreškom ovisne veličine?

Iz zakona rasprostranjenja pogreške znamo da je uz  $F = y_1 + y_2 + \dots + y_k$ :

$$m^2(F) = m^2(y_1) + m^2(y_2) + \dots + m^2(y_k)$$

Što možemo pisati i kao:

$$m^2(F) = \left[ \frac{\partial F}{\partial y_1} m(y_1) \right]^2 + \left[ \frac{\partial F}{\partial y_2} m(y_2) \right]^2 + \dots + \left[ \frac{\partial F}{\partial y_k} m(y_k) \right]^2 = \sum_{j=1}^k \left[ \frac{\partial F}{\partial y_j} m(y_j) \right]^2$$



# Pogreška ovisnih veličina

Pokazuje se da se ovaj pristup može poopćiti i na nelinearne funkcije  $F(y_1, y_2, \dots, y_k)$ , te da vrijedi:

$$m^2(F) = \sum_{j=1}^k \left[ \frac{\partial F}{\partial y_j} \right]^2 m^2(y_j) \quad \text{srednja pogreška}$$

$$\sigma^2(F) = \sum_{j=1}^k \left[ \frac{\partial F}{\partial y_j} \right]^2 \sigma^2(y_j) \quad \text{standardna devijacija}$$

$$M^2(F) = \sum_{j=1}^k \left[ \frac{\partial F}{\partial y_j} \right]^2 M^2(y_j) \quad \text{nepouzdanost}$$





# Pogreška ovisnih veličina

Primjer:  $F = a \cdot b + c = F(a, b, c)$

$$M^2(F) = \left\{ \left( \frac{\partial F}{\partial a} M_a \right)^2 + \left( \frac{\partial F}{\partial b} M_b \right)^2 + \left( \frac{\partial F}{\partial c} M_c \right)^2 \right\}_{\substack{a=\bar{a} \\ b=\bar{b} \\ c=\bar{c}}} =$$
$$= (\bar{b} \cdot M_a)^2 + (\bar{a} \cdot M_b)^2 + (1 \cdot M_c)^2$$





# Pogreška ovisnih veličina

Primjer iz praktikuma (vježba 8):

Modul elastičnosti čelika određuje se mjerenjem dimenzija čelične šipke, savijenosti čelične šipke po jedinici sile i udaljenosti potporanja. Teorijska formula za savijenost šipke je:

$$\lambda = \frac{1}{4E} \frac{L^3}{ab^3} F,$$

gdje je  $E$  modul elastičnosti čelika,  $L$  je udaljenost potporanja,  $a$  je širina šipke,  $b$  je debljina šipke, a  $F$  je sila teže utega.

Izmjerene su sljedeće veličine:

$$A = \frac{\lambda}{F} = (0,76 \pm 0,01) \frac{\text{mm}}{\text{N}} = (7,6 \pm 0,1) \cdot 10^{-4} \frac{\text{m}}{\text{N}} = \bar{A} \pm M_A$$

$$a = (10,26 \pm 0,05) \text{ mm} = (1,026 \pm 0,005) \cdot 10^{-2} \text{ m} = \bar{a} \pm M_a$$

$$b = (1,53 \pm 0,03) \text{ mm} = (1,53 \pm 0,03) \cdot 10^{-3} \text{ m} = \bar{b} \pm M_b$$

$$L = (29,0 \pm 0,1) \text{ cm} = (2,90 \pm 0,01) \cdot 10^{-1} \text{ m} = \bar{L} \pm M_L$$

Najvjerojatnija vrijednost za modul elastičnosti čelika je:

$$\bar{E} = \frac{1}{4} \frac{\bar{L}^3}{\bar{a}\bar{b}^3} \frac{1}{\bar{A}} = 2,183226 \text{ N/m}^2$$

Varijanca modula elastičnosti je:

$$M_E^2 = \left( \frac{\partial E}{\partial A} \right)^2 M_A^2 + \left( \frac{\partial E}{\partial a} \right)^2 M_a^2 + \left( \frac{\partial E}{\partial b} \right)^2 M_b^2 + \left( \frac{\partial E}{\partial L} \right)^2 M_L^2 =$$





