

# Metoda najmanjih kvadrata

## Inženjerski problem.

Predpostavimo da imamo dvije veličine:  $x$  i  $y$  ( na primjer, masa i temperatura, tlak i obujam, prva i druga koordinata točke u ravnini itd.).

Tada postoje dvije mogućnosti. Prva je da su te veličine (u određenim uvjetima) nezavisne, a druga je da su zavisne.

Nezavisnost znači da uz svaku očitane vrijednost veličine  $x$  teoretski možemo očitati bilo koju vrijednost veličine  $y$  (takve su, na primjer, masa zlata i temperatura zlata pri uobičajenim masama i temperaturama; očitana masa zlata ne daje nam nikakvu informaciju o temperaturi).

Zavisnost pak znači da očitana vrijednost veličine  $x$  potpuno određuje (ili ograničuje) vrijednost veličine  $y$ .

Na primjer, ako se točka giba po kružnici u koordinatnoj ravnini, onda očitana vrijednost varijable  $x$  općenito uvjetuje dvije moguće vrijednosti varijable  $y$ .

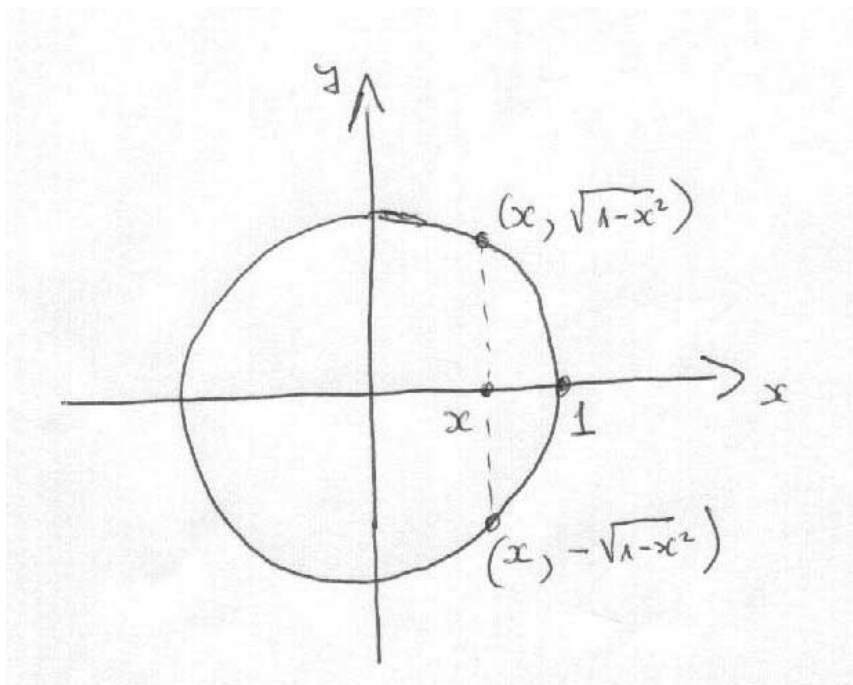
Da bi bili još konkretniji zamislimo da se točka giba po jediničnoj kružnici sa središtem u ishodištu. Predpostavimo da smo očitali prvu koordinatu i da smo dobili  $x=0.8$ .

Tada postoje dvije mogućnosti za vrijednost veličine  $y$ :  $0.6$  i  $-0.6$ .

Općenito, u ovim okolnostima veličina  $x$  može poprimiti bilo koju vrijednost između  $-1$  i  $1$  (uključujući i njih). Isto je s veličinom  $y$ . Međutim, te su dvije veličine povezane relacijom:

$$x^2 + y^2 = 1 \quad (\text{jednadžba kružnice}).$$

Zato je  $y = \pm\sqrt{1-x^2}$  pa svakoj vrijednosti veličine  $x$  odgovaraju dvije vrijednosti veličine  $y$  (osim ako je  $x=1$  ili  $x=-1$ ; tada je  $y=0$ ).



U inženjerstvu je važan slučaj kad **svaka vrijednost veličine x, uvjetuje točno jednu vrijednost veličine y**. To se kraće zapisuje pomoću funkcija:

$$y = f(x),$$

gdje je  $f$  pravilo prema kojemu  $y$  ovisi o  $x$ .

