

# Curs 20

Cristian Niculescu

## 1 Regresia liniară

### 1.1 Scopurile învățării

1. Să poată folosi metoda celor mai mici pătrate pentru a potrivi o dreaptă cu date bivariate.
2. Să poată da o formulă pentru eroarea pătratică totală la potrivirea oricărui tip de curbă cu datele.
3. Să poată spune cuvintele homoscedasticitate și heteroscedasticitate.

### 1.2 Introducere

Presupunem că avem colectate date bivariate  $(x_i, y_i), i = 1, \dots, n$ . Scopul regresiei liniare este modelarea relației dintre  $x$  și  $y$  prin aflarea unei funcții  $y = f(x)$  care dă o potrivire apropiată cu datele. Presupunerile modelării pe care le vom folosi sunt că  $x_i$  **nu** sunt aleatoare și că  $y_i$  este o funcție de  $x_i$  plus un zgomot aleator. Cu aceste presupuneri,  $x$  este numită **variabilă independentă** sau **predictor** și  $y$  este numită variabilă **dependentă** sau **răspuns**.

**Exemplul 1.** Costul unui timbru de clasa întâi în dolari de-a lungul timpului este dat în lista următoare:

.05 (1963)	.06 (1968)	.08 (1971)	.10 (1974)	.13 (1975)	.15 (1978)	.20 (1981)	.22 (1985)
.25 (1988)	.29 (1991)	.32 (1995)	.33 (1999)	.34 (2001)	.37 (2002)	.39 (2006)	.41 (2007)
.42 (2008)	.44 (2009)	.45 (2012)	.46 (2013)	.49 (2014)			

Folosind codul R:

```
x=c(3,8,11,14,15,18,21,25,28,31,35,39,41,42,46,47,48,49,52,53,54)
```

```
y=c(5,6,8,10,13,15,20,22,25,29,32,33,34,37,39,41,42,44,45,46,49)
```

```
lm(y~x),
```

```
obținem:
```

```
Call:
```

```
lm(formula = y ~ x)
```

```
Coefficients:
```

```
(Intercept) x
```

-0.1324 0.8791.

Am aflat că dreapta care dă "potrivirea celor mai mici pătrate" cu aceste date (dreapta de regresie) este

$$y = -0.1324 + 0.8791x,$$

unde  $x$  este numărul de ani de la 1960, iar  $y$  este în cenți.

Folosind acest rezultat "prezicem" că în 2021 ( $x = 61$ ), costul unui timbru va fi 53 de cenți (deoarece  $-0.1324 + 0.8791 \cdot 61 = 53.4927$ ).

Folosind codul R (în continuarea celui de mai sus):

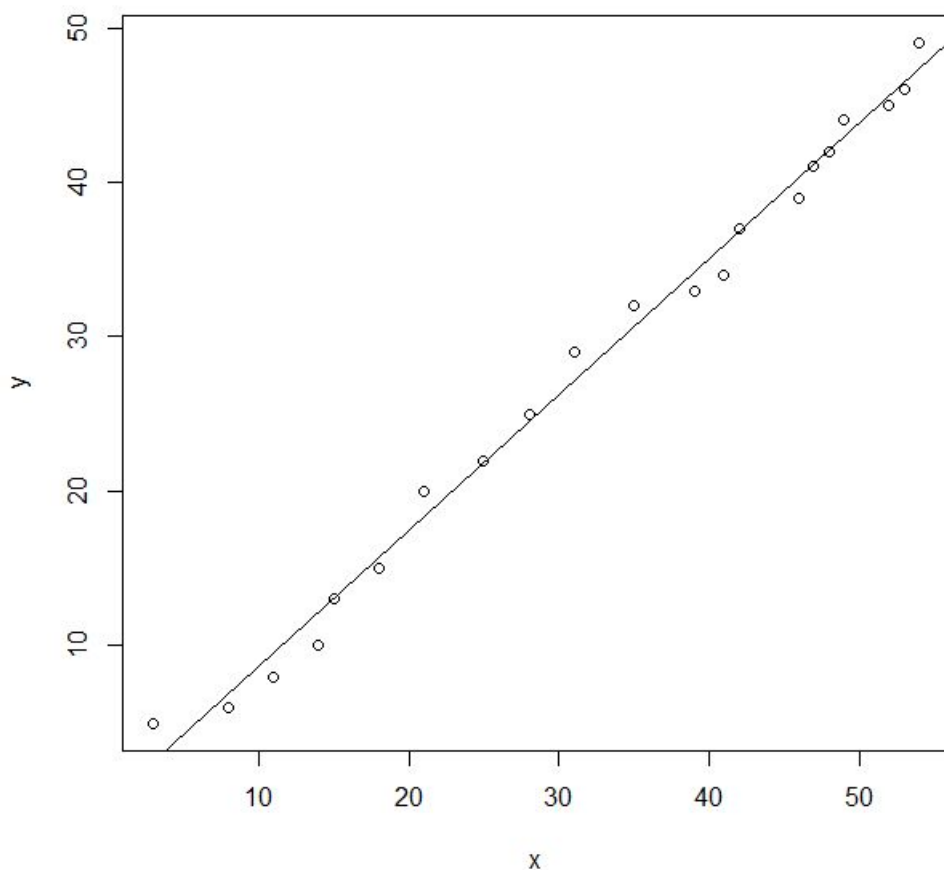
```
a=-0.1324
```

```
b=0.8791
```

```
plot(x,y)
```

```
abline(a,b)
```

