



**UNIVERSITATEA "AUREL VLAICU" DIN ARAD**  
**FACULTATEA DE ȘTIINȚE EXACTE**  
**DOMENIUL: INFORMATICĂ**  
**PROGRAMUL DE STUDIU: INFORMATICĂ**  
**FORMA DE ÎNVĂȚĂMÂNT CU FRECVENȚĂ**

**LUCRARE DE LICENȚĂ**

**ÎNDRUMĂTOR ȘTIINȚIFIC:**

**Lector univ. dr. – VLAD F. DRĂGOI**

**ABSOLVENT:**

**DREGHICI GH. NICOLETA**

**ARAD**

**Iunie 2021**



**UNIVERSITATEA "AUREL VLAICU" DIN ARAD**

**FACULTATEA DE ȘTIINȚE EXACTE**

**DOMENIUL: INFORMATICĂ**

**PROGRAMUL DE STUDIU: INFORMATICĂ**

**FORMA DE ÎNVĂȚĂMÂNT CU FRECVENȚĂ**

**MODELE DE REGRESIE ȘI RANDOM FOREST  
PENTRU DATA SCIENCE**

**ÎNDRUMĂTOR ȘTIINȚIFIC:**

**Lector univ. dr. – VLAD F. DRĂGOI**

**ABSOLVENT:**

**DREGHICI GH. NICOLETA**

**ARAD**

**Iunie 2021**

# Cuprins

<b>Introducere .....</b>	<b>1</b>
<b>1. Econometrie – studiu de caz.....</b>	<b>5</b>
<b>2. Etapele și conținutul unui proiect Data Science.....</b>	<b>14</b>
2.1. Planificarea proiectului.....	15
2.2. Obținerea/colectarea datelor .....	15
2.3. Pregatirea datelor .....	16
2.3.1. Tratarea erorilor.....	16
2.3.2. Tratarea valorilor care lipsesc .....	16
2.3.3. Conversii între tipuri de date .....	17
2.3.4. Tratarea valorilor extreme (outliers) .....	18
2.3.5. Transformarea datelor .....	18
2.4. Analiza datelor.....	20
2.4.1. Reducerea dimensionalității .....	21
2.5. Modelarea datelor .....	24
2.6. Prezentarea rezultatelor .....	26
<b>3. Aspecte teoretice și practice ale unui proiect data science bazat pe analiza de regresie .....</b>	<b>27</b>
3.1. Consideratii statistice.....	27
3.2. Regresia liniară.....	31
3.2.1. Ipoteze statistice asupra modelului regresiei liniare .....	33
3.2.2. Predicția punctuală .....	35
3.2.3. Estimarea parametrilor pe bază de interval de încredere .....	35
3.3. Regresia RANSAC (random simple consensus) .....	36
3.4. Regresia Decision Tree și Regresia Random Forest .....	37
3.5. Evaluarea performanței algoritmilor de machine learning .....	40
3.6. Metode de evaluare a performanței modelelor .....	41
3.6.1. Coeficientul de determinare .....	41
3.6.2. Coeficientul de determinare ajustat.....	42
3.6.3. Eroarea medie patratică (mean squared error) .....	43
3.6.4. Eroarea rădăcinii medie pătratică (root mean squared error).....	43
3.6.5. Media erorilor absolute (mean absolute error).....	43