### Najjednostavnija pravila zavisnosti (funkcije)

1. linearna zavisnost  $y = ax + b$  (to je  $f(x) := ax + b$ , a grafički je prikaz te zavisnosti pravac)
2. kvadratna zavisnost  $y = ax^2 + bx + c$  (grafički prikaz je parabola).
3. kubna zavisnost  $y = ax^3 + bx^2 + cx + d$
4. recipročna zavisnost (obrnuta proporcionalnost)  $y = a/x$  (grafički prikaz je hiperbola)
4. eksponencijalna zavisnost  $y = ae^{bx}$
5. potencijnska zavisnost  $y = a \cdot b^x$   
itd.

Uočite da ima više linearnih zavisnosti (kvadratnih zavisnosti i sl.). Linearna je zavisnost poznata (određena) onda ako znamo realne brojeve  $a, b$  (koeficijente) itd. Da naglasimo kako linearna funkcija zavisi o svojim koeficijentima pišemo:  
 $f(x, a, b) := ax + b$ .

Slično, za kvadratnu funkciju:

$$f(x, a, b, c) := ax^2 + bx + c.$$

Brojeve  $a, b$  odnosno  $a, b, c$  zovemo i **parametrima**. To je i općenito, a ne samo za linearne ili kvadratne veze. Na primjer,

$$f(x, a, b) := a \cdot e^{bx}$$

je familija funkcijskih veza (eksponencijalnog tipa) ovisna o parametrima  $a, b$ .

Zamislimo pokus u kojemu mijenjamo po volji **obujam** plina, a pri svakoj konkretnoj vrijednosti obujma očitavamo **tlak**. Tu je uobičajeno da veličina  $x$  bude obujam (nezavisna veličina, veličina koju po volji mijenjamo) a da veličina  $y$  bude tlak (zavisna veličina, veličina koju dobijemo mjerenjem).

**Primjer 1.** Predpostavimo da smo mijenjali veličinu  $x$  i pri tom mjerili odgovarajuće vrijednosti veličine  $y$ . Rezultate zapišimo u obliku tablice.

$x $	1	2	3	4	5	6
$y $	2.1	5.1	8.2	11.0	14.5	17.4

Iz prvog pogleda na tablicu uočavamo da se povećanjem veličine  $x$ , povećava i veličina  $y$ . Nas zanima pravilo prema kojem se to događa. Brzo uočavamo da smo veličinu  $x$  povećavali za jednu mjernu jedinicu. Pri tom se veličina  $y$  povećavala redom za:

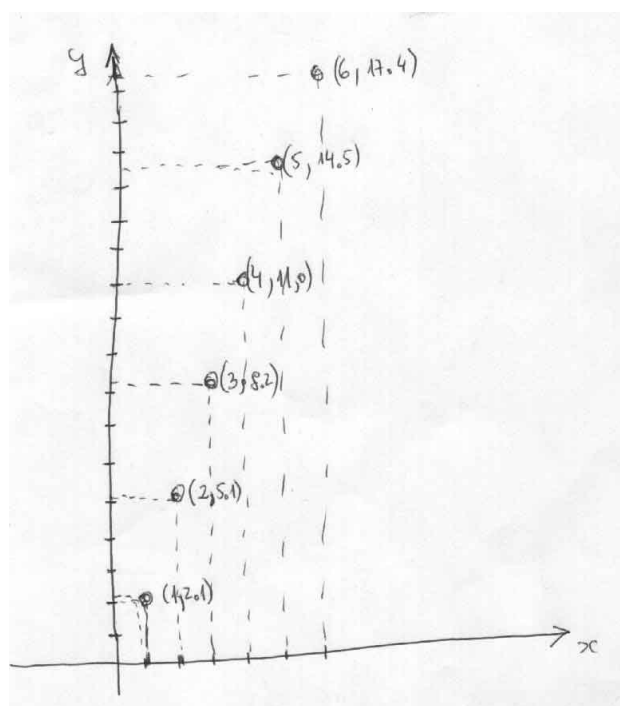
3.0 pa za 3.1 pa za 2.8 pa za 2.5 pa za 2.9.  
(vidimo da su te promjene, iako različite, ipak prilično bliske).

