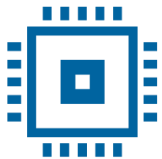


Analiza și procesarea datelor prin tehnici de Învățare Automată

4. Tehnici de învățare supervizată 2: Regresie



Universitatea
Transilvania
din Brașov

FACULTATEA DE INGINERIE ELECTRICĂ
ȘI ȘTIINȚA CALCULATOARELOR

Șef Lucrări Dr. Ing. Horia Modran

Contact: horia.modran@unitbv.ro / modranhoria@gmail.com

Tel: 0770171577

2024 - 2025



Clasificare şi Regresie

- Învăţarea unei funcţii discrete: **Clasificare**
 - Clasificare binară:
 - fiecare exemplu este clasificat ca adevărat (pozitiv) sau fals (negativ)
 - poate, de asemenea, clasifica în mai multe clase (3, 4, 5...)
- Învăţarea unei funcţii continue: **Regresie**
 - Regresie liniară
 - Regresie logistică



Regresie

- Tehnică de învăţare supervizată
- Urmăreşte să determine corelaţia între variabile
- În contextul învăţării automate: prezicerea unei variabile ţintă continue în funcţie de una sau mai multe variabile
- Utilizată în
 - predicţie/prognoză
 - modelarea seriilor de timp
 - determinarea relaţiei cauză-efect



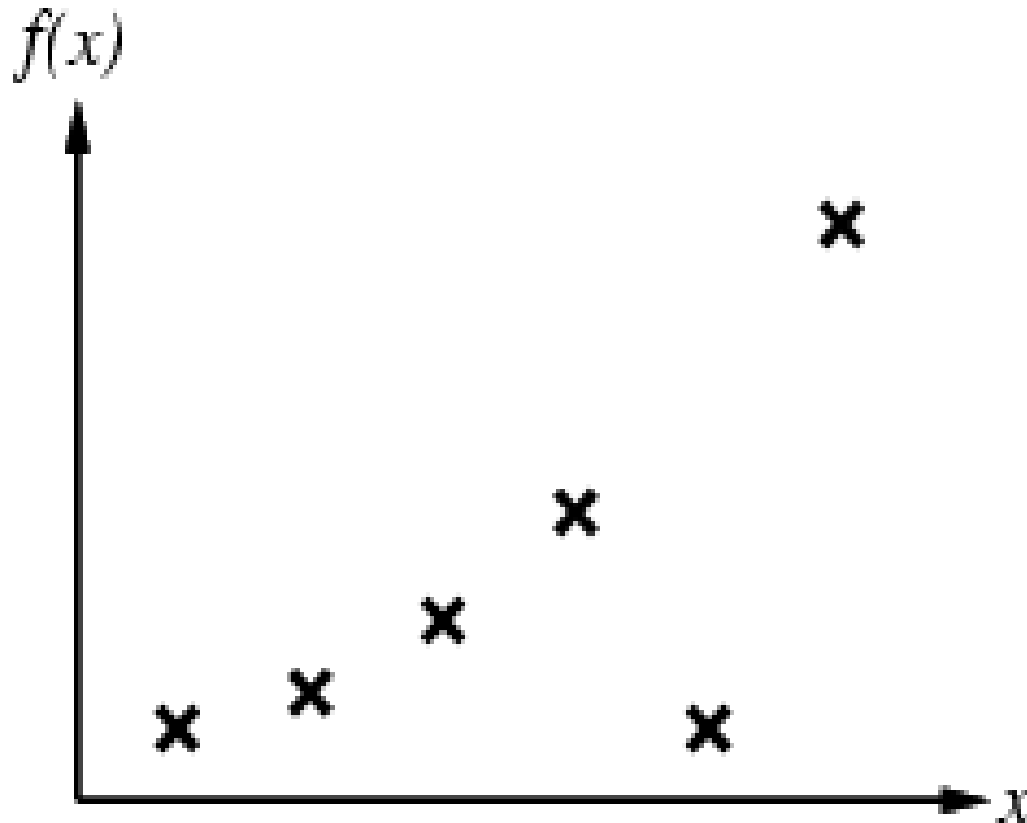
Corelaţie

- Asocierie liniară între două variabile
- Arată cum să determinăm atât natura, cât şi puterea relaţiei dintre două variabile
- Corelaţia este între -1 şi $+1$
- Corelaţia zero indică faptul că nu există nicio relaţie între variabile
- Coeficientul de corelaţie Pearson
 - cea mai cunoscută măsură a dependenţei dintre două mărimi



Aproximarea unei funcţii

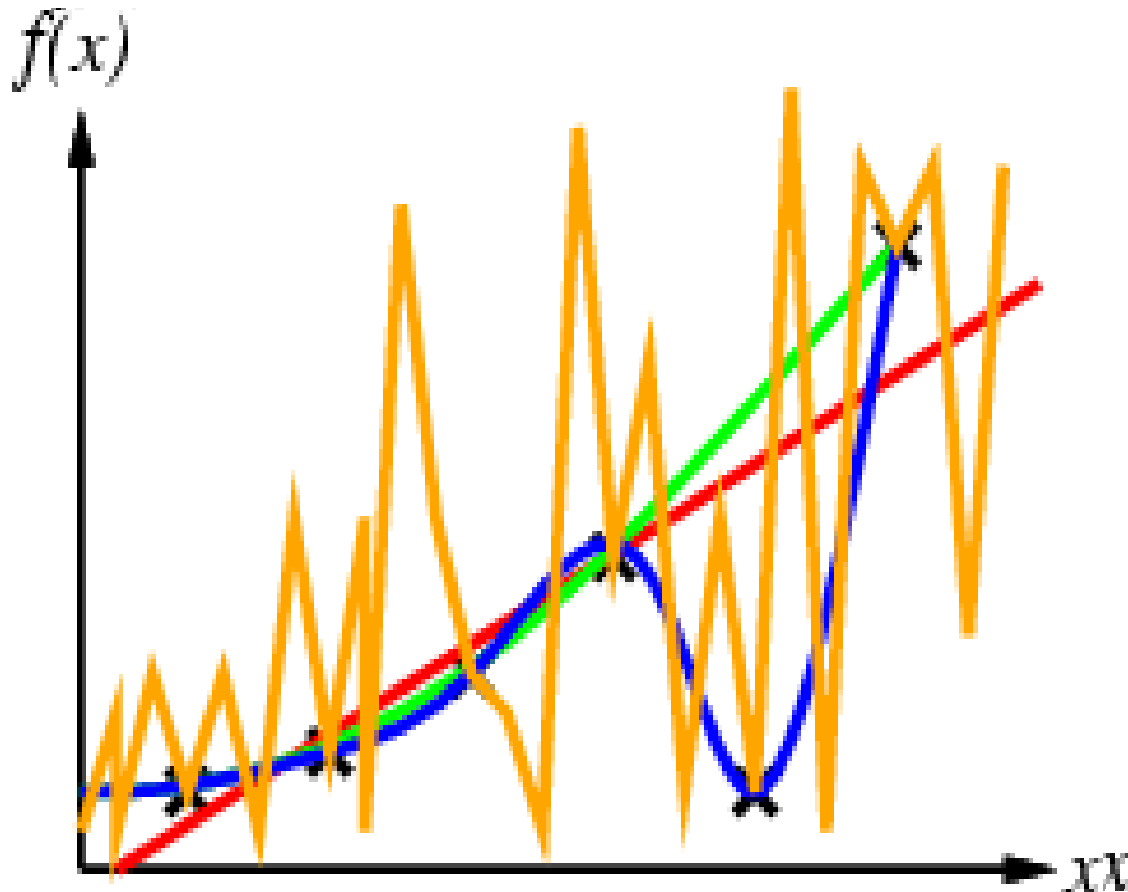
- Să se determine funcţia care trece prin punctele de mai jos