3.6.6.	Graficul rezidurilor din regresie .....	44
3.6.7.	Testarea modelului de regresie pe baza statisticii $t$ Student .....	44
3.6.8.	Testarea modelului de regresie pe baza statisticii test $F$ .....	45
3.7.	Strategii de validare a modelelor .....	47
3.7.1.	Validarea simplă (seturi de date disjuncte de antrenare și testare) .....	47
3.7.2.	Validarea încrucișată (k-fold cross validation) .....	47
3.7.3.	Compensare bias - varianță .....	48
3.7.4.	Regresia Ridge .....	49
3.7.5.	Regresia Lasso.....	50
3.7.6.	Implementarea comparativă a algoritmilor de machine learning.....	51
<b>4.</b>	<b>Studiu de caz – predictia prețului autoturismelor prin analiza de regresie.....</b>	<b>52</b>
4.1.	Tehnologii utilizate pentru analiza și modelarea datelor .....	52
4.2.	Analiza și modelarea datelor .....	55
4.2.1.	Analiza preliminară .....	55
4.2.2.	Reducerea dimensionalității datelor .....	56
4.2.3.	Determinarea strategiei și a tehnicilor de modelare .....	59
4.2.4.	Regresia liniară simplă .....	61
4.2.5.	Regresia random forest cu o variabilă explicativă .....	64
4.2.6.	Regresia liniară multiplă .....	65
4.2.7.	Regresia random forest cu patru variabile explicative .....	69
4.2.8.	Model alternativ bazat pe regresia random forest .....	70
4.2.9.	Compararea rezultatelor modelelor dezvoltate .....	72
	<b>Bibliografie.....</b>	<b>77</b>

## Lista figurilor

Figura 1.1: Tabel greutate-consum.....	6
Figura 1.2: Scatterplot greutate si consum pt. autovehicule .....	7
Figura 1.3: Graficul funcției de regresie .....	8
Figura 1.4: Rezultate regresie liniara simpla OLS .....	9
Figura 1.5: Graficul valorilor reziduale.....	10
Figura 1.6: Intervalul de predicție/Intervalul de încredere.....	13
Figura 2.1: Etapele unui proiect data science.....	14
Figura 3.1: Tabelul ANOVA.....	46
Figura 3.2: Tabel validare încrucișată .....	48
Figura 4.1: Variabile inițiale car-price .....	55
Figura 4.2: Variabile esențiale car-price .....	56
Figura 4.3: Random Forest – importanța atributelor.....	57
Figura 4.4: Lasso – importanța atributelor .....	58
Figura 4.5: Reducerea nr. de observații în urma reducerii dimensionalității datelor .....	59
Figura 4.6: Matricea de corelație.....	59
Figura 4.7: Graficul perechilor de variabile .....	60
Figura 4.8: Rezultate regresia liniară simplă OLS .....	62
Figura 4.9: Graficul funcției de regresie liniară simplă .....	62
Figura 4.10: Normalitatea erorilor - regresia liniară simplă.....	63
Figura 4.11: Graficul valorilor reziduale – regresia liniară simplă.....	63
Figura 4.12: Intervalul de predicție și intervalul de încredere .....	64
Figura 4.13: Regresia random forest simplă: enginesize → price .....	65
Figura 4.14: Rezultate regresia liniară multiplă OLS .....	66
Figura 4.15: Normalitatea erorilor – regresia liniară multiplă .....	67
Figura 4.16: Graficul valorilor reziduale – regresia liniară multiplă .....	67
Figura 4.17: Evoluția coeficientului de determinare ajustat .....	69
Figura 4.18: Scatterplot regresia random forest cu 2 variabile explicative.....	70
Figura 4.19: Regresia random forest simplă: curbweight → price .....	71
Figura 4.20: Regresia random forest simplă: horsepower → price .....	71
Figura 4.21: Regresia random forest simplă: highwaympg → price .....	72
Figura 4.22: Tabel compararea rezultatelor modelelor .....	74
Figura 4.23: Punctaj comparativ modele implementate .....	74

## Lista cu notații/abrevieri

În lucrarea de față au fost utilizate următoarele notații/abrevieri:

$y$  = variabilă dependentă/target/aleatoare, valoarea inițială a variabilei de explicat ( $y_i$  reprezintă un vector)

$\hat{y}$  = valoarea estimată a variabilei dependente  $y$  prin aplicarea modelelor ( $\hat{y}_i$  reprezintă un vector)

$\bar{y}$  = media variabilei  $y$

$x$  = variabilă independentă/explicativă, predictor, atribut, variabilă care nu este aleatoare și pe baza căreia/căroră se încearcă explicarea variabilei dependente  $y$  ( $x_i$  reprezintă un vector)

$\bar{x}$  = media variabilei  $x$

$a, b$  = parametri/coeficienți de estimat ( $a$  = constanta,  $b$  = panta/coeficientul de regresie)