**Kad god pri jednakim promjenama jedne veličine uočimo odprilike jednake promjene druge veličine, moramo posumljati u linearnu vezu među njima.** Dakle:

$$y = ax + b.$$

Drugi je način uočavanja možebitne linearne veze grafički. U tu svrhu podatke zapišimo u obliku uređenih parova i ucrtajmo ih u koordinatni sustav kojemu je  $x$  horizontalna os, a  $y$  vertikalna:

(1; 2.1) , (2; 5.1), (3; 8.2), (4; 11.0), (5; 14.5), (6; 17.4)



Uočavamo da su točke odprilike na pravcu. Grafički uvid u oblik zavisnosti možemo provesti i onda ako razmaci među vrijednostima veličine  $x$  nisu jednaki (nisu ekvidistantni), dok je takav uvid računanjem promjene veličine  $y$  otežan (iako se i on može provesti).

Sjetite se da su za crtanje pravca potrebne samo dvije točke, a mi ih imamo 5. Lako se vidi da ne postoji pravac koji prolazi kroz sve te točke. Zato od svih pravaca treba izabrati onaj koji *najbolje* prolazi pokraj tih točaka. Kriterij odabiranja tog pravca, tj. pripadajućih koeficijenata  $a, b$  daje *metoda najmanjih kvadrata*. Tako je i općenito, samo što je onda kad

uočena zavisnost odudara od linearne, otežano biranje tipa zavisnosti. Često to i nije problem. Naime tip veze je u mnogim primjerima predviđen pripadajućim teorijama. Na primjer, veza između obujma i tlaka plina pri fiksiranoj temperaturi predviđena je u obliku Van-der- Wallsove jednadžbe, Peng-Robinsonove jednadžbe i sl. (a mi samo moramo odrediti parametre koji se pojavljuju u tim jednadžbama).

Općenito imamo  $n$  vrijednosti veličine  $x$  i  $n$  vrijednosti zavisne veličine  $y$ . To zadajemo tablicom:

$x $	$x_1$	$x_2$	....	$x_n$
$y $	$y_1$	$y_2$	.....	$y_n$

### Načelo na kojemu se zasniva metoda najmanjih kvadrata.

Metoda najmanjih kvadrata zasniva se na načelu da su najbolji oni parametri  $a, b$  za koje je suma kvadrata razlika između mjerenih vrijednosti  $y_i$ ,  $i=1,2,\dots,n$  i izračunatih vrijednosti  $f(x_i, a, b)$  minimalna.

Ima više matematičkih razloga za prihvaćanje ovog načela, a tu ih nećemo spominjati. Napomenimo da nije dobro razmatrati zbroj razlika eksperimentalnih i teoretskih podataka jer se pozitivne i negativne razlike (**odstupanja**) poništavaju. Da bi uzeli u obzir i pozitivna i negativna odstupanja, matematičari su na početku razmatrali apsolutne vrijednosti razlika i tražili da njihova suma bude minimalna. To nije loš kriterij, ali su apsolutne vrijednosti nepogodne jer se ne mogu općenito derivirati. Taj, ali i neki drugi razlozi, prevagnuli su u korist sume kvadrata.

### Postupak određivanja parametara metodom najmanjih kvadrata.

Označimo  $i$ -to odstupanje kao  $D_i := y_i - f(x_i, a, b)$

To je razlika između mjerene (eksperimentalne) vrijednosti  $y_i$  i teoretske vrijednosti  $f(x_i, a, b)$ , tj. vrijednosti funkcije  $f(x, a, b)$  za  $x=x_i$ .

Prema metodi najmanjih kvadrata, parametre određujemo tako da suma

$$D_1^2 + D_2^2 + \dots + D_n^2$$

bude minimalna.

Taj izraz ovisi o nepoznatim parametrima  $a, b$ . Zato pišemo:

$$F(a, b) := [y_1 - f(x_1, a, b)]^2 + [y_2 - f(x_2, a, b)]^2 + \dots + [y_n - f(x_n, a, b)]^2$$

ili, kraće

$$F(a, b) := \sum_{i=1}^n [y_i - f(x_i, a, b)]^2$$

( $F$  ovisi i o vrijednostima  $x_i, y_i$ , za  $i=1,2,\dots,n$ , međutim te su vrijednosti poznate). Funkcija  $F$  često se naziva **funkcija cilja** (općenito, funkcija cilja može biti i neka druga pogodna funkcija, na primjer, odstupanja se mogu množiti nekim težinama).