Aproximarea unei funcţii

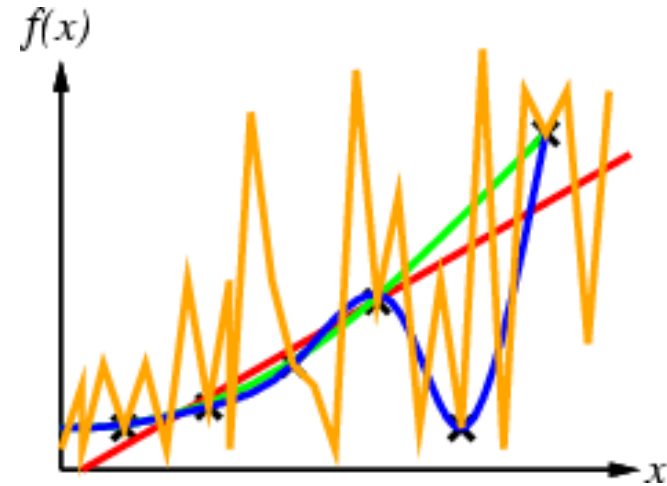
▣ Modalităţi de abordare





Briciul lui Ocam

- Trebuie preferată cea mai simplă ipoteză consistentă cu datele
- În caz de egalitate, modelele mai simple tind să generalizeze mai bine decât cele complexe



William of Occam (1285-1347) – filozof englez
Entia non sunt multiplicanda praeter necessitatem
~ Cea mai simplă explicație este cea mai bună



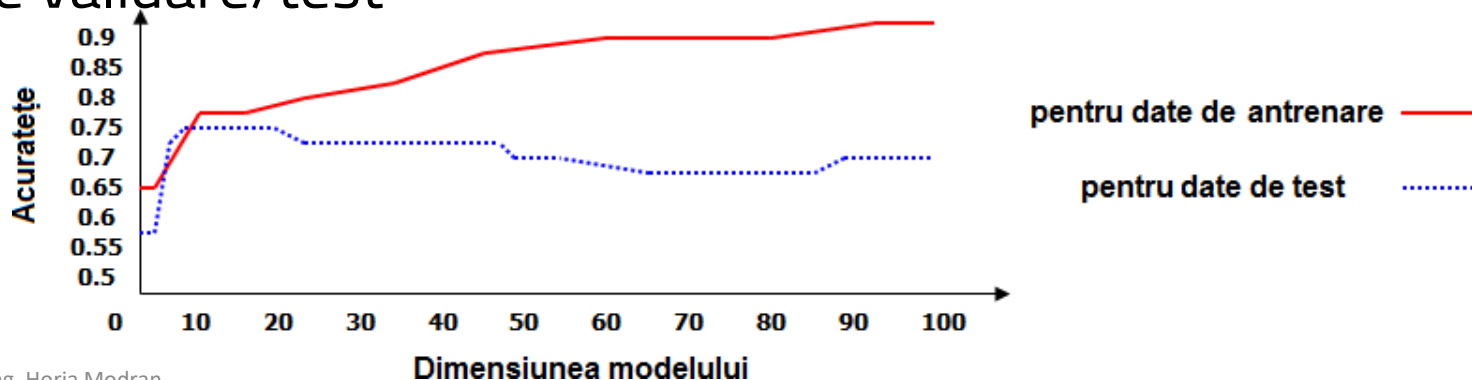
Generalizarea

- Scopul este găsirea unui model pentru determinarea valorii clasei în funcţie de valorile celorlalte attribute cu eroare cât mai mică
- Modelul trebuie să aibă **capacitate de generalizare**, care arată cât de bun este modelul pentru date **noi**
- De obicei există o **mulţime de antrenare** pentru crearea modelului şi o **mulţime de test** pentru verificarea capacităţii de generalizare



Generalitatea unui model

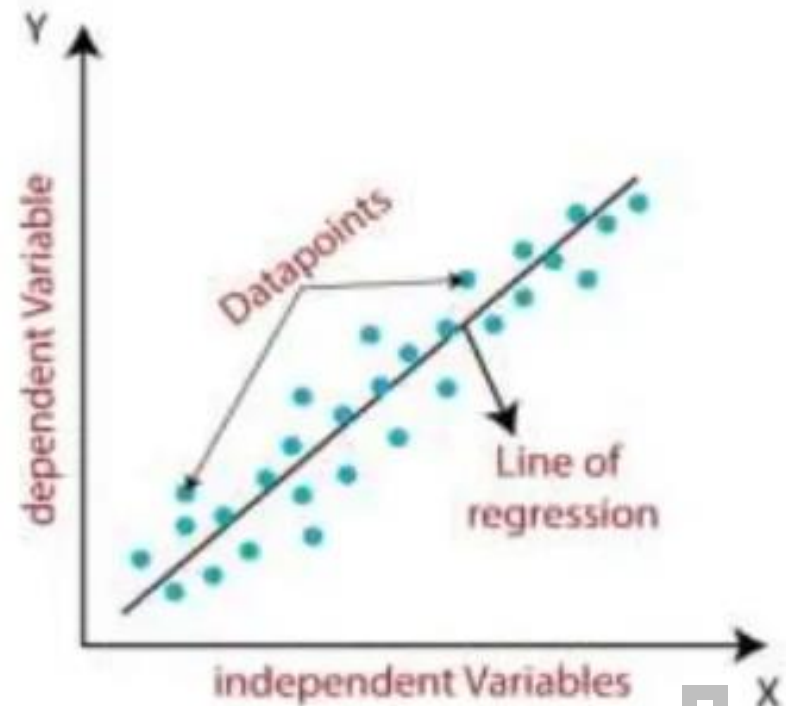
- **Subpotrivirea** (*underfitting*): modelul este prea simplu şi nu poate învăţa distribuţia datelor
- **Suprapotrivirea** (*overfitting*): modelul este prea complex şi poate fi influenţat de zgomot şi date irelevante
- Un model suprapotrivit are performanţe foarte bune pe mulţimea de antrenare, dar performanţe slabe pe mulţimea de validare/test