# Pogreška ovisnih veličina

$$\begin{aligned} &= \left( \frac{\bar{L}^3}{4\bar{a}\bar{b}^3\bar{A}^2} \right)^2 M_A^2 + \left( \frac{\bar{L}^3}{4\bar{a}^2\bar{b}^3\bar{A}} \right)^2 M_a^2 + \left( \frac{3\bar{L}^3}{4\bar{a}\bar{b}^4\bar{A}} \right)^2 M_b^2 + \left( \frac{3\bar{L}^2}{4\bar{a}\bar{b}^3\bar{A}} \right)^2 M_L^2 = \\ &= \bar{E}^2 \left[ \left( \frac{M_A}{\bar{A}} \right)^2 + \left( \frac{M_a}{\bar{a}} \right)^2 + \left( 3 \frac{M_b}{\bar{b}} \right)^2 + \left( 3 \frac{M_L}{\bar{L}} \right)^2 \right] = \\ &= \bar{E}^2 \left[ 1,73 \cdot 10^{-4} + 2,37 \cdot 10^{-5} + 3,46 \cdot 10^{-3} + 1,07 \cdot 10^{-4} \right] \end{aligned}$$

Standardna pogreška je:

$$M_E = \bar{E} \sqrt{3,76 \cdot 10^{-3}} = 0,13 \cdot 10^{11} \text{ N/m}^2$$

Rezultat za modul elastičnosti pišemo:

$$E = (2,2 \pm 0,1) \cdot 10^{11} \text{ N/m}^2$$





# Pogreška ovisnih veličina

Što ako imamo samo procijenjene pogreške  $\Delta y_i$ ?

Tada ne možemo govoriti o rasprostranjenju pogreške, što je zakon definiran na statističkom skupu.

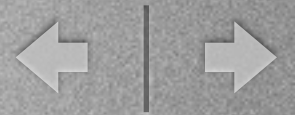
Umjesto toga razvijamo  $F$  u red oko  $\bar{F}$  i računamo *maksimalnu pogrešku*.

$$F = F(\bar{y}_1, \bar{y}_2, \dots, \bar{y}_k) \pm \sum_{j=1}^k \left( \frac{\partial F}{\partial y_j} \right)_{y_j=\bar{y}_j} \times \Delta y_j = \bar{F} \pm \Delta F$$

Budući da je  $\Delta F \geq 0$ , uzimamo apsolutne vrijednosti derivacija, pa je:

$$\Delta F = \sum_{j=1}^k \left| \frac{\partial F}{\partial y_j} \right|_{y_j=\bar{y}_j} \times \Delta y_j$$





# Pogreška ovisnih veličina

Primjer:

$$F = a \cdot b + c$$

$$\Delta F = \left\{ \left| \frac{\partial F}{\partial a} \right| \Delta a + \left| \frac{\partial F}{\partial b} \right| \Delta b + \left| \frac{\partial F}{\partial c} \right| \Delta c \right\}_{\substack{a=\bar{a} \\ b=\bar{b} \\ c=\bar{c}}} = \bar{b} \Delta a + \bar{a} \Delta b + \Delta c$$



## Mjerenja s različitim statističkim težinama

Pretpostavimo da je u više različitih laboratorija neovisno mjerena ista veličina  $X$ . Ostvareni su rezultati:

$$X_1 = \bar{x}_1 \pm M_1, X_2 = \bar{x}_2 \pm M_2, \dots, X_l = \bar{x}_l \pm M_l$$

Kako će netko tko dobije samo listu rezultata, a ne zna ništa o pojedinostima eksperimenta, izraziti rezultat ovih mjerenja u obliku:

$$X = \bar{\bar{x}} + \bar{M}$$

gdje je  $\bar{\bar{x}}$  *opća aritmetička sredina*, a  $\bar{M}$  *opća nepouzdanost*.

Mjerenja se obično izvode na različite načine i s različitim pouzdanostima pa ih ne smijemo tretirati ravnopravno. Kažemo da mjerenja imaju različite statističke težine.





## Mjerenja s različitim statističkim težinama

Pretpostavimo da svakom rezultatu pridijelimo *statističku težinu*  $w_j$ , tako da je:

$$\bar{\bar{x}} = \sum_{j=1}^l w_j \bar{x}_j; \quad \sum_{j=1}^l w_j$$

Dakle,  $w_j$  ima slično značenje kao i relativna frekvencija pojavljivanja događaja, a  $\bar{x}_j$  kao pojedinačno mjerenje.

U praksi se najčešće uzima (veća težina za manji  $M_j$ ):

$$w_j = \frac{1}{M_j} \bigg/ \sum_{k=1}^l \frac{1}{M_k^2} \quad \longrightarrow \quad \text{zadovoljava} \quad \sum_{j=1}^l w_j = 1$$



## Mjerenja s različitim statističkim težinama

Prema tome, opća aritmetička sredina je:

$$\bar{\bar{x}} = \frac{\sum_{j=1}^l \frac{\bar{x}_j}{M_j^2}}{\sum_{j=1}^l \frac{1}{M_j^2}}$$





## Mjerenja s različitim statističkim težinama

Koliki je  $M$ ?

Zamislimo jedno mjerenje koje ima statističku težinu  $w = 1$  i nepouzdanost  $\bar{M}$ , a u potpunosti zamjenjuje  $l$  rezultata.