Treba odrediti parametre  $a, b$  u kojima funkcija cilja  $F$  postiže minimum.

Uvjeti lokalnog ekstrema (pomoću parcijalnih derivacija) za funkciju  $F$  jesu:

$$\frac{\partial F}{\partial a} = 0 \quad \text{i} \quad \frac{\partial F}{\partial b} = 0$$

Dakle

$$\frac{\partial(\sum_{i=1}^n [y_i - f(x_i, a, b)]^2)}{\partial a} = 0 \quad \text{i} \quad \frac{\partial(\sum_{i=1}^n [y_i - f(x_i, a, b)]^2)}{\partial b} = 0$$

Koristeći svojstva derivacija (derivacija zbroja je zbroj derivacija) dobijemo:

$$\sum_{i=1}^n 2[y_i - f(x_i, a, b)] \cdot \left(-\frac{\partial f(x_i, a, b)}{\partial a}\right) = 0 \quad \text{i} \quad \sum_{i=1}^n 2[y_i - f(x_i, a, b)] \cdot \left(-\frac{\partial f(x_i, a, b)}{\partial b}\right) = 0$$

Nakon sređivanja dobijemo sustav dviju jednačja s dvjema nepoznicama a, b.

$$\sum_{i=1}^n [y_i - f(x_i, a, b)] \cdot \frac{\partial f(x_i, a, b)}{\partial a} = 0$$

$$\sum_{i=1}^n [y_i - f(x_i, a, b)] \cdot \frac{\partial f(x_i, a, b)}{\partial b} = 0 \quad (*)$$

Iz tog sustava određujemo nepoznate parametre a, b. Općenito sustav može imati više rješenja. Također može se dogoditi da neka rješenja odgovaraju maksimumu ili sedlastoj točki, a ne minimumu. Može se dogoditi i to da neka rješenja nemaju fizikalna značenja. Na sreću, u najvažnijem slučaju, slučaju linearnih veza i njima srodnih, rješenje tog sustava je jedinstveno, tj. parametri se mogu odrediti jednoznačno.

### Linearna regresija.

**Određivanje parametara a, b za linearnu vezu (određivanje regresijskog pravca).**

Tu je  $f(x, a, b) := ax + b$ , pa je

$$f(x_i, a, b) = ax_i + b$$

$$\frac{\partial f(x_i, a, b)}{\partial a} = x_i \quad \text{i} \quad \frac{\partial f(x_i, a, b)}{\partial b} = 1$$

Ako to uvrstimo u sustav (\*), dobijemo

$$\sum_{i=1}^n (y_i - ax_i - b) \cdot x_i = 0$$

$$\sum_{i=1}^n (y_i - ax_i - b) \cdot 1 = 0$$

nakon raspisivanja dobijemo sustav dviju linearnih jednačja s dvjema nepoznicama:

$$\left(\sum_{i=1}^n x_i^2\right) \cdot a + \left(\sum_{i=1}^n x_i\right) \cdot b = \sum_{i=1}^n x_i y_i$$

$$\left(\sum_{i=1}^n x_i\right) \cdot a + n \cdot b = \sum_{i=1}^n y_i$$

U tom su sustavu parametri a, b nepoznanice, a brojevi

$\sum_{i=1}^n x_i^2$ ,  $\sum_{i=1}^n x_i$ ,  $n$ ,  $\sum_{i=1}^n x_i y_i$ ,  $\sum_{i=1}^n x_i$  **koeficijenti sustava** (oni su poznati, jer se dobiju iz mjernih podataka).