obținem:



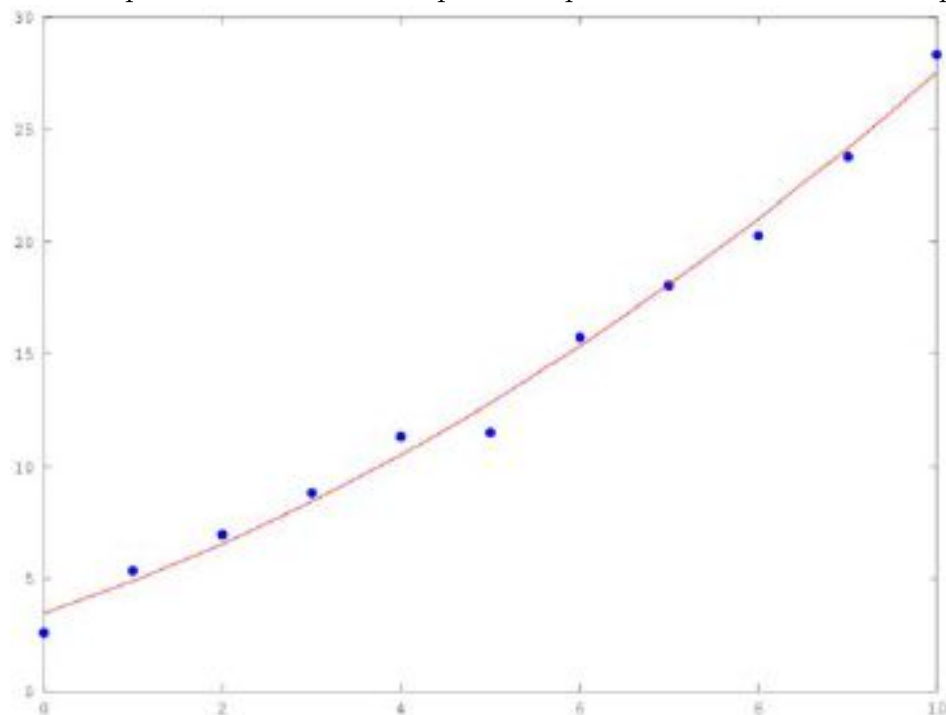
Costul timbrului (cenți) vs. timp (ani din 1960)

Niciuna din date nu se află chiar pe dreaptă. Mai degrabă această dreaptă are "cea mai bună potrivire" în raport cu [toate datele](#), cu o mică eroare pentru fiecare dată.

**Exemplul 2.** Presupunem că avem  $n$  perechi de tați și fii adulți. Fie  $x_i$  și  $y_i$  înălțimile celui de-al  $i$ -lea tată, respectiv fiu. Dreapta celor mai mici pătrate

pentru aceste date poate fi folosită pentru a prezice înălțimea de adult a unui băiat tânăr din cea a tatălui lui.