$\varepsilon$  = eroarea aleatoare neobservată sau reziduu, care conține toată informația care nu este explicată de relația liniară între  $x$  și  $y$

$\bar{\varepsilon}$  = media erorilor aleatoare / reziduuri

$i$  = numărul variabilelor cuprinse într-un vector

$n$  = numărul de observații

$k$  = numărul de variabile explicative

$\sigma$  = abaterea standard

$r_{x,y}$  = coeficientul de corelație între  $x$  și  $y$

$SST$  = variația totală a datelor  $y$  de intrare / suma pătratelor abaterilor individuale de la medie

$SSR$  = variația totală a datelor  $y$  estimate / suma pătratelor abaterilor de regresie

$SSE$  – variația între datele de intrare și cele estimate / suma pătratelor valorilor reziduale (erorilor)

$R^2$  = coeficientul de determinare

$R^2_{adj}$  = coeficientul de determinare ajustat

$MSE$  = eroarea medie patrată (mean squared error)

$RMSE$  = eroarea rădăcinii medie pătratică (root mean squared error)

$MAE$  = media erorilor absolute (mean absolute error)

$t$  = testul t Student

$F$  = testul F Fisher

## Introducere

Într-o lume aflată în continuă mișcare și evoluție permanentă, informația este resursa de bază care determină dezvoltarea societății umane. Forma brută a informației se regăsește în date. Datele reprezintă forma fizică și suporturile de reprezentare a informațiilor în cifre, litere, cuvinte, diverse simboluri și alte însemne. Datele sunt ansambluri de simboluri și nu au o semnificație în sine. Este necesară culegerea datelor, prelucrarea și interpretarea lor, pentru a fi transformate în informații cu sens, care să poată fi utilizate în diverse scopuri și să aducă utilitate în luarea deciziilor care ne influențează viața la nivel individual și social.

În societatea actuală, având în vedere dezvoltarea fără precedent a tehnologiei informațiilor și comunicațiilor, se generează și se vehiculează o cantitate uriașă de date, care crește exponențial.

Pentru a identifica elemente importante în vastul univers al datelor și pentru a da un sens și un scop fluxului complex de date, s-au dezvoltat instrumente, metode și tehnici de interpretare a datelor și de transformare a lor în informații utile.

Data science a apărut ca urmare a acestor necesități, ca urmare a nevoilor de tratare a unor cantități mari de date complexe și de extragere a informațiilor utile pe care acestea le conțin.

Data science este un domeniu interdisciplinar, care utilizează metode științifice pentru a extrage cunoștințe și perspective din date structurate și nestructurate, și aplică cunoștințele extrase într-o gamă largă de domenii. Data science este un concept revoluționar, capabil să abordeze cantitățile uriașe de date care sunt generate în zilele noastre, un concept care utilizează statistica și matematica, analiza datelor și informatica, ca instrumente și metode de măsurare și interpretare obiectivă a realității. Data science este un concept care utilizează instrumente și tehnici moderne pentru a identifica anumite tipare și relații existente între diverse date sau seturi de date, pentru a extrage informații semnificative și a descoperi noi structuri și relații, în scopul construirii unor decizii fundamentate științific. Data science se bazează pe inteligența artificială și pe subdomeniile ei machine learning și deep learning, pentru a crea algoritmi capabili să extragă și să modeleze seturi complexe de date, să facă clasificări și grupări în funcție de diverse criterii, să facă predicții și să identifice tendințe.

Data science are întrebuințări multiple, în diverse domenii de activitate, la nivel organizațional și instituțional. Rezultatele obținute în ultimii ani prin utilizarea data science au arătat că a devenit un domeniu indispensabil pentru cunoașterea realității economice și sociale și pentru elaborarea de decizii și strategii de acțiune. Toți actorii economici relevanți, din

principalele ramuri economice, utilizează într-o formă sau alta data science pentru a se informa și a lua decizii. Data science are beneficii majore la nivel instituțional, în business, în cercetare și inovare, iar companiile și statele care nu vor folosi avantajele aduse de data science riscă să rămână în urmă și să devină incapabile să se adapteze în timp util noilor provocări ale mediului competițional.