Rješavanjem tog linearnog sustava dobijemo konačne formule (indekse ispuštamo!):

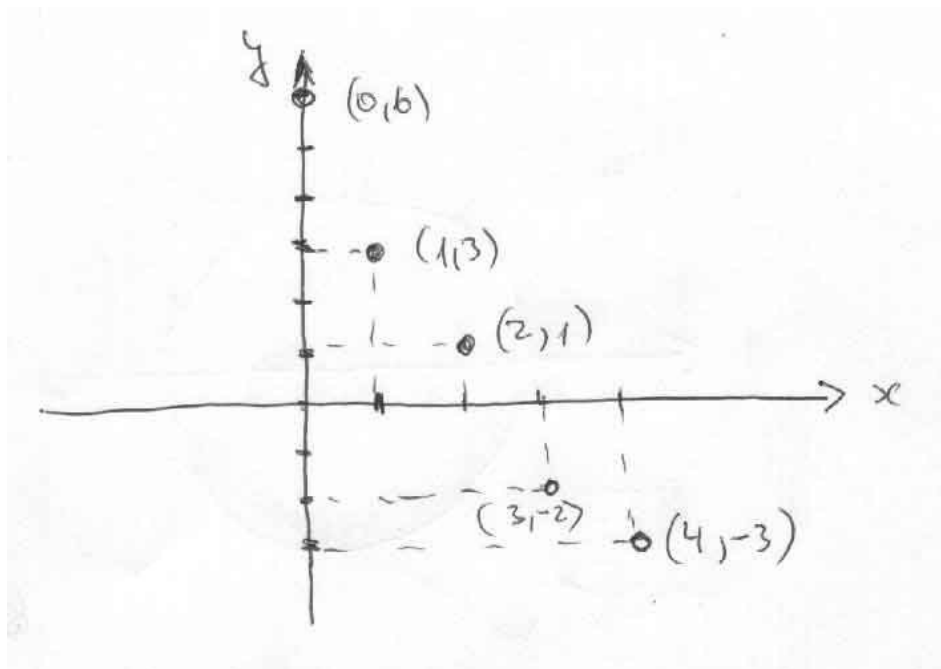
$$a = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \cdot \sum x_i^2 - (\sum x_i)^2}, \quad b = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \cdot \sum x_i^2 - (\sum x_i)^2} \quad (**)$$

Dobiveni pravac s jednadžbom  $y = ax+b$  zove se **regresijski pravac**.

**Primjer 2.** Procijenimo oblik veze, odredimo parametre, jednadžbu regresijskog pravca, odstupanja i vrijednost funkcije cilja ako je:

$x_i$		0	1	2	3	4
$y_i$		6	3	1	-2	-3

Točke (0,6), (1,3), (2,1), (3,-2), (4,-3) predočimo u koordinatnom sustavu.



Vidimo da su točke približno na pravcu, pa tražimo linearnu vezu, tj. vezu oblika  $y=ax+b$ . Unaprijed vidimo da je  $a < 0$  i da je  $b \approx 6$ . Da bismo lakše računali  $a, b$  iz (\*\*), izradimo ovakvu tablicu (u posljednjem su stupci zbrojevi elemenata u odgovarajućem reduku):

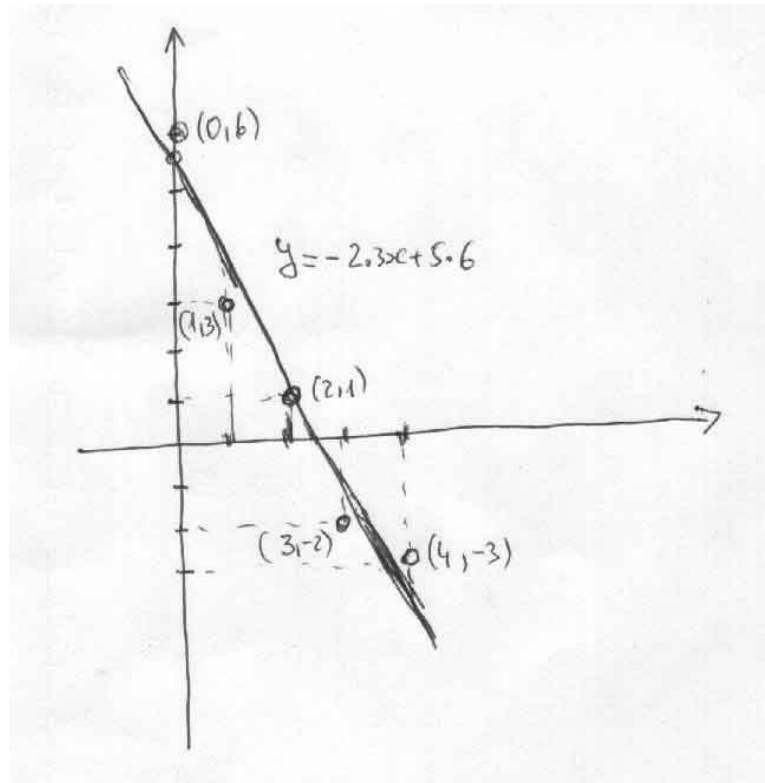
$x_i$		0	1	2	3	4		10
$y_i$		6	3	1	-2	-3		5
$x_i^2$		0	1	4	9	16		30
$x_i y_i$		0	3	2	-6	-12		-13