**Exemplul 3.** Nu suntem limitați la drepte cu cea mai bună potrivire.  $\forall d \in \mathbb{N}^*$ , metoda celor mai mici pătrate poate fi folosită pentru a afla un polinom de grad  $d$  cu "cea mai bună potrivire cu datele". Iată o figură arătând potrivirea datelor cu o parabolă prin metoda celor mai mici pătrate:



Potrivirea unei parabole,  $y = b_2x^2 + b_1x + b_0$  cu datele

### 1.3 Potrivirea unei drepte folosind cele mai mici pătrate

Presupunem că avem datele  $(x_i, y_i)$  ca mai sus. Scopul este să aflăm dreapta

$$y = \beta_1x + \beta_0,$$

care "se potrivește cel mai bine" cu datele. Modelul nostru spune că fiecare  $y_i$  este prezis de  $x_i$  până la o eroare  $\epsilon_i$ :

$$y_i = \beta_1x_i + \beta_0 + \epsilon_i.$$

Deci,

$$\epsilon_i = y_i - \beta_1x_i - \beta_0.$$

Metoda celor mai mici pătrate află valorile  $\hat{\beta}_0$  și  $\hat{\beta}_1$  ale lui  $\beta_0$  și  $\beta_1$  care minimizează suma pătratelor erorilor:

$$S(\beta_0, \beta_1) = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_1 x_i - \beta_0)^2.$$

Folosind analiza matematică (detalii în adaos), aflăm

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \quad (1)$$

unde

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad s_{xx} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

Aici,  $\bar{x}$  este media de selecție a lui  $x$ ,  $\bar{y}$  este media de selecție a lui  $y$ ,  $s_{xx}$  este dispersia de selecție a lui  $x$  și  $s_{xy}$  este covarianța de selecție a lui  $x$  și  $y$ .

**Exemplul 4.** Folosiți cele mai mici pătrate pentru a potrivi cu o dreaptă următoarele date: (0,1), (2,1), (3,4).

**Răspuns.** În cazul nostru,  $(x_1, y_1) = (0, 1)$ ,  $(x_2, y_2) = (2, 1)$  și  $(x_3, y_3) = (3, 4)$ . Deci

$$\bar{x} = \frac{1}{3}(x_1 + x_2 + x_3) = \frac{1}{3}(0 + 2 + 3) = \frac{5}{3},$$

$$\bar{y} = \frac{1}{3}(y_1 + y_2 + y_3) = \frac{1}{3}(1 + 1 + 4) = \frac{6}{3} = 2,$$

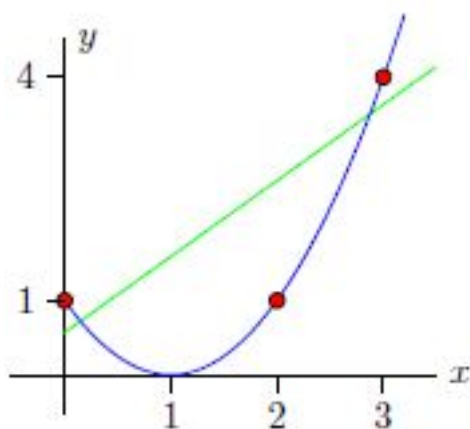
$$s_{xx} = \frac{1}{3-1} \sum_{i=1}^3 (x_i - \bar{x})^2 = \frac{1}{2} \left[ \left(0 - \frac{5}{3}\right)^2 + \left(2 - \frac{5}{3}\right)^2 + \left(3 - \frac{5}{3}\right)^2 \right] = \frac{7}{3},$$

$$s_{xy} = \frac{1}{2} \left[ \left(0 - \frac{5}{3}\right)(1 - 2) + \left(2 - \frac{5}{3}\right)(1 - 2) + \left(3 - \frac{5}{3}\right)(4 - 2) \right] = 2;$$

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} = \frac{2}{\frac{7}{3}} = \frac{6}{7};$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 2 - \frac{6}{7} \cdot \frac{5}{3} = 2 - \frac{10}{7} = \frac{4}{7}.$$

Deci dreapta de regresie a celor mai mici pătrate are ecuația  $y = \frac{4}{7} + \frac{6}{7}x$ . Aceasta este arătată ca dreapta verde din figura următoare.