În prezenta lucrare este tratată analiza de regresie ca instrument de studiu în data science.

Introducerea în analiza de regresie s-a făcut printr-un studiu de caz de econometrie, în care a fost utilizată regresia liniară simplă pentru a arăta puterea unui model explicativ în identificarea unor estimatori care pot explica anumite fenomene și relații economice. În acest scop a fost preluat un set de date din cartea “Econometrie – La regression lineaire simple et multiple” publicată de Ricco Rakotomalala. Pe baza unui set de date simplu, care conține 28 de observații, o variabilă explicativă și o variabilă de explicat, s-a realizat estimarea consumului de combustibil al autovehiculelor în funcție de greutatea lor, prin modelarea cu ajutorul regresiei liniare simple prin metoda celor mai mici pătrate. Predicțiile elaborate sunt susținute de premise științifice și au la bază confirmarea unor ipoteze teoretice și calcularea unor indicatori statistici care arată performanțele modelului și susțin validitatea lui.

Pentru consolidarea analizei de regresie ca și model explicativ s-a procedat în continuare la efectuarea unui studiu de caz mai complex, care respectă etapele și principiile unui studiu de data science, pe baza unui set de date mai complicat. Setul de date a fost preluat de pe internet de pe site-ul Kaggle (<https://www.kaggle.com/hellbuoy/car-price-prediction>), care pune la dispoziție diverse seturi de date bazate pe măsurători practice, pentru studii de data science.

Obiectivul studiului este predicția prețului autovehiculelor pe baza unor attribute/variabile explicative preluate din viața reală prin măsurare. O companie din China - Geely Auto - dorește să intre pe piața auto din SUA și să înființeze o unitate de producție. A fost contactată o firmă de consultanță care a efectuat studii de piață și a obținut un set de date despre diverse tipuri de autoturisme de pe piața americană. Scopul studiului este de a afla care sunt factorii principali care influențează prețurile auto de pe piața americană și cum variază prețurile în funcție de variabilele explicative. Compania poate adapta designul produselor și poate să stabilească o strategie de business, în funcție de rezultatele studiului, care vor arăta dinamica prețurilor. Setul de date, conține 205 observații și are 26 de attribute (proprietăți măsurate), care pot explica mai mult sau mai puțin nivelul prețului de vânzare.



În etapa de analiză a datelor au fost identificate principalele atribute care pot explica prețul autovehiculelor, prin implementarea unor metode de analiză. Pentru selecția atributelor explicative relevante s-a utilizat filtrul corelației slabe între predictorii și variabila target, algoritmul random forest pentru determinarea importanței atributelor, regresia lasso (regularizarea L1) pentru selectarea atributelor. Ca urmare a aplicării acestor metode, din cele 26 de atribute au fost selectate patru, respectiv capacitatea motorului (enginesize), greutatea mașinii (curbweight), randamentul motorului/caii putere (horsepower) și indicatorul consumului de combustibil pe autostrada (highwaympg), toate cele patru fiind foarte corelate cu prețul.

Pentru stabilirea strategiei și a tehnicilor de modelare s-a utilizat graficul perechilor de variabile, care arată la nivel preliminar tipul de relație existentă între variabilele implicate și matricea de corelație, care arată un nivel de corelație ridicat între toate cele patru variabile explicative ( $> 0.75$ ). Această situație a corelației ridicate între variabilele explicative determină multicolinearitatea variabilelor, care duce la instabilitatea modelului, motiv pentru care se recomandă evitarea ei. Prin urmare dacă respectăm principiul multicolinearității trebuie să eliminăm toate variabilele înalt corelate între ele și să păstrăm doar o singură variabilă, cea mai bine corelată cu prețul, respectiv variabila enginesize.