Regresie Liniară

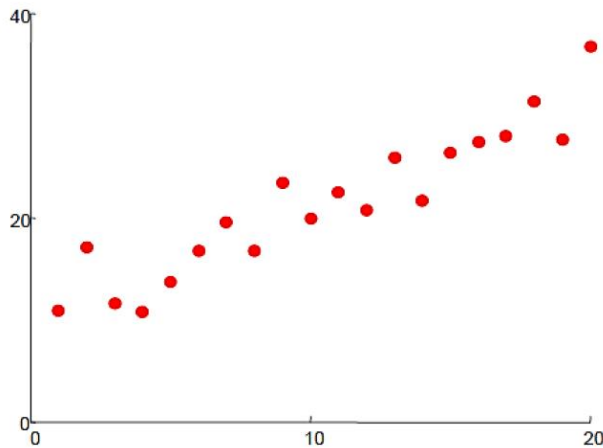
- **Regresia Liniară** este utilizată pentru a prezice valoarea unei variabile pe baza valorilor uneia/mai multor variabile
- Este o metodă de învăţare supervizată
- Scopul: determinarea relaţiei liniare dintre o variabilă dependentă şi una/mai multe variabile independente



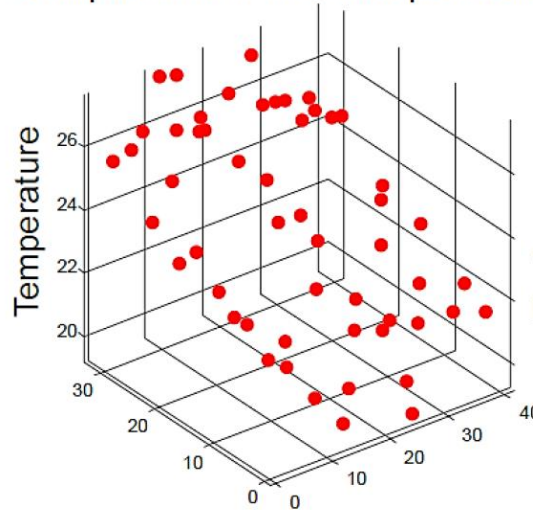


Regresie Liniară

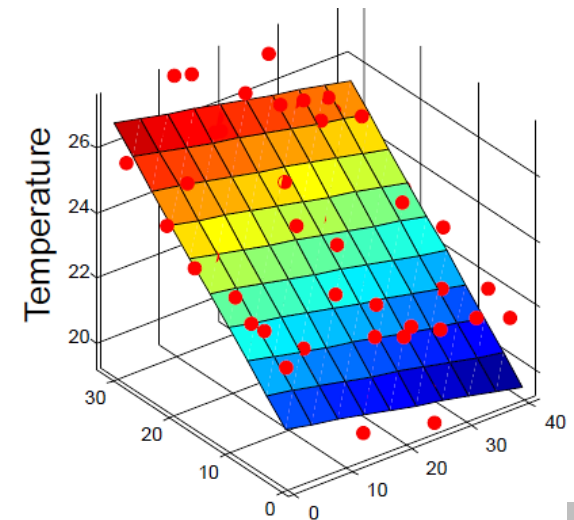
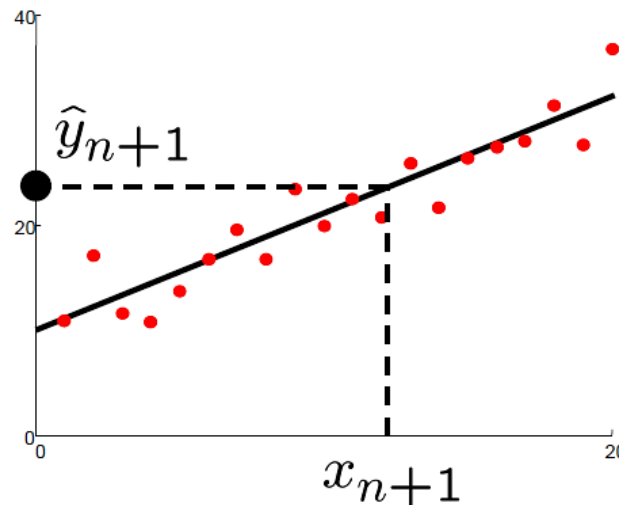
Samples with ONE independent variable



Samples with TWO independent variables



Given examples $(x_i, y_i)_{i=1 \dots n}$
Predict y_{n+1} given a new point x_{n+1}





Regresia Liniară

- Considerăm un model

$$y = a + bX + \varepsilon$$

, unde **a** şi **b** sunt interceptarea şi panta (cunoscuţi drept coeficienţi sau parametri), iar ε este termenul de eroare