Potrivirea celor mai mici pătrate a unei drepte (verde) și a unei parabole (albastru)

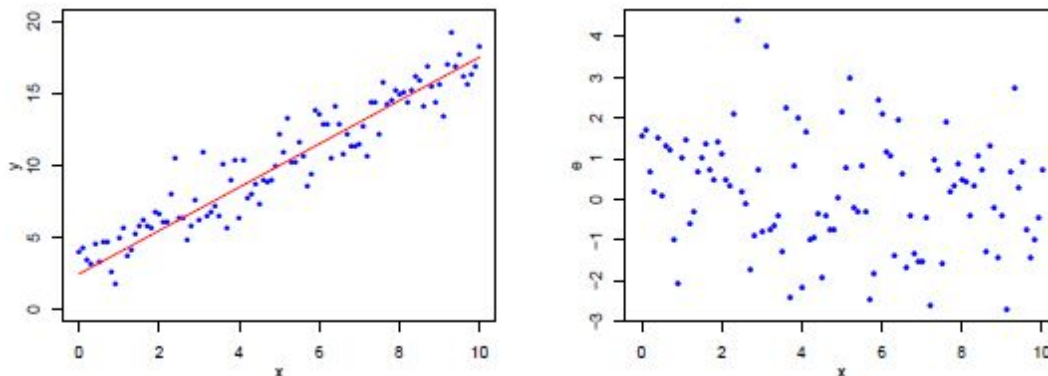
**Regresie liniară simplă:** Este puțin confuz, dar cuvântul "liniară" din "regresie liniară" nu se referă la potrivirea unei drepte. Totuși, cea mai uzuală curbă pentru potrivire este o dreaptă. Potrivirea unei drepte la date bivariate este numită [regresie liniară simplă](#).

### 1.3.1 Reziduuri

Pentru o dreaptă, modelul este

$$y_i = \hat{\beta}_1 x_i + \hat{\beta}_0 + \epsilon_i.$$

Gândim  $\hat{\beta}_1 x_i + \hat{\beta}_0$  ca prezicând sau explicând  $y_i$ . Termenul rămas  $\epsilon_i$  este numit [reziduul](#), pe care-l gândim ca pe un zgomot aleator sau o eroare de măsurare. O verificare vizuală folositoare a modelului de regresie liniară este reprezentarea reziduurilor. Datele ar trebui să fie lângă dreapta de regresie. Reziduurile ar trebui să arate cam la fel de-a lungul domeniului lui  $x$ .

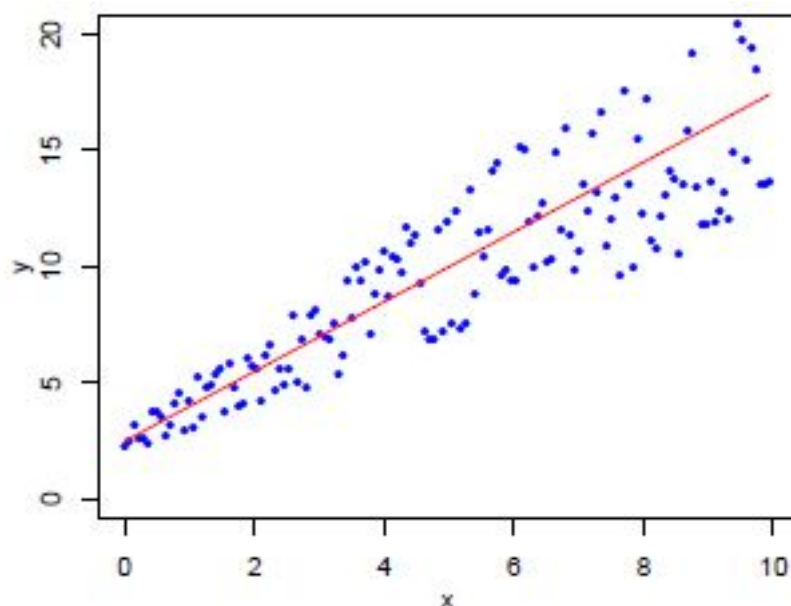


Date cu dreapta de regresie (stânga) și reziduuri (dreapta). Observați homoscedasticitatea.