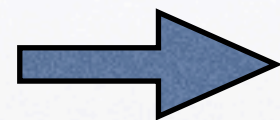
Znamo otprije:

$$\frac{w_i}{w_j} = \frac{M_j^2}{M_i^2}$$

i stavimo

$$w_i = 1$$
$$M_i = \bar{M}$$

$$\frac{1}{w_j} = \frac{M_j^2}{\bar{M}^2}$$



$$\bar{M}^2 = w_j M_j^2 = \frac{1}{\sum_{k=1}^l \frac{1}{M_k^2}}$$



# Mjerenja s različitim statističkim težinama

Prema tome, opća nepouzdanost je:

$$\bar{M} = \sqrt{\frac{1}{\sum_{j=1}^l \frac{1}{M_j^2}}}$$





## Mjerenja s različitim statističkim težinama

### Kako računati?

(i)  $\bar{x}_i \approx \bar{x}_j$  i  $M_i \approx M_j$ ,  $\forall i, j$

koristiti gornje formule za  $\bar{\bar{x}}$  i  $\bar{M}$

(ii)  $M_k \ll M_i$ , za neki  $k$ , tj. rezultat  $X_k = \bar{x}_k \pm M_k$  puno je precizniji od svih ostalih.

$w_k \gg w_i$  za sve  $i$ , pa uzimamo  $w_k = 1$ :

$$\bar{\bar{x}} = \bar{x}_k; \quad \bar{M} = M_k$$

(iii)  $M_i \approx M_j$  za sve  $i, j$ , ali  $\bar{x}_i$  i  $\bar{x}_j$  se znatno razlikuju za sve  $i, j$

zanemarujemo  $M_i$ , te smatramo  $\bar{x}_i$  neposredno mjerenim veličinama:

$$\bar{\bar{x}} = \frac{1}{l} \sum_{j=1}^l \bar{x}_j; \quad \bar{M} = \sqrt{\frac{\sum_{j=1}^l (\bar{\bar{x}} - \bar{x}_j)^2}{l(l-1)}}$$



# Korelacija u linearnoj regresiji

Zamislimo neki događaj na kojem su definirana dva obilježja  $x$  i  $y$ .

Svako promatranje tog događaja daje par brojeva  $(x, y)$ .

Npr. elastična sila, obilježja: masa utega i duljina opruge

slobodni pad, obilježja: vrijeme i prijeđeni put

studenti, obilježja: uspjeh na prijemnom i uspjeh na studiju

Pitamo se: jesu li obilježja  $x$  i  $y$  na bilo koji način povezana, tj. *korelirana*?

Drugim riječima, postoji li funkcionalna veza  $y(x)$  koja može dobro opisati parove  $(x, y)$ ?

Mi ćemo se zadržati na najjednostavnijem slučaju, tj. istraživanjem postoji li zadovoljavajući opis u okvirima *linearne regresije*.