- Regresie liniară simplă

- o singură variabilă independentă este utilizată

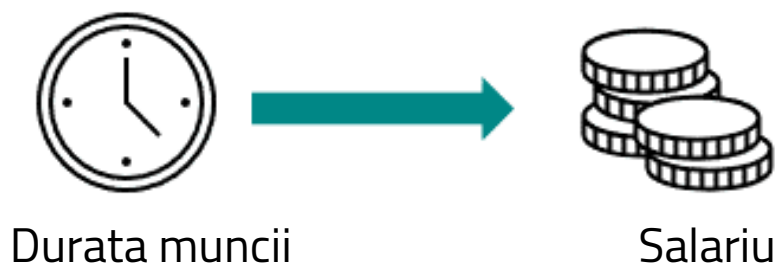
- Regresie liniară multiplă

- două sau mai multe variabile independente sunt folosite pentru predicţie

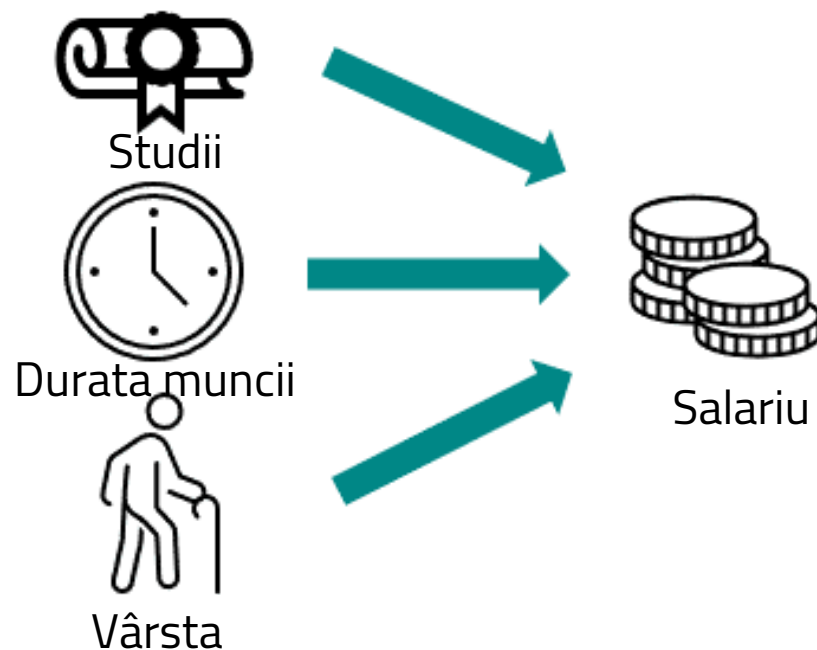


Tipuri de regresie liniară

Simple Linear Regression



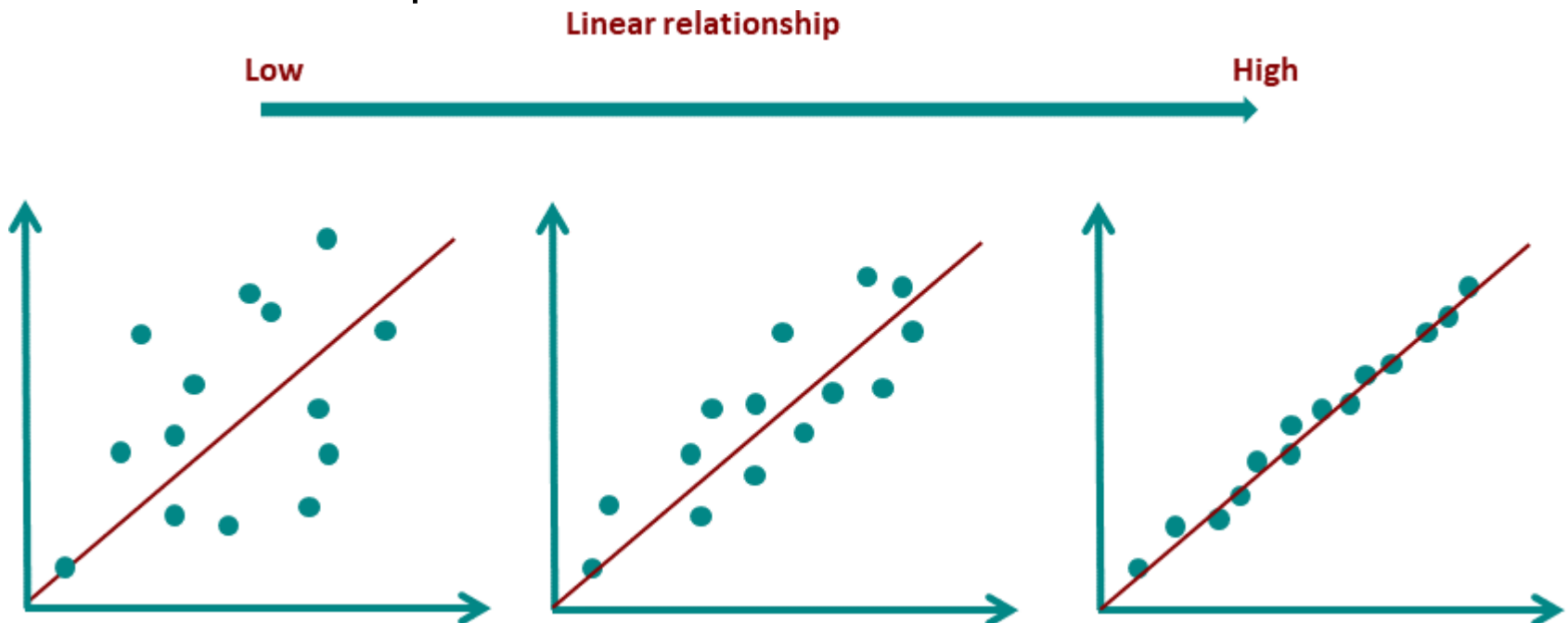
Multiple Linear Regression





Regresia Liniară simplă

- Exemplu: Are **înălţimea** influenţă asupra **greutăţii** unei persoane?
Variabilă independentă
Variabilă dependentă
- Scop: prezicerea valorii variabilei dependente (y) pe baza variabilei independente (X)





Regresia Liniară simplă

$$y_i = a + bX_i + \varepsilon_i$$

, unde y = variabila dependentă, x – variabila independentă,

a – intercept, b – panta liniei, ε – termenul de eroare

- **Model Simplu:** doar un X
- **Liniar în parametri:** niciun parametru nu apare ca exponent sau este înmulțit/împărțit cu un alt parametru
- **Liniar în variabila predictor (X):** X apare doar la prima putere



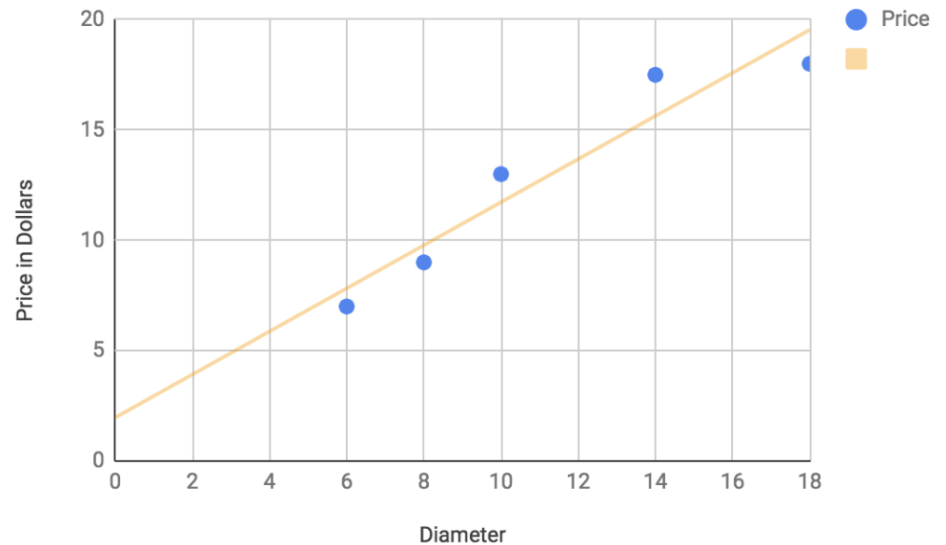
Exemplu

- Predicţia preţului pentru pizza
- Care este preţul unei pizza de 20 inch?



Diameter	Price
6	7
8	9
10	13
14	17.5
18	18

Pizza Price by Diameter





Calculare

x	y	x^2	xy
3	8	9	24
9	6	81	54
5	4	25	20
3	2	9	6
$\Sigma x = 20$	$\Sigma y = 20$	$\Sigma x^2 = 124$	$\Sigma xy = 104$

x	y
3	8
9	6
5	4
3	2

$$a = \frac{\Sigma y \Sigma x^2 - \Sigma x \Sigma xy}{n \Sigma x^2 - (\Sigma x)^2} = \frac{20 * 124 - 20 * 104}{4 * 124 - 20^2} = \frac{400}{96} = 4.17$$