Tu je još  $n=5$ , pa iz (\*\*) dobijemo:

$$a = \frac{5 \cdot (-13) - 10 \cdot 5}{5 \cdot 30 - 10^2} = \frac{-115}{50} = -2.3$$

$$b = \frac{30 \cdot 5 - 10 \cdot (-13)}{5 \cdot 30 - 10^2} = \frac{280}{50} = 5.6$$

Dakle, tražena je linearna veza:  $y = -2.3x + 5.6$  (to je jednadžba regresijskog pravca).



Da bismo odredili odstupanja i vrijednost funkcije cilja, izračunajmo najprije vrijednosti funkcije  $f(x) := -2.3x + 5.6$ , redom za  $x=0,1,2,3,4$ .

$$f(0) = 5.6$$

$$f(1) = 3.3$$

$$f(2) = 1.0$$

$$f(3) = -1.3$$

$$f(4) = -3.6$$

Vidimo da su teoretski rezultati bliski eksperimentalnim podacima, što pokazuje i tablica (u posljednjem su reduku **odstupanja**  $D_i := y_i - f(x_i)$ ).

$x_i$	0	1	2	3	4
$y_i$	6	3	1	-2	-3
$f(x_i)$	5.6	3.3	1.0	-1.3	-3.6
$D_i$	0.4	-0.3	0.0	-0.7	0.6

Uočite da je **zbroj odstupanja jednak nuli** (to vrijedi općenito:  $\sum D_i = 0$ ).

Minimalna vrijednost funkcije cilja jednaka je zbroju kvadrata odstupanja:

$$\sum D_i^2 = 0.4^2 + (-0.3)^2 + 0.0^2 + (-0.7)^2 + 0.6^2 = 1.10.$$

**Napomena.** Za veliki broj podataka provođenje metode najmanjih kvadrata može biti mukotržno. Zato je dobro naučiti primjenu grafičkog kalkulatora ili nekog dostupnog kompjutorskog programa. Na primjer, naredba LinReg na grafičkom kalkulatoru, za podatke iz Primjera 1. daje nam (zaokruženo na tri decimale)  $a=3.071$  i  $b = -1.033$ .

Dobivena linearna veza može nam poslužiti za procjenu (približno određivanje) vrijednosti veličine  $y$  za vrijednosti  $x$  unutar područja mjerenja (**interpolacija**) ili izvan njega (**ekstrapolacija**)

**Primjer 3.** Procijenimo vrijednost veličine  $x$  iz predhodnog primjera za  $x=2.5$  i  $x=5$ .

Iz  $f(x)=-2.3x+5.6$ , dobijemo  $f(2.5)=-2.3 \cdot 2.5+5.6 = -0.15$ .

Dakle interpolacijom smo dobili da vrijednosti  $x=2.5$  približno odgovara vrijednost  $y= -0.15$  (koju, inače, nismo mjerili). Uočite da se ta vrijednost razlikuje od srednje vrijednosti u  $x=2$  i  $x=3$  (koja je jednaka  $\frac{2+(-3)}{2}=-0.5$ ).

Takodjer se dobije  $f(5)=-2.3 \cdot 5+5.6= -5.9$ , pa je  $y=-5.9$  približna vrijednost veličine  $y$  za  $x=5$  (kažemo da smo je dobili ekstrapolacijom).