$$y = ax + b$$

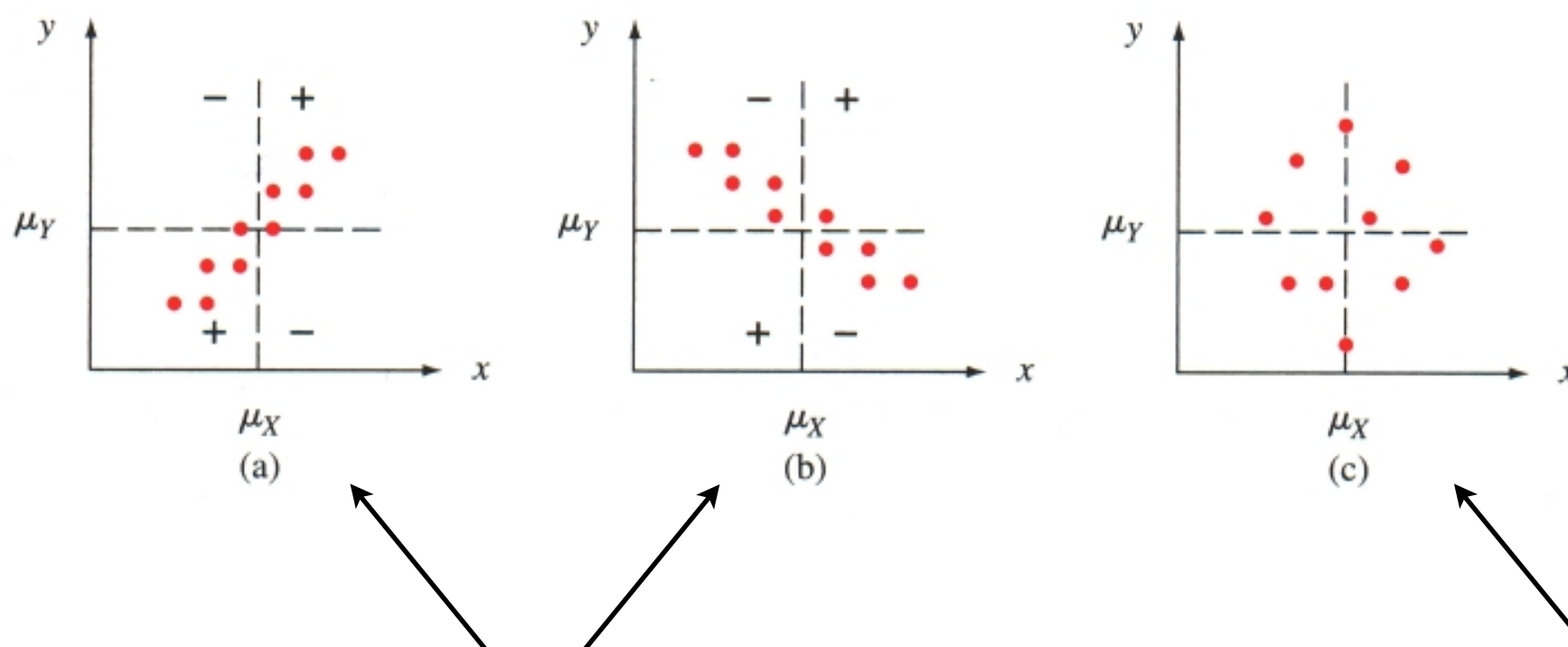




# Korelacija u linearnoj regresiji

Prva ocjena može se donijeti crtanjem grafa, tj. točaka u  $x$ - $y$  ravnini:

Primjer:



- možemo “od oka” povući pravac;  
korelacija vjerojatno postoji:  
1. koja je jednačba pravca?  
2. koja je kvantitativna ocjena  
postojanja korelacije

- ne možemo “od oka” povući pravac;  
korelacija vjerojatno ne postoji:  
1. koja je kvantitativna ocjena  
nepostojanja korelacije



# Korelacija u linearnoj regresiji

Događaj možemo tretirati kao 2D slučajnu varijablu, pri čemu frekvencija pojavljivanja nekog događaja ima značenje (nenormirane) vjerojatnosti:

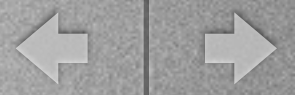
$$P(x_i, y_j) \Leftrightarrow \frac{f_{i,j}}{N}, \quad N = \sum_{i=1}^n \sum_{j=1}^m f_{i,j}$$

$$f_i = \sum_{j=1}^m f_{ij}; \quad f_j = \sum_{i=1}^n f_{ij}$$

Npr.  $\{(1,1), (1,1), (1,2), (2,2), (2,2), (2,2)\}$

$$f_{11} = 2, f_{12} = 1, f_{22} = 2; \quad N = 2 + 1 + 2 = 5$$





# Korelacija u linearnoj regresiji

Također vrijedi:

$$\bar{x} = \frac{1}{N} \sum_{i,j} f_{i,j} x_i = \frac{1}{N} \sum_{i=1}^n f_i x_i$$

$$\bar{y} = \frac{1}{N} \sum_{i,j} f_{i,j} y_j = \frac{1}{N} \sum_{j=1}^m f_j y_j$$

$$\sigma_x^2 = \frac{1}{N} \sum_{i,j} f_{i,j} (x_i - \bar{x})^2 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x})^2$$

$$\sigma_y^2 = \frac{1}{N} \sum_{i,j} f_{i,j} (y_j - \bar{y})^2 = \frac{1}{N} \sum_{j=1}^m f_j (y_j - \bar{y})^2$$

$$\sigma_{xy} = \frac{1}{N} \sum_{i,j} f_{i,j} (x_i - \bar{x})(y_j - \bar{y}) = \frac{1}{N} \sum_{i,j} f_{i,j} x_i y_j - \bar{x} \cdot \bar{y}$$



# Korelacija u linearnoj regresiji

Sjetimo se da je za nezavisna svojstva  $x$  i  $y$ ,  $\sigma_{xy} = 0$

Definiramo *koeficijent korelacije*  $r$ :

$$r = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

$\sigma_{xy} \rightarrow 0$ ;  $r$  je mali, korelacija je slaba

Najčešći kriterij dobrote korelacije:

- korelacija je dobra ukoliko je  $|r| \geq 0.5$

U praksi je kriterij često stroži, npr  $|r| \geq 0.7$

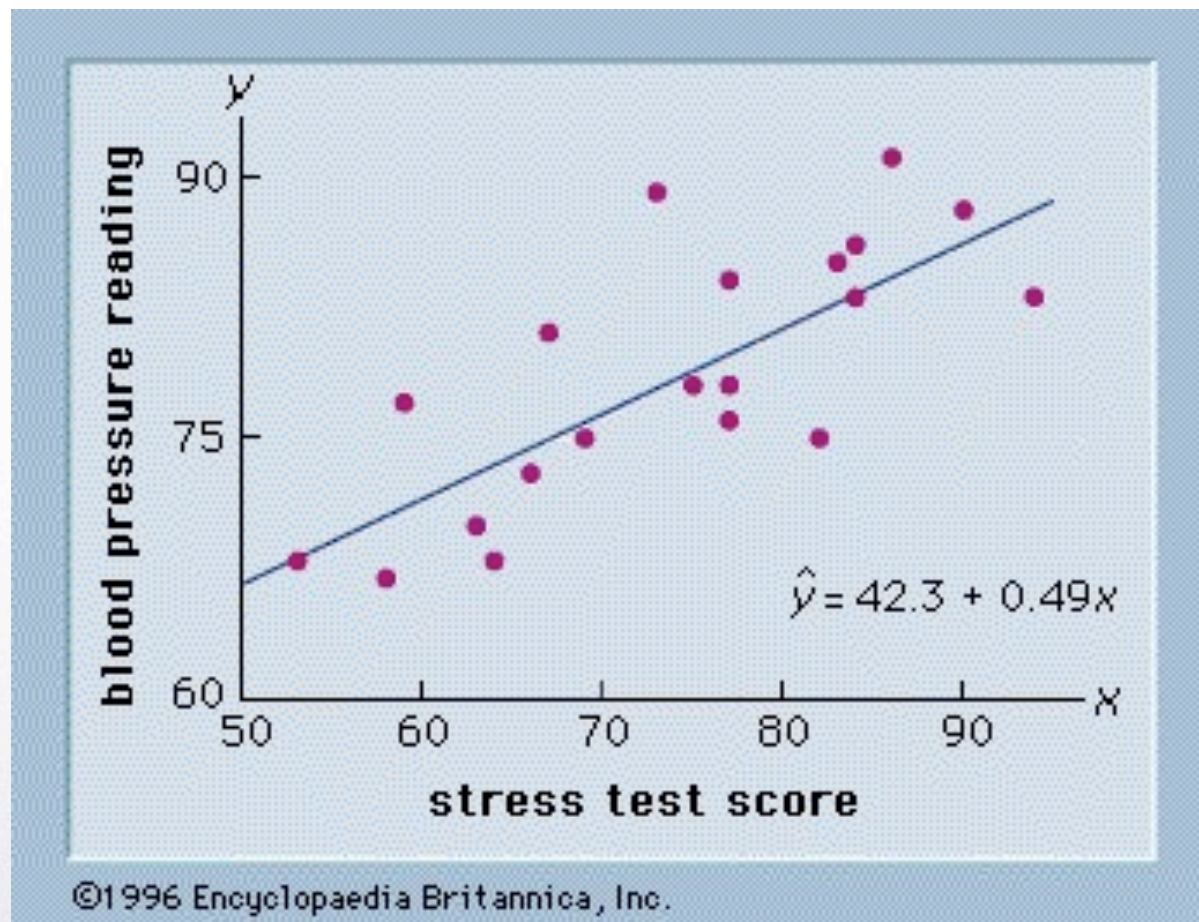
Dakle,  $r$  je *kvantitativni* pokazatelj prisustva ili odsustva korelacije veličina  $x$  i  $y$





# Korelacija u linearnoj regresiji

No, kako povući najbolji pravac  $y(x)$  ili  $x(y)$ ?