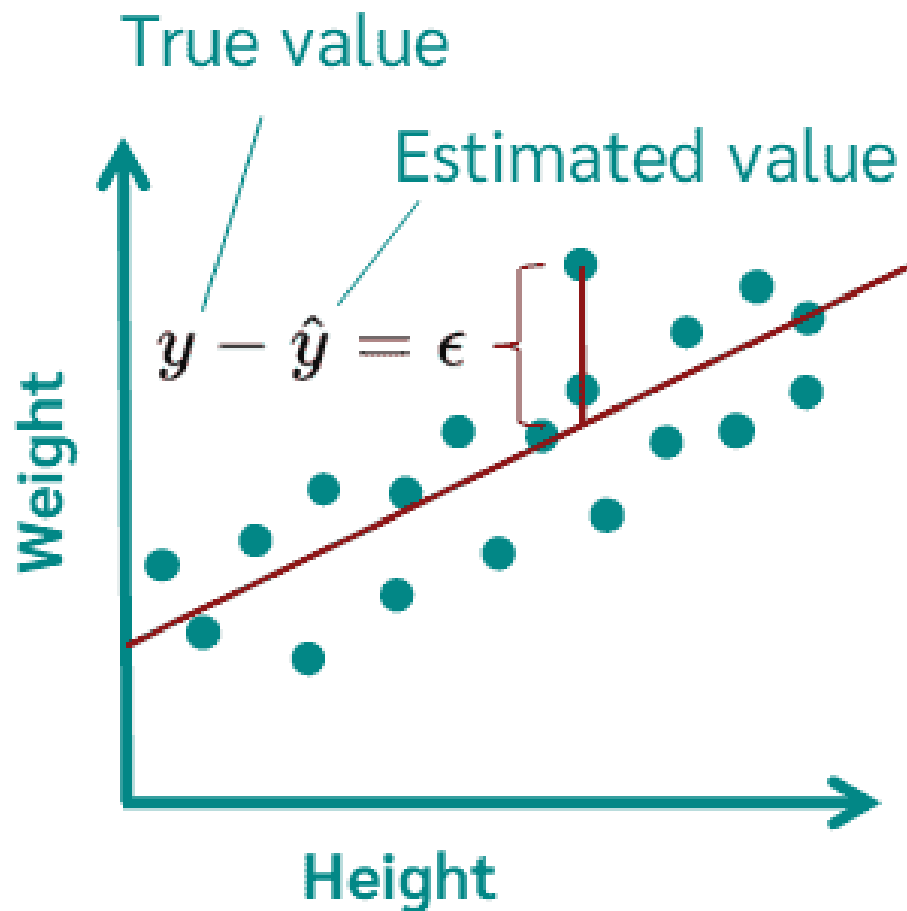
$$b = \frac{n \Sigma xy - \Sigma x \Sigma y}{n \Sigma x^2 - (\Sigma x)^2} = \frac{4 * 104 - 20 * 20}{4 * 124 - 20^2} = \frac{16}{96} = 0.166$$

$$y = 4.17 + 0.166x$$





Eroare



Error epsilon

$$y = b \cdot x + a + \epsilon$$





Eroare

- Graficul *Scatterplot* arată că punctele nu sunt pe o linie şi astfel, pe lângă relaţie, descriem şi eroarea ε :

$$y_i = a + b_1 X_i + \varepsilon_i$$

- distanţa dintre valoarea estimată şi valoarea adevărată trebuie să fie cât mai mică (implicit şi eroarea ε)
- Coeficientul de regresie b poate avea acum semne diferite:
 - $b > 0$: există o corelaţie pozitivă între x şi y (cu cât x mai mare, cu atât y mai mare)
 - $b < 0$: există o corelaţie negativă între x şi y
 - $b = 0$: nu există nicio corelaţie între x şi y



Regresia Liniară multiplă

$$y = a + b_1 X_1 + b_2 X_2 + \cdots + b_n X_n + \varepsilon$$

, unde y = variabila dependentă, x – variabila independentă,

a – intercept, b_1, b_2, \dots, b_n – pantele, ε – termenul de eroare

- Coeficienţii pot fi interpretaţi ca la regresia simplă
- Mai multe variabile independente
- O singură variabilă dependentă
- Utilizată adesea în cercetări de piaţă



Regresia Liniară multiplă

■ Predicție Y după x_1 și x_2

$$X = \begin{pmatrix} 1 & 1 & 4 \\ 1 & 2 & 5 \\ 1 & 3 & 8 \\ 1 & 4 & 2 \end{pmatrix} \quad y = \begin{pmatrix} 1 \\ 6 \\ 8 \\ 12 \end{pmatrix}$$

x1 Product 1 Sales	x2 Product 2 Sales	Y Weekly Sales
1	4	1
2	5	6
3	8	8
4	2	12

$$y = a + b_1x_1 + b_2x_2$$

$$b = ((X^T X)^{-1} X^T) Y, \text{ unde } a = \begin{pmatrix} a \\ b_1 \\ b_2 \end{pmatrix}$$



Exemplu calcul

$$X^T X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 & 4 \\ 1 & 2 & 5 \\ 1 & 3 & 8 \\ 1 & 4 & 2 \end{pmatrix} = \begin{pmatrix} 4 & 10 & 19 \\ 10 & 30 & 46 \\ 19 & 46 & 109 \end{pmatrix}$$

$$(X^T X)^{-1} = \begin{pmatrix} 4 & 10 & 19 \\ 10 & 30 & 46 \\ 19 & 46 & 109 \end{pmatrix}^{-1} = \begin{pmatrix} 3.15 & -0.59 & -0.30 \\ -0.59 & 0.2 & 0.016 \\ -0.3 & 0.016 & 0.054 \end{pmatrix}$$

$$(X^T X)^{-1} X^T = \begin{pmatrix} 3.15 & -0.59 & -0.30 \\ -0.59 & 0.2 & 0.016 \\ -0.3 & 0.016 & 0.054 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{pmatrix} = \begin{pmatrix} 0.05 & 0.47 & -1.02 & 0.19 \\ -0.32 & -0.098 & 0.155 & 0.26 \\ -0.065 & 0.05 & 0.185 & -0.125 \end{pmatrix}$$

$$= ((X^T X)^{-1} X^T) Y = \begin{pmatrix} 0.05 & 0.47 & -1.02 & 0.19 \\ -0.32 & -0.098 & 0.155 & 0.26 \\ -0.065 & 0.05 & 0.185 & -0.125 \end{pmatrix} \begin{pmatrix} 1 \\ 6 \\ 8 \\ 12 \end{pmatrix} = \begin{pmatrix} -1.69 \\ 3.48 \\ -0.05 \end{pmatrix}$$

$$\mathbf{y} = -1.69 + 3.48x_1 - 0.05x_2$$



Regresie cu variabile categoriale

- Presupunem că toate variabilele sunt variabile continue
- Variabile categoriale:
 - Variabile ordinale - codifică datele cu valori continue
 - Evaluation: Excellent (5), Very good (4), Good (3), Poor (2), Very poor (1)
 - Variabile nominale – de utilizează *dummy variables*
 - Department: Computer, Biology, Physics

	Computer	Biology	Physics
Computer	1	0	0
Biology	0	1	0
Physics	0	0	1



Funcţie cost

- Funcţia de cost este necesară pentru a calcula diferenţa dintre valorile reale şi cele prezise
- Pentru modelul de regresie liniară, funcţia de cost va fi minimul Erorii Pătratică Medie Rădăcină (engl. *Root Mean Square Error* – **RMSE**) a modelului, obţinută prin scăderea valorilor prezise din valorile reale
- Este o valoare numerică care indică succesul sau eşecul unui anumit model, fără a fi nevoie să înţelegem funcţionarea interioară a unui model



Tipuri de funcţie cost

- **Eroare Medie (*Mean Error*)** – aceste erori pot fi negative sau pozitive
 - prin urmare, se pot anula reciproc în timpul însumării, iar **eroarea medie a modelului poate fi zero**
- **Eroare medie pătratică (*Mean Squared Error* - MSE)** - MSE reprezintă diferenţa medie pătrată dintre predicţii şi rezultatele aşteptate

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

n – nr. de puncte
 y_i - valoarea reală
 \hat{y}_i - valoarea prezisă



Tipuri de funcţie cost

- Eroare absolută medie (*Mean Absolute Error* – MAE) – MAE calculează valoarea absolută a diferenţei (nu există posibilitate de erori negative)