Uočimo da u Primjeru 2. dobiveni regresijski pravac prolazi podatkom  $(2,1)$ . Općenito, on ne mora prolaziti ni kroz jednu točku. Jedan od razloga zašto se tu tako dogodilo jest taj što je 2 aritmetička sredina vrijednosti  $x$  veličine, a 1 aritmetička sredina vrijednosti  $y$  veličine (provjerite; to je slučaj, sredina skupa podataka općenito nije podatak). Naime, regresijski pravac uvijek prolazi točkom  $(\bar{x}, \bar{y})$ , gdje je

$$\bar{x} := \frac{x_1 + x_2 + \dots + x_n}{n}, \quad \bar{y} := \frac{y_1 + y_2 + \dots + y_n}{n}. \quad \text{Dakle } \bar{y} = a\bar{x} + b.$$

## Linearna korelacija

Korelacija je mjera linearne zavisnosti dviju serija podataka  $x_1, x_2, \dots, x_n$  i  $y_1, y_2, \dots, y_n$ .

Drugim riječima, ako su točke  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  grupirane oko regresijskog pravca, onda govorimo da su podatci **korelirani (linearno korelirani)**. Na osnovi toga govori se da su pripadne veličine  $x, y$  korelirane. Razina koreliranosti mjeri se **koeficijentom korelacije**

$$r := \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \cdot \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

Vrijedi:

(1)  $-1 \leq r \leq 1$

(2) Ako je  $r$  pozitivan onda je  $i$  koeficijent smjera regresijskog pravca pozitivan i obratno (ako je  $r > 0$  onda je  $a > 0$ , a ako je  $r < 0$  onda je  $a < 0$ )

(3) Što je  $r$  bliži 1 ili -1 to su veličine značajnije linearno korelirane, tj. podatci su bliže regresijskom pravcu, a što je  $r$  bliži 0, podatci su razbacaniji.

(4) Ako je  $r = 1$  ili  $r = -1$  onda su sve točke  $(x_i, y_i)$  na regresijskom pravcu, tj. podatci su potpuno linearno zavisni.

(5) Ako je  $r = 0$  onda nema nikakve linearne zavisnosti među veličinama.



**Primjer 1.** Odredimo koeficijent korelacije za podatke

$x_i$		0	1	2	3	4
$y_i$		6	3	1	-2	-3

iz Primjera 2. u lekciji Metoda najmanjih kvadrata.

Napravimo tablicu poput one za metodu najmanjih kvadrata, samo dodajmo redak s  $y_i^2$ .

$x_i$		0	1	2	3	4		10
$y_i$		6	3	1	-2	-3		5
$x_i^2$		0	1	4	9	16		30
$x_i y_i$		0	3	2	-6	-12		-13
$y_i^2$		36	9	1	4	9		59

$$r = \frac{5 \cdot (-13) - 5 \cdot 10}{\sqrt{5 \cdot 30 - 10^2} \cdot \sqrt{5 \cdot 59 - 5^2}} = \frac{-115}{30\sqrt{15}} \approx -0.98976 \text{ (na pet decimala)}$$

što je visoka razina linearne zavisnosti. Uočite da je  $r < 0$ , što je u skladu s tim da je  $a < 0$  (sjetite se da je jednadžba regresijskog pravca  $y = -2.3x + 5.6$ ).

**Napomena.** Naredba LinReg uz podatke o parametrima također bi nam izbacila i vrijednost koeficijenta  $r$ .

Podatci u sljedećem primjeru su vrlo nisko linearno korelirani.

**Primjer 2.** Odredimo koeficijent korelacije za podatke

$x_i$		-3	-2	-1	0	1	2	3
$y_i$		0	2	-2	3	2	1	-1

Unesimo podatke u kvadratnu  $7 \times 7$  mrežu kao na slici. Vidimo da podatci ne prate ni jedan pravac, pa procjenjujemo da je koeficijent korelacije blizu nule.