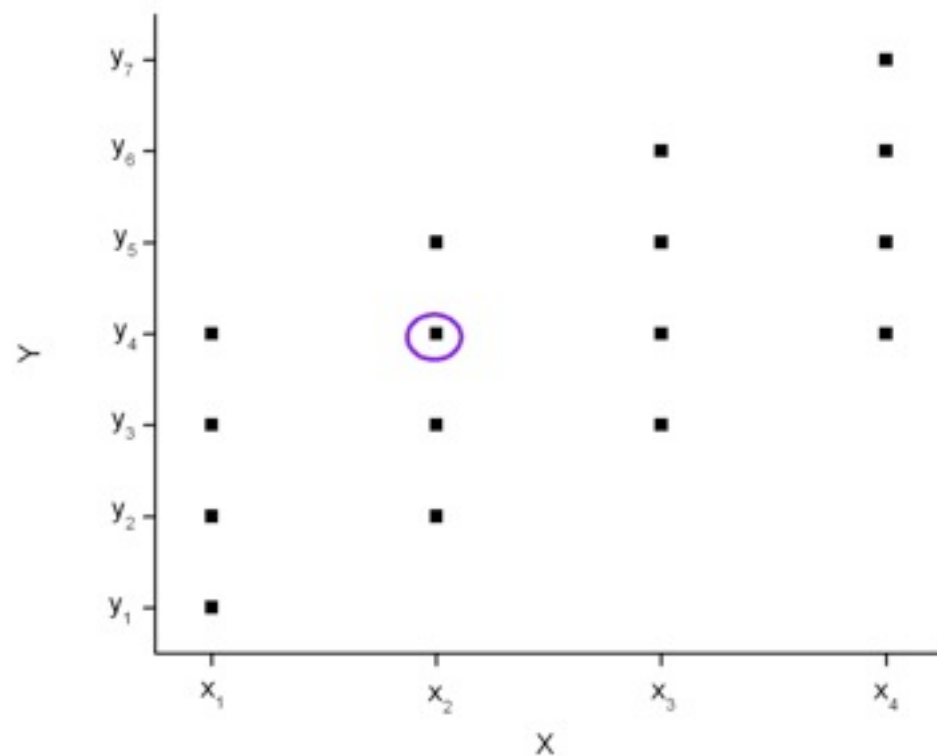
Od oka?  
Previše proizvoljno





# Korelacija u linearnoj regresiji

Pretpostavimo da postoji funkcionalna veza  $y(x) = ax+b$ , gdje su  $a$  i  $b$  koeficijenti koji nabolje opisuju eksperimentalne rezultate



Preuredimo malo notaciju:

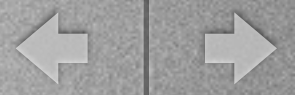
$$(x_i, y_j) \leftrightarrow (x_i, y_{ij})$$

npr. zaokruženu točku  
možemo pisati

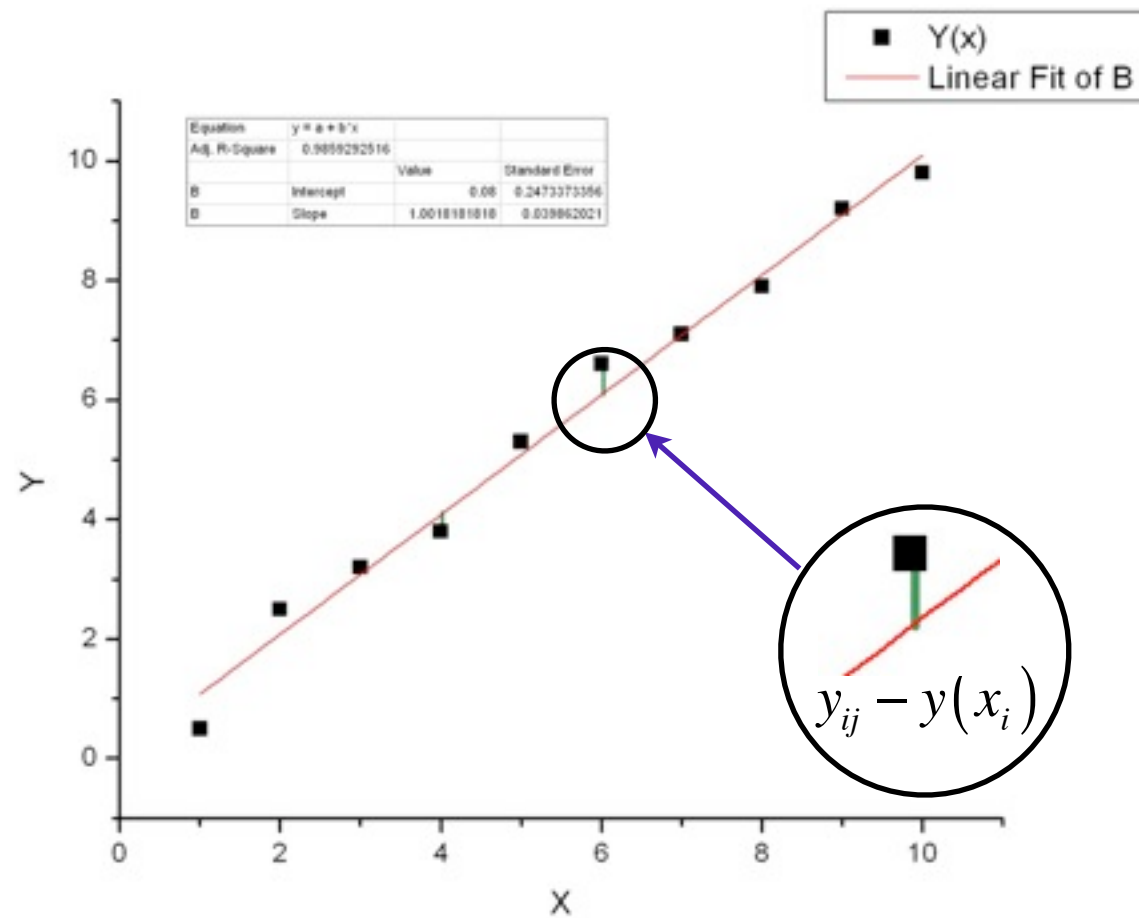
$$(x_2, y_4) \text{ ili } (x_2, y_{24})$$