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

n – nr. de puncte
 y_i – valoarea reală
 \hat{y}_i – valoarea prezisă

- Root Mean Squared Error – unul dintre cei doi indicatori principali de performanţă pentru un model de regresie

$$RMSE = \sqrt{\frac{\sum_{i=1}^n |\hat{y}_i - y_i|^2}{n}}$$



Coeficient de determinare

- Coeficientul de determinare (R^2) este proporția variației variabilei dependente care este datorată variabilelor independente

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

$$SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SS_{tot} = \sum_{i=1}^n (y_i - \bar{y}_i)^2$$

y_i - valoarea reală

\hat{y}_i - valoarea estimate de model

\bar{y}_i - valoarea medie

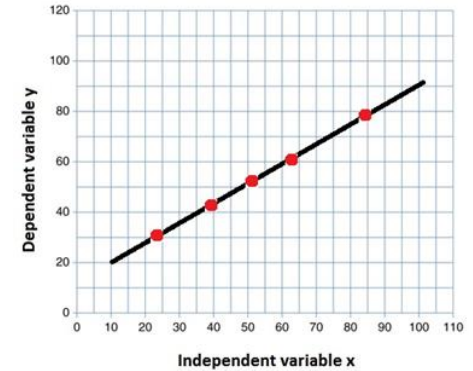
$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y}_i)^2}$$



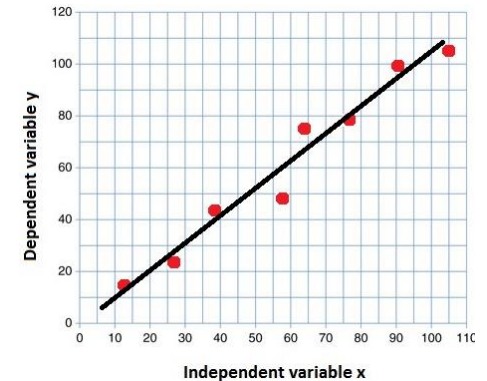
Coeficient de determinare

■ Exemplu de valori pentru R^2

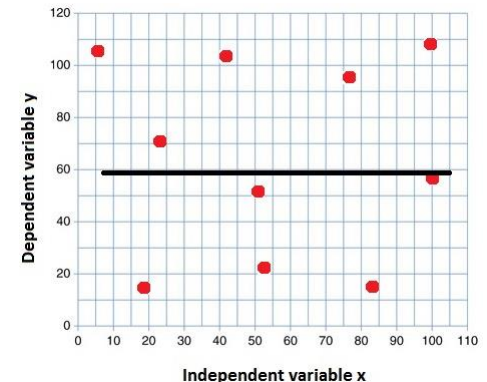
$R^2 = 1$ Toată variaţia valorilor y este explicată de valorile x



$R^2 = 0.83$ 83 % din variaţia valorilor y este explicată de valorile x



$R^2 = 0$ Variaţia valorilor y nu este influenţată deloc de valorile x





Regresie Liniară pentru clasificare

- Pentru clasificare binară
 - se codifică etichetele clasei ca $y=0,1$ sau $\{-1,1\}$
 - se aplică: $\mathbf{y} = \mathbf{X} * \mathbf{b} + \mathbf{e}$
 - se determină de care clasă este mai apropiată predicția
 - dacă clasa 1 este codificată la 1 și clasa 2 este - 1
- class 1 if $f(x) \geq 0$*
class 2 if $f(x) < 0$
- Modelel liniare NU sunt optimizate pentru clasificare



Regresie logistică



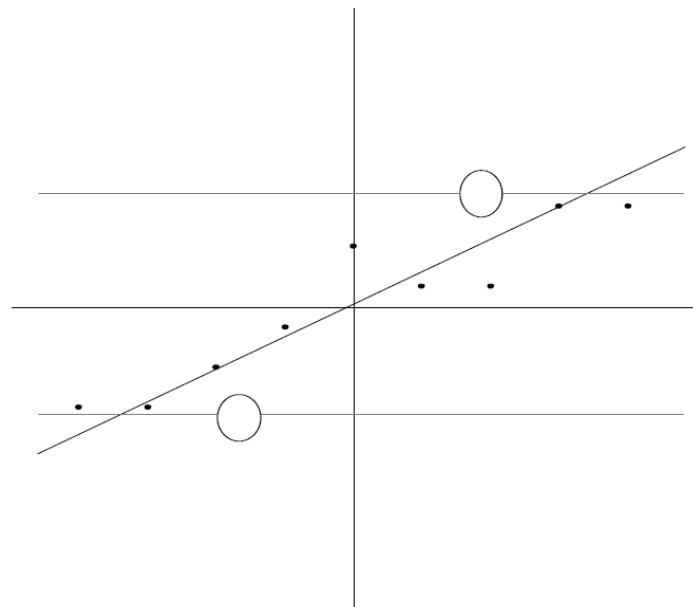
Regresie Logistică

- Prezice rezultatele pe o variabilă de rezultat binară
 - de exemplu, dacă un pacient are sau nu o boală
 - dacă un nou solicitant va reuşi sau nu
 - rezultatul nu este continuu sau distribuit normal
- când avem un răspuns de tip variabilă binară
 - codificăm „boală” ca 1 şi „fără boală” ca 0, putem doar să potrivim o linie prin acele puncte aşa cum am face cu regresia liniară? **Posibil! Dar există anumite probleme.**



Regresie Liniară: Probleme

- Problema potrivirii unei linii regulate de regresie la o variabilă dependentă binară
 - linia pare să simplifice prea mult relația
 - oferă predicții care nu pot fi valori observabile ale lui Y pentru valorile extreme ale lui X
 - abordarea este analogă cu potrivirea unui model liniar la probabilitatea evenimentului
- Produce predicții neobservabile pentru valori extreme ale variabilei dependente





Regresie Logistică

- Se începe cu ecuaţia Regresiei Liniare
- Se vor determina coeficienţii de regresie