$x_i$		-3	-2	-1	0	1	2	3		0
$y_i$		0	2	-2	3	2	1	-1		5
$x_i^2$		9	4	1	0	1	4	9		28
$x_i y_i$		0	-4	2	0	2	2	-3		-1
$y_i^2$		0	4	4	9	4	1	1		23

$$\text{Sad je } r = \frac{7 \cdot (-1) - 0 \cdot 0}{\sqrt{7 \cdot 28 - 0^2} \cdot \sqrt{7 \cdot 23 - 5^2}} \approx -0.042875 \text{ (na šest decimala)}$$

što je praktički jednako nuli. Također, može se provjeriti da je pripadna suma kvadrata odstupanja za linearnu regresiju jednaka 19.392857 (na šest decimala), što također upućuje na vrlo slabu linearnu vezu.

Linearnu korelaciju ne treba shvatiti kao jedini oblik zavisnosti dviju veličina (serija podataka). Dvije veličine mogu biti vrlo jasno zavisne, a da im je koeficijent (linearne) korelacije jednak nuli; to samo znači da su one linearno nekorelirane. To pokazuje sljedeći primjer.

**Primjer 3.** Odredimo koeficijent korelacije za podatke

$x_i$		-3	-2	-1	0	1	2	3
$y_i$		9	4	1	0	1	4	9

Ucrtavanjem podataka vidimo da oni ne prate ni jedan pravac. Kako je  $\sum x_i = 0$  i  $\sum x_i y_i = 0$ , vidimo da je  $r=0$ . Dakle, podaci su linearno nekorelirani. S druge strane, oni su zavisni. Naime, povezani su relacijom  $y = x^2$  (točke su na paraboli).

Često se postavlja pitanje koji koeficijenti znače visoku, koji nisku, a koji srednju linearnu koreliranost. Na to pitanje nema jasnog odgovora. On ovisi i o znanstvenom području na koje se primjenjuje, a unutar znanstvenog područja na konkretan problem koji se razmatra. Na primjer, u psihologijskim istraživanjima, u pravilu, čim je  $r > 0.5$  smatra se da je koreliranost značajna, a ako je  $r > 0.8$  vrlo značajna, dok u preciznim fizikalnim ili kemijskim istraživanjem često niti  $r = 0.9$  ne upućuje na značajnu koreliranost.

**Primjer 4.** U sljedećoj tablici su u prvom redu bodovi prvih 9 najboljih rezultata postignutih iz kolokvija na Matematici 1, a u drugoj su odgovarajući bodovi iz Matematike 2.

$x_i$		103	93	84	81	81	80	79	79	78
$y_i$		99	73	82	85	77	79	73	55	83

Odredimo regresijski pravac i koeficijent korelacije. Komentirajmo rezultate.

Da dobijemo predodžbu, podatke predočavamo u koordinatnom sustavu.

Predviđamo pozitivan koeficijent korelacije jer podatci prate glavnu dijagonalu (ali ne visok, jer podatci variraju). Procijenjujemo da je koeficijent regresijskog pravca nešto manji od 1. Zadatak se može izraditi prema uzoru na prijašnje primjere. Mi ćemo se poslužiti grafičkim kalkulatorom, koji ima gotov program za metodu najmanjih kvadrata i linearnu korelaciju. Dobijemo, zaokružujući na dvije decimale,  $y = 0.79x + 12.29$  i  $r = 0.56$ .

**Komentar.** Dobili smo  $a = 0.79 < 1$ , što je u skladu s činjenicom da su rezultati Matematike 2, nešto niži od rezultata Matematike 1. Koeficijent korelacije nije blizu 1, ali je veći od 0.5, što, pri ovakvoj problematici upućuje na nezanemarivu korelaciju. Suma kvadrata odstupanja jednaka je, na četiri decimale, 763.7222 što izgleda veliko, ali taj rezultat treba tumačiti tako da je prosječno odstupanje oko 9 bodova, što i nije tako veliko.

**Zašto nas je u ovom primjeru zanimala linearna korelacija među rezultatima?**

Zato što intuitivno prihvaćamo da će rezultati iz Matematike 2 biti približno proporcionalni onima iz Matematike 1, tj. da odprilike jednake razine usvajanja nekog znanja uvjetuju odprilike jednake razine usvajanja novog znanja koje počiva na starom. Naravno da to ne vrijedi za svakog konkretnog pojedinca, već u prosjeku.