Primijetimo da vrijedi:  $\sum_{i,j} f_{ij} x_i y_j = \sum_{i,j} f_{ij} x_i y_{ij}$





# Korelacija u linearnoj regresiji

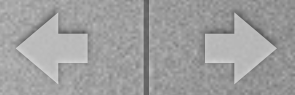


Uvodimo sume kvadrata odstupanja:

$$S^2 = \frac{1}{N} \sum_{i,j} f_{ij} (y_{ij} - y(x_i))^2$$

Moramo kvadrirati razlike, jer dobar opis  $y(x) = ax + b$  postiže se samo ako je podjednako pozitivnih i negativnih razlika, pa je:

$$\sum_{i,j} f_{ij} (y_{ij} - y(x_i)) \approx 0$$



# Korelacija u linearnoj regresiji

$$S^2 = \frac{1}{N} \sum_{i,j} f_{ij} (y_{ij} - ax_i - b)^2 \quad \longrightarrow \quad \text{za najbolji odabir, ovaj izraz je minimalan}$$

$$\frac{\partial S^2}{\partial a} = 0; \quad \frac{\partial S^2}{\partial b} = 0$$

$$\frac{\partial S^2}{\partial a} = -\frac{2}{N} \sum_{i,j} f_{ij} x_i (y_{ij} - ax_i - b) = 0$$

$$\frac{\partial S^2}{\partial b} = -\frac{2}{N} \sum_{i,j} f_{ij} (y_{ij} - ax_i - b) = 0$$





# Korelacija u linearnoj regresiji

To nam daje sustav jednažbi za  $a$  i  $b$ :

$$(1) \quad \sum_{i,j} f_{ij} y_{ij} = a \sum_{i,j} f_{ij} x_i + b \sum_{i,j} f_{ij}$$

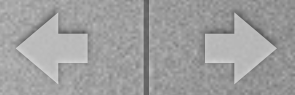
$$(2) \quad \sum_{i,j} f_{ij} y_{ij} x_i = a \sum_{i,j} f_{ij} x_i^2 + b \sum_{i,j} f_{ij} x_i$$

Prepoznamo sljedeće:

$$\sum_{i,j} f_{ij} y_{ij} = \sum_j f_j y_j = N\bar{y}; \quad \sum_{i,j} f_{ij} x_i = N\bar{x}$$

$$\sum_{i,j} f_{ij} y_{ij} x_i = N(\sigma_{xy} + \bar{x} \cdot \bar{y}); \quad \sum_{i,j} f_{ij} = N$$

$$\sum_{i,j} f_{ij} x_i^2 = \sum_i f_i x_i^2 = N\bar{x}^2$$



# Korelacija u linearnoj regresiji

(1)  $N\bar{y} = aN\bar{x} + bN \Rightarrow \bar{y} = a\bar{x} + b$   pravac ide kroz točku  $(\bar{x}, \bar{y})$

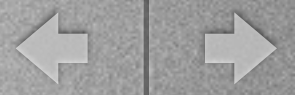
(2)  $N(\sigma_{xy} + \bar{x}\bar{y}) = (a\bar{x}^2 + b\bar{x})N$

$$\sigma_{xy} + \bar{x}\bar{y} = a\bar{x}^2 + b\bar{x}$$

To sve skupa daje:

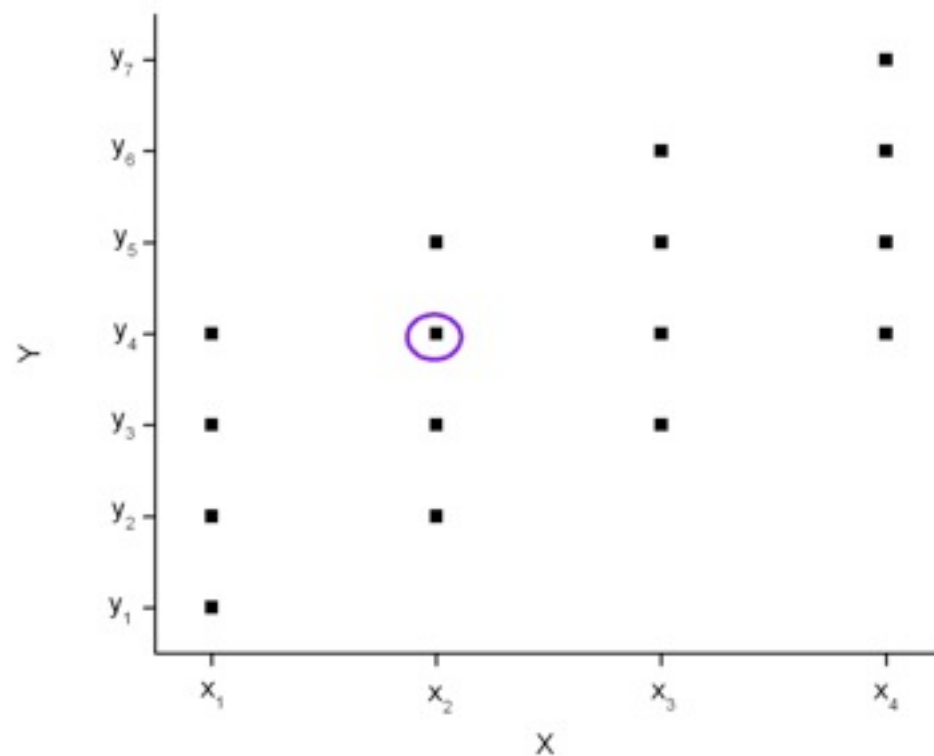
$$a = \frac{\sigma_{xy}}{x^2 - \bar{x}^2} = \frac{\sigma_{xy}}{\sigma_x^2}, \quad b = \bar{y} - \frac{\sigma_{xy}}{\sigma_x^2} \bar{x}$$





# Korelacija u linearnoj regresiji

Pravac  $y(x)=ax+b$  zovemo 'pravac  $y$  s obzirom na  $x$ ', i on pokazuje kako aritmetičke sredine točaka  $y_{ij}$  ovise o  $x$



- zaokružena točka je  $\bar{y}_2$ , tj. aritmetička sredina svih točaka  $y_{2j}$

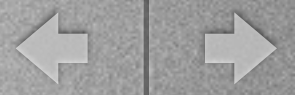


# Korelacija u linearnoj regresiji

Ekvivalentno, možemo promatrati kako  $\bar{x}_j$  ovisi o  $y$ , i tad dobivamo pravac 'x s obzirom na y'

$$x = Ay + B; \quad A = \frac{\sigma_{xy}}{\sigma_y^2} \quad B = \bar{x} - \frac{\sigma_{xy}}{\sigma_y^2} \bar{y}$$





# Korelacija u linearnoj regresiji

Pogledajmo još kolika je suma kvadrata odstupanja za  $y(x)$ :

$$\begin{aligned} S^2 &= \frac{1}{N} \sum_{i,j} f_{ij} (y_{ij} - ax_i - b)^2 = \frac{1}{N} \sum_{i,j} f_{ij} \left( y_{ij} - \frac{\sigma_{xy}}{\sigma_x^2} x_i - \bar{y} + \frac{\sigma_{xy}}{\sigma_x^2} \bar{x} \right)^2 = \\ &= \frac{1}{N} \sum_{i,j} f_{ij} \left( (y_{ij} - \bar{y}) - \frac{\sigma_{xy}}{\sigma_x^2} (x_i - \bar{x}) \right)^2 = \\ &= \frac{1}{N} \sum_{i,j} f_{ij} (y_{ij} - \bar{y})^2 + \frac{\sigma_{xy}^2}{\sigma_x^4} \frac{1}{N} \sum_{i,j} f_{ij} (x_i - \bar{x})^2 - 2 \sum_{i,j} f_{ij} (y_{ij} - \bar{y})(x_i - \bar{x}) \frac{\sigma_{xy}}{\sigma_x^2} = \\ &= \sigma_y^2 + \frac{\sigma_{xy}^2}{\sigma_x^4} \sigma_x^2 - 2 \frac{\sigma_{xy}}{\sigma_x^2} \sigma_{xy} = \sigma_y^2 - \frac{\sigma_{xy}^2}{\sigma_x^2} = \sigma_y^2 \left( 1 - \frac{\sigma_{xy}^2}{\sigma_x^2 \sigma_y^2} \right) = \sigma_y^2 (1 - r^2) \end{aligned}$$

$$S^2 = \sigma_y^2 (1 - r^2)$$



# Korelacija u linearnoj regresiji

Očito je sljedeće:

(1)  $S^2 > 0$ , jer  $|r| \leq 1$

(2) Za  $r = 1$ ,  $S^2 = 0$ , tj. korelacija je potpuna i sve točke leže na pravcu  $y = ax + b$