Variabilă dependentă

Variabile independente

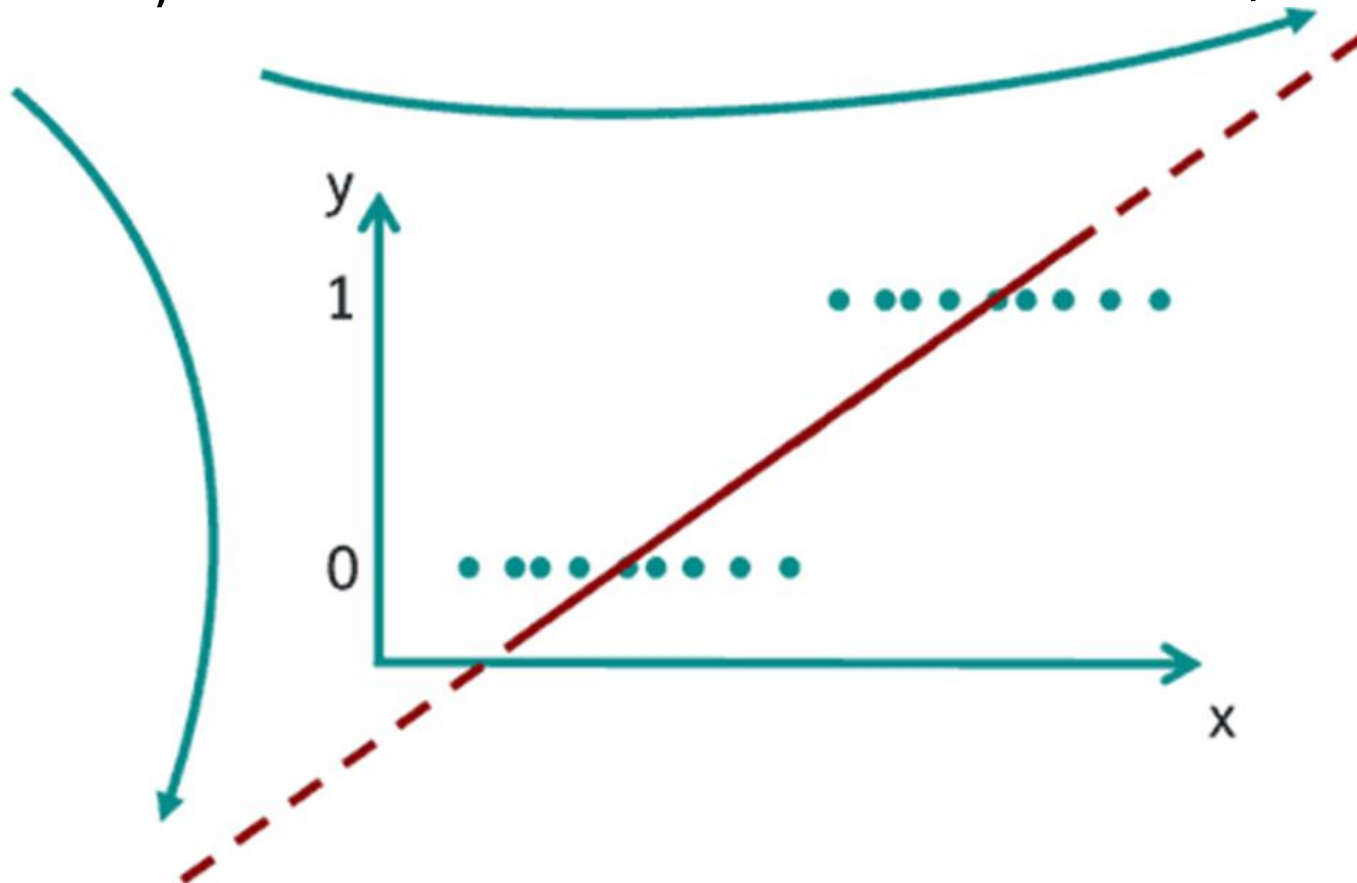
$$\hat{y} = b_1 \cdot x_1 + b_2 \cdot x_2 + \dots + b_k \cdot x_k + a$$

Coeficienţi de regresie



Regresie Logistică

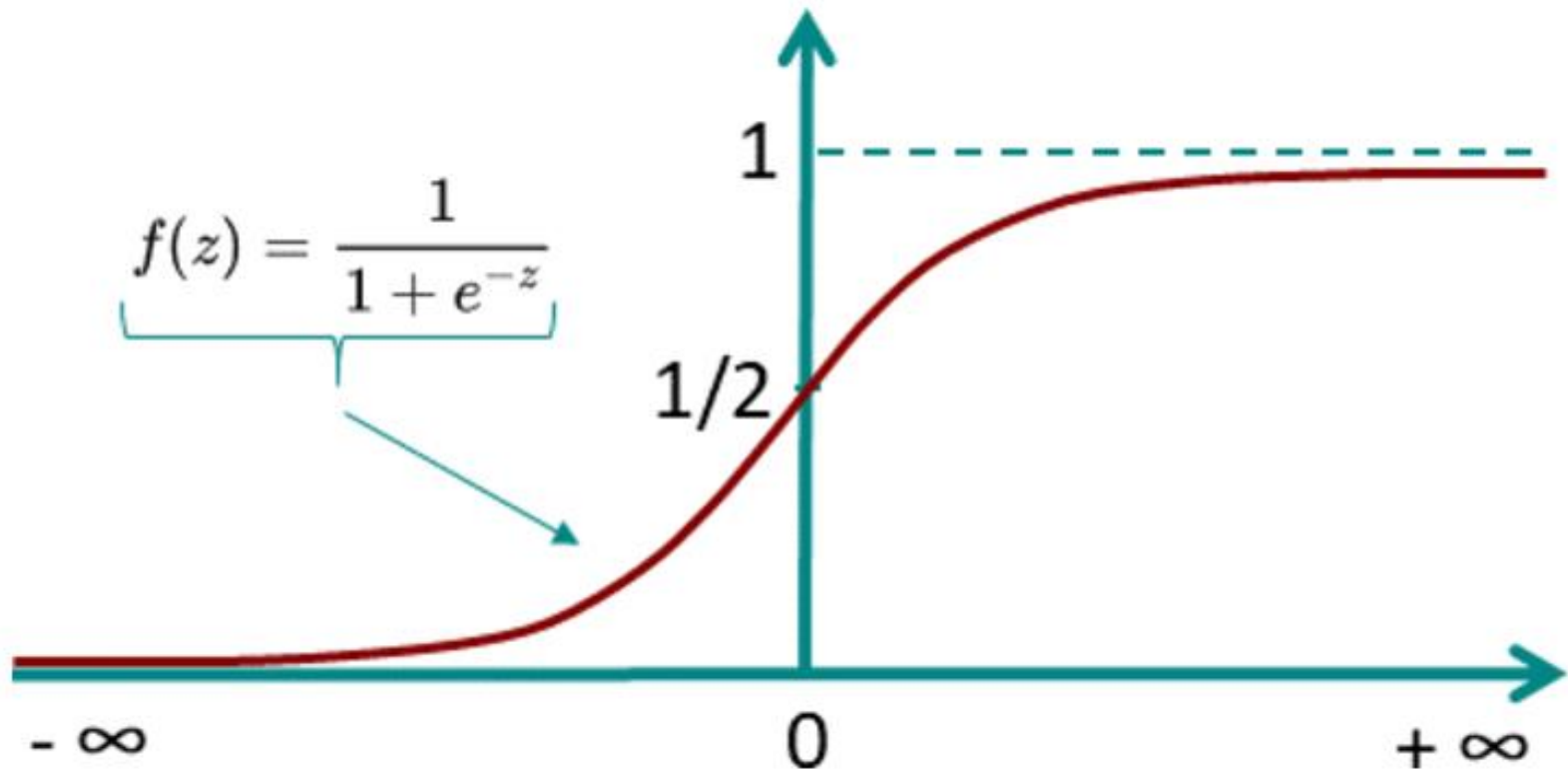
- Rezultatul calcului ar fi un model liniar (similar cu Regresia Linieară) \rightarrow intervalul de valori este de la $-\infty$ la $+\infty$





Funcția logistică

- Modelul logistic se bazează pe funcția logistică
- Are valori doar între 0 și 1