### 1.3.2 Homoscedasticitatea

O presupunere importantă a modelului de regresie liniară este că reziduurile  $\epsilon_i$  au aceeași dispersie  $\forall i$ . Această presupunere este numită **homoscedasticitate**. Puteți vedea aceasta în cazul ambelor figuri de mai sus. Datele sunt în banda de lățime fixă în jurul dreptei de regresie și la fiecare  $x$  reziduurile au cam aceeași împrăștiere verticală.

Mai jos este o figură arătând date **heteroscedastice**. Împrăștierea verticală a datelor crește când  $x$  crește. Înainte de a folosi cele mai mici pătrate pe aceste date ar trebui să transformăm datele pentru a fi homoscedastice.



Date heteroscedastice

## 1.4 Regresie liniară pentru potrivirea polinoamelor

Potrivirea unei drepte la date este numită **regresie liniară simplă**. Putem de asemenea folosi regresia liniară pentru a potrivi polinoame cu datele. Folosirea cuvântului "liniară" în ambele cazuri poate părea confuză. Aceasta este deoarece cuvântul "liniară" din "regresia liniară" nu se referă la potrivirea unei drepte. Mai degrabă se referă la ecuațiile algebrice liniare pentru parametrii necunoscuți  $\beta_i$ , i.e. fiecare  $\beta_i$  are exponentul 1.

**Exemplul 5.** Luați aceleași date ca în exemplul 4 și folosiți cele mai mici pătrate pentru a afla parabola cu cea mai bună potrivire pentru date.

**Răspuns.** O parabolă are formula  $y = \beta_0 + \beta_1 x + \beta_2 x^2$ . Eroarea pătratică este

$$S(\beta_0, \beta_1, \beta_2) = \sum_{i=1}^3 (y_i - (\beta_0 + \beta_1 x_i + \beta_2 x_i^2))^2.$$

După substituirea valorilor date pentru  $x_i$  și  $y_i$  putem folosi analiza matematică (egalăm derivatele parțiale în raport cu  $\beta_0, \beta_1, \beta_2$  cu 0, obținând un sistem de 3 ecuații liniare cu 3 necunoscute) pentru a afla tripletul  $(\beta_0, \beta_1, \beta_2)$  care minimizează  $S$ . Sau putem folosi codul R

```
x=c(0,2,3)
y=c(1,1,4)
C=cbind(1,x,x^ 2)
solve(t(C)%*%C,t(C)%*%y).
```

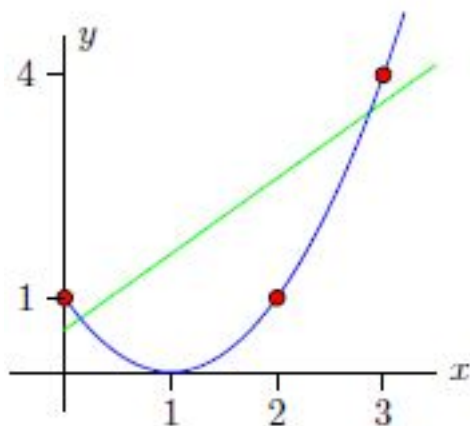
Sau codul R

```
x=c(0,2,3)
y=c(1,1,4)
x1=x
x2=x^ 2
lm(y~x1+x2)
```

Cu aceste date, parabola celor mai mici pătrate are ecuația

$$y = 1 - 2x + x^2.$$

Pentru 3 puncte, potrivirea pătratică este perfectă.



Potrivirea celor mai mici pătrate a unei drepte (verde) și a unei parabole (albastru)

**Exemplul 6.** Perechile  $(x_i, y_i)$  pot da vârsta și mărimea vocabularului a  $n$  copii. Deoarece copiii mici dobândesc cuvinte noi într-un ritm accelerat, putem ghici că un polinom de grad mai mare poate fi cea mai bună potrivire



pentru date.

**Exemplul 7.** (Transformarea datelor). Uneori este necesar să transformăm datele înainte de a folosi regresia liniară. De exemplu, presupunem că relația este exponențială, i.e.  $y = ce^{ax}$ . Atunci

$$\ln(y) = ax + \ln(c).$$

Deci putem folosi regresia liniară simplă pentru a obține un model

$$\ln(y_i) = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

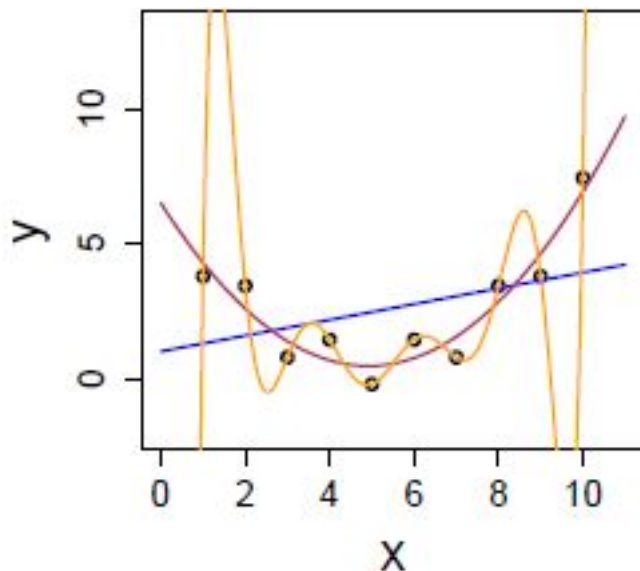
și apoi obținem modelul exponențial

$$y_i = e^{\hat{\beta}_0} e^{\hat{\beta}_1 x_i}.$$

#### 1.4.1 Suprapotrivirea

Putem totdeauna obține o potrivire mai bună folosind un polinom de ordin mai mare. De exemplu, date fiind 6 date bivariate (cu  $x_i$  distincte) se poate totdeauna afla un polinom de grad 5 care trece prin toate. Aceasta poate duce la [suprapotrivire](#). Adică, potrivirea zgomotului la fel de bine ca adevărata relație între  $x$  și  $y$ . Un model suprapotrivit va potrivi datele originale mai bine, dar va prezice mai puțin bine  $y$  pentru noi valori ale lui  $x$ . O povocare a modelării statistice este echilibrarea potrivirii modelului cu complexitatea modelului.

**Exemplul 8.** În reprezentarea de mai jos potrivim polinoame de gradul 1, 2 și 9 la 10 date bivariate. Modelul de gradul 2 (maro) dă o potrivire semnificativ mai bună decât modelul de gradul 1 (albastru). Modelul de gradul 9 (portocaliu) dă o potrivire exactă cu datele, dar dintr-o privire am ghici că este suprapotrivit. Adică, nu ne așteptăm că potrivească bine următoarea dată bivariată pe care o vedem. De fapt, datele au fost generate folosind un model pătratic, deci modelul de gradul 2 va tinde să facă cea mai bună potrivire cu date noi.



#### 1.4.2 Funcția R `lm`

Nu facem regresia liniară cu mâna. Regresia liniară se reduce la rezolvarea sistemelor de ecuații liniare, i.e. la calcul matriceal. Funcția R `lm` poate fi folosită la potrivirea unui polinom de orice grad cu datele. ([lm înseamnă model liniar](#)). De fapt, `lm` poate potrivi multe tipuri de funcții, exceptând polinoamele, după cum puteți explora folosind ajutorul lui R sau google.

### 1.5 Regresie liniară multiplă

Datele nu sunt totdeauna bivariante. Pot fi trivariate sau chiar de o dimensiune mai mare. Presupunem că avem datele de forma

$$(y_i, x_{1i}, x_{2i}, \dots, x_{mi}).$$

Putem analiza aceste date într-o manieră foarte similară cu datele bivariante. Adică, putem folosi cele mai mici pătrate pentru a potrivi modelul

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m.$$

Aici fiecare  $x_j$  este o variabilă predictor și  $y$  este variabila de răspuns. De exemplu, putem fi interesați de cum variază o populație de pești în funcție nivelele măsurate ale câtorva poluanți, sau am vrea să prezicem înălțimea de adult a unui fiu pe baza înălțimilor tatălui și a mamei.