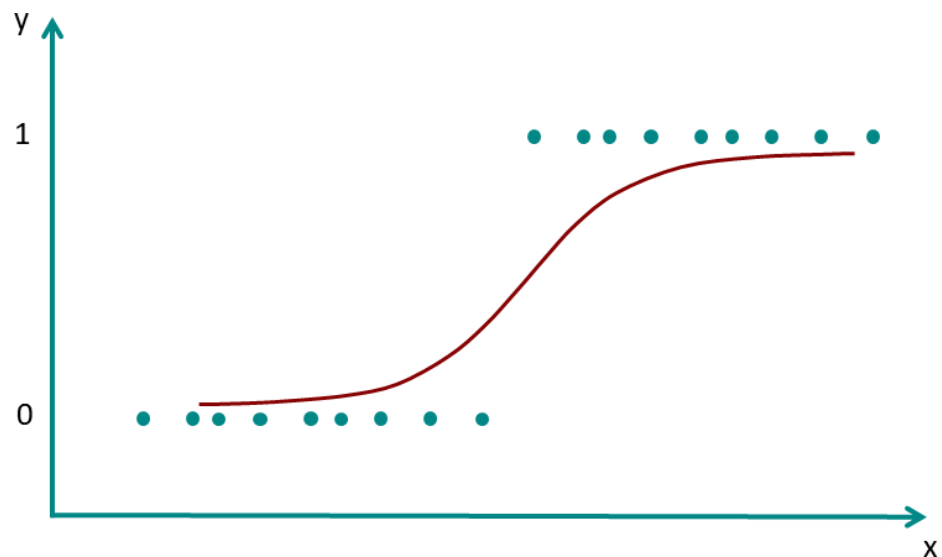
Funcția logistică

■ Ecuația funcției logistice

$$\hat{y} = b_1 \cdot x_1 + b_2 \cdot x_2 + \dots + b_k \cdot x_k + a$$

$$f(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-(b_1 \cdot x_1 + \dots + b_k \cdot x_k + a)}}$$

■ Graficul funcției va arăta astfel:





Abordare probabilistică

- Funcţia logistică este perfectă pentru a descrie probabilitatea $P(y=1|x)$
- Probabilitatea ca pentru valori date ale variabilei independente, variabila dependentă binară y să fie 0 sau 1 este dată de ecuaţia:

$$P(y = 1|x_1, x_2, \dots, x_n) = \frac{1}{1 + e^{-(b_1x_1 + b_2x_2 + \dots + b_nx_n + a)}}$$

$$\begin{aligned} P(y = 0|x_1, x_2, \dots, x_n) &= 1 - P(y = 1) \\ &= 1 - \frac{1}{1 + e^{-(b_1x_1 + b_2x_2 + \dots + b_nx_n + a)}} \end{aligned}$$



Exemplu

- Exemplu - tabelul arată numărul de ore de studiu pentru un student şi dacă a promovat examenul (1) sau nu (0)

Study Hours (x_k)	0.5	0.75	1.00	1.25	1.50	1.75	2.00	2.25	2.50	2.75	3.00	3.25	3.50	4.00	4.25	4.50	4.75	5.00	5.50
Pass (y_k)	0	0	0	0	0	0	0	1	0	1	0	1	0	1	1	1	1	1	1

Ecuatie de regresie: $\hat{y}_i = a + b_1 X_i$ Determinăm valorile pentru:
 $a = -4.1$
 $b_1 = 1.5$

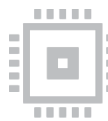
De exemplu, pentru $x = 2$: $\hat{y}_i = a + b_1 * 2 = -4.1 + 1.5 * 2 = -1.1$

Aplicăm funcția logistică: $P(y = 1|x) = \frac{1}{1 + e^{-\hat{y}}} = \frac{1}{1 + e^{1.1}} \approx 0.25 \rightarrow \text{clasa 0}$

Similar pentru $x = 4$: $\hat{y}_i = a + b_1 * 4 = -4.1 + 1.5 * 4 = 1.9$

$$P(y = 1|x) = \frac{1}{1 + e^{-\hat{y}}} = \frac{1}{1 + e^{-1.9}} \approx 0.87 \rightarrow \text{clasa 1}$$





Regresie Logistică Multinomială

- De regulă regresia logistică este aplicată în probleme de clasificare binară
- Regresia logistică multinomială este o metodă de clasificare care generalizează regresia logistică la problema de clasificare multclasă
- Modul de funcționare este același ca și în regresia logistică, singura diferență fiind că variabilele dependente sunt mai degrabă categorice decât binare, adică există n rezultate posibile și nu doar două





Indicatori de performanţă

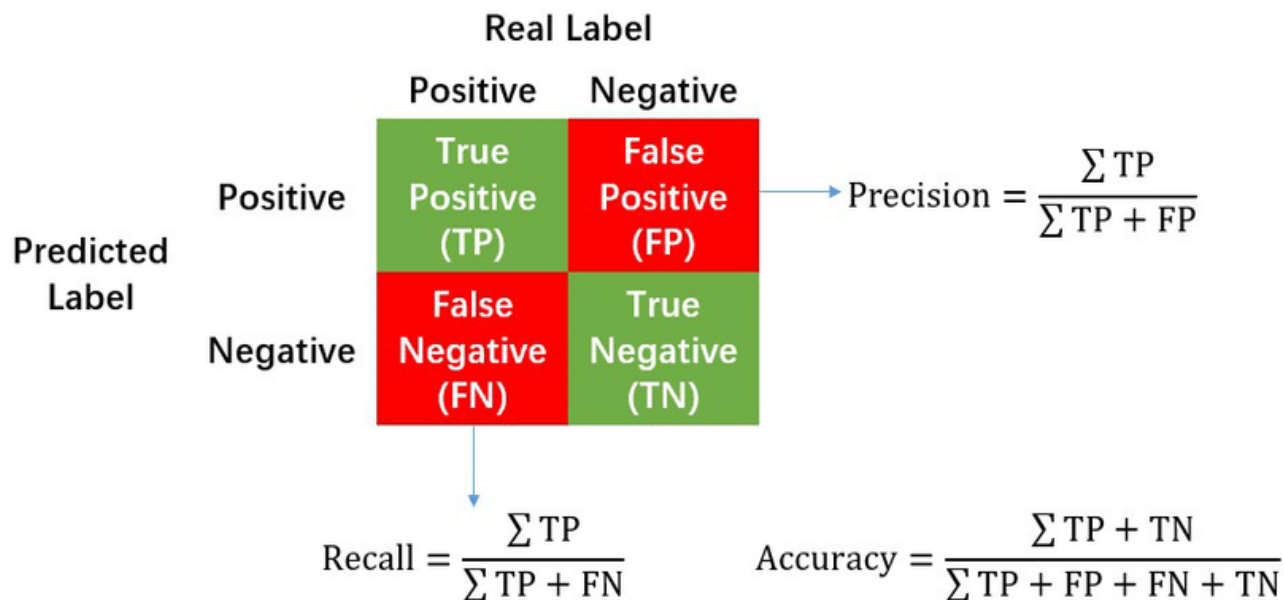
- Fiind algoritm de clasificare, se aplică aceiaşi indicatori de performanţă specifici algoritmilor de clasificare:

- Acurateţe

- Precizie

- Recall

- Scor F1



$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall}$$



ÎNTREBĂRI ?