## 1.6 Cele mai mici pătrate ca un model statistic

Modelul de regresie liniară pentru potrivirea unei drepte spune că valoarea  $y_i$  din perechea  $(x_i, y_i)$  este extrasă dintr-o variabilă aleatoare

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i,$$

unde termenii de "eroare"  $\varepsilon_i$  sunt variabile aleatoare independente cu media 0 și deviația standard  $\sigma$ . Presupunerea standard este că  $\varepsilon_i$  sunt i.i.d. cu repartiția  $N(0, \sigma^2)$ . În orice caz, media lui  $Y_i$  este dată de:

$$E(Y_i) = \beta_0 + \beta_1 x_i + E(\varepsilon_i) = \beta_0 + \beta_1 x_i.$$

Din această perspectivă, metoda celor mai mici pătrate alege valorile lui  $\beta_0$  și  $\beta_1$  care minimizează dispersia de selecție din jurul drepte.

De fapt, estimarea celor mai mici pătrate  $(\hat{\beta}_0, \hat{\beta}_1)$  coincide cu estimarea de verosimilitate maximă pentru parametrii  $(\beta_0, \beta_1)$ ; adică, dintre toți coeficienții posibili,  $(\hat{\beta}_0, \hat{\beta}_1)$  sunt cei care fac datele observate cele mai probabile.

## 1.7 Regresia la medie

Motivul termenului "regresie" este că variabila de răspuns prezisă  $y$  va tinde să fie "mai aproape" de (i.e. să regreseze la) media ei decât variabila predictor  $x$  este față de media ei. Aici "mai aproape" este în ghilimele deoarece trebuie să controlăm scala (i.e. deviația standard) fiecărei variabile. Modul de a controla scala este mai întâi să standardizăm fiecare variabilă.

$$u_i = \frac{x_i - \bar{x}}{\sqrt{s_{xx}}}, \quad v_i = \frac{y_i - \bar{y}}{\sqrt{s_{yy}}}.$$

Standardizarea schimbă media în 0 și dispersia în 1:

$$\bar{u} = \bar{v} = 0, \quad s_{uu} = s_{vv} = 1.$$

Proprietățile algebrice ale covarianței arată că

$$s_{uv} = \frac{s_{xy}}{\sqrt{s_{xx}s_{yy}}} = \rho,$$

coeficientul de corelație. Astfel, potrivirea celor mai mici pătrate pentru  $v = \beta_0 + \beta_1 u$  are

$$\hat{\beta}_1 = \frac{s_{uv}}{s_{uu}} = \rho \text{ și } \hat{\beta}_0 = \bar{v} - \hat{\beta}_1 \bar{u} = 0.$$

Deci dreapta celor mai mici pătrate este  $v = \rho u$ . Deoarece  $\rho$  este coeficientul de corelație, el este între -1 și 1. Presupunem că este pozitiv și mai mic ca 1 (i.e.,  $x$  și  $y$  sunt pozitiv, dar nu perfect corelate). Atunci formula  $v = \rho u$  înseamnă că, dacă  $u$  este pozitiv, atunci valoarea prezisă a lui  $v$  este mai mică decât  $u$ . Adică,  $v$  este mai aproape de 0 decât  $u$ . Echivalent,

$$\frac{y - \bar{y}}{\sqrt{s_{yy}}} < \frac{x - \bar{x}}{\sqrt{s_{xx}}},$$

i.e.,  $y$  regresează la  $\bar{y}$ . Standardizarea are grijă de scală.

Considerăm cazul extrem al corelației 0 între  $x$  și  $y$ . Atunci, indiferent de valoarea lui  $x$ , valoarea prezisă a lui  $y$  este totdeauna  $\bar{y}$ . Adică,  $y$  a regresat până la media lui.

Dreapta de regresie trece totdeauna prin punctul  $(\bar{x}, \bar{y})$ .

**Exemplul 9.** Regresia la medie este importantă în studiile longitudinale. Rice (*Mathematical Statistics and Data Analysis*) dă următorul exemplu. Presupunând că li se dau copiilor un test IQ la vârsta de 4 ani și altul la vârsta de 5 ani, ne așteptăm ca rezultatele să fie pozitiv corelate. Analiza de mai sus spune că, în medie, acei copii care au făcut slab la primul test tind să arate îmbunătățire (i.e. regresează la medie) în al 2-lea test. Astfel, o intervenție inutilă poate fi interpretată greșit ca utilă deoarece pare a îmbunătăți scorurile.

**Exemplul 10.** Alt exemplu cu consecințe practice este recompensa și pedeapsa. Imaginați-vă o școală unde performanța înaltă la un examen este recompensată și performanța slabă este pedepsită. Regresia la medie ne spune că (în medie) studenții foarte performanți vor face puțin mai slab la următorul examen și studenții puțin performanți vor face puțin mai bine. O viziune ne-sofisticată a datelor va face să pară că pedeapsa a îmbunătățit performanța și recompensa de fapt a scăzut performanța. Sunt reale consecințe dacă cei cu autoritate acționează după această idee.

## 1.8 Adaos

### 1.8.1 Demonstrația formulei pentru potrivirea celor mai mici pătrate a unei drepte

Cea mai directă demonstrație este cu analiza matematică. Suma erorilor pătrate este

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_1 x_i - \beta_0)^2.$$

Luând derivatele parțiale (și reamintind că  $x_i$  și  $y_i$  sunt date, deci constante)

$$\begin{aligned}\frac{\partial S}{\partial \beta_0} &= \sum_{i=1}^n -2(y_i - \beta_1 x_i - \beta_0) = 0 \\ \frac{\partial S}{\partial \beta_1} &= \sum_{i=1}^n -2x_i(y_i - \beta_1 x_i - \beta_0) = 0.\end{aligned}$$

Se obține următorul sistem de 2 ecuații liniare în necunoscutele  $\beta_0$  și  $\beta_1$ :

$$\begin{aligned}\left(\sum_{i=1}^n x_i\right) \beta_1 + n\beta_0 &= \sum_{i=1}^n y_i \\ \left(\sum_{i=1}^n x_i^2\right) \beta_1 + \left(\sum_{i=1}^n x_i\right) \beta_0 &= \sum_{i=1}^n x_i y_i\end{aligned}$$

Rezolvând sistemul obținem formulele (1).

Pentru multe aplicații între discipline vezi:

[http://en.wikipedia.org/wiki/Linear\\_regression#Applications\\_of\\_linear\\_regression](http://en.wikipedia.org/wiki/Linear_regression#Applications_of_linear_regression).

### 1.8.2 Măsurarea potrivirii

Odată ce se calculează coeficienții de regresie, este important să verificăm cât de bine modelul de regresie se potrivește cu datele (i.e., cât de aproape cea mai potrivită dreaptă urmărește datele). O măsură uzuală dar brută a ”bunătății de potrivire” este **coeficientul de determinare**, notat  $R^2$ . **Suma totală a pătratelor** este dată de:

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2.$$

**Suma reziduală a pătratelor** este dată de suma pătratelor reziduurilor. Când potrivim o dreaptă, aceasta este:

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2.$$

RSS este porțiunea ”neexplicată” a sumei totale a pătratelor, i.e. neexplicată de ecuația de regresie. Diferența TSS–RSS este porțiunea ”explicată” a sumei totale a pătratelor. **Coeficientul de determinare**  $R^2$  este raportul dintre porțiunea ”explicată” și suma totală a pătratelor:

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}}.$$

Cu alte cuvinte,  $R^2$  măsoară proporția variabilității datelor care este contabilizată pentru modelul de regresie. O valoare aproape de 1 indică o potrivire bună, în timp ce o valoare aproape de 0 indică o potrivire slabă. În cazul regresiei liniare simple,  $R^2$  este pur și simplu pătratul coeficientului de corelație dintre valorile observate  $y_i$  și valorile prezise  $\beta_0 + \beta_1 x_i$ .

**Exemplul 11.** În exemplul 8 de suprapotrivire, valorile lui  $R^2$  sunt ("degree"="grad"):

degree	$R^2$
1	0.3968
2	0.9455
9	1.0000

Măsura bunătații potrivirii crește când  $n$  (gradul) crește. Potrivirea este mai bună, dar modelul devine de asemenea mai complex, deoarece este nevoie de mai mulți coeficienți pentru a descrie polinoame de grad mai mare.