Totuși în practica economică pentru elaborarea unei strategii de business (care se bazează pe întocmirea unui buget și a unui plan de producție) este nevoie de cât mai multă informație și trebuie utilizate toate variabilele semnificative. S-a pornit de la ipotezele că modelul ideal conform normelor teoretice, cu o variabilă explicativă, nu este folositor în realizarea obiectivului studiului, iar o variabilă explicativă nu poate aduce cantitatea de informație pe care o pot aduce mai multe variabile explicative.

Prin urmare, având în vedere complexitatea contextului, s-a stabilit utilizarea mai multor modele și compararea performanțelor lor, pentru a se putea trage o concluzie pertinentă cu privire la cel mai eficient model explicativ. Totodată s-au analizat în profunzime influențele multicolinearității asupra modelelor, pentru a vedea dacă modelele influențate de multicolinearitate sunt valide.

S-a utilizat regresia liniară simplă, regresia lasso, regresia ridge și regresia random forest. Au fost aplicate modelele simple, cu o variabilă explicativă și modelele multiple, cu patru variabile explicative. Totodată a fost dezvoltat și un model alternativ bazat pe patru modele de regresie random forest simplă, aplicate separat pentru fiecare variabilă explicativă. Rezultatul

final este determinat de media celor patru modele, ponderată cu coeficientul de determinare aferent fiecărui model.

Pentru determinarea celor mai performante modele au fost comparate rezultatele obținute de acestea. Indiferent de modelul implementat, scopul este ca distanța între valorile inițiale ale variabilei target (prețul în cazul nostru) și valorile estimate prin aplicarea modelului să fie cât mai mică. Întrucât valorile inițiale și valorile estimate se prezintă sub forma unor vectori, s-au comparat p-norme vectorilor rezultați din diferența între valorile inițiale și cele generate prin aplicarea modelelor.

Analiza și modelarea datelor pentru cele două studii s-a făcut în Python, cu utilizarea Jupyter Notebook și a pachetelor Numpy, Pandas, Matplotlib și Sklearn.

## 1. Econometrie – studiu de caz

Econometria s-a dezvoltat ca o necesitate în investigarea fenomenelor și proceselor economice, atât la nivel macroeconomic, cât și la nivel microeconomic. „*Econometria înseamnă aplicarea metodelor statistice datelor economice, pentru a da conținut empiric relațiilor economice*”. (M Hashem Pesaran, 1987, „Econometrics”)

Econometria utilizează teoria statistică, statisticile matematice și teoria probabilităților, aplicând modele statistice și matematice asupra datelor din lumea reală, ca unelte pentru identificarea unor estimatori care să exprime în mod imparțial și eficient diverse fenomene economice și relațiile dintre ele [1].

La nivel macroeconomic, econometria furnizează economiștilor instrumente prin care să analizeze istoria unor fenomene economice și evoluția lor în timp, să măsoare și să determine relații de dependență existente între diverși indicatori sau variabile economice, să facă previziuni aplicate în mediul economic. La nivel microeconomic, fiecare entitate economică poate utiliza econometria pentru analizarea mediului în care își desfășoară activitatea, pentru analizarea propriei activități economice, pentru determinarea unor strategii de management și marketing care să asigure funcționarea în condiții de profitabilitate.

Una dintre metodele statistice fundamentale utilizate în econometrie este analiza de regresie, utilizată pentru estimarea relațiilor între o variabilă dependentă și una sau mai multe variabile independente. Cea mai simplă formă a analizei de regresie este regresia liniară simplă, care presupune determinarea unei linii drepte care reprezintă cel mai bine relația între două variabile cantitative discrete, în funcție de un criteriu matematic specific [2].

Un model de regresie liniară simplă, bazat pe metoda celor mai mici pătrate, este reprezentat de următoarea ecuație:

$$y = a + bx + \varepsilon$$

Pentru ilustrarea modelului regresiei liniare simple, vom utiliza un exemplu, un **studiu de caz: estimarea consumului de combustibil al autovehiculelor în funcție de greutatea lor**.

Este un model simplu, pe baza unui eșantion cu 28 de observații, (datele fiind preluate din cartea “Econometrie La regression lineaire simple et multiple” publicată de Ricco Rakotomalala) conform Figura 1.1, în care avem:

➤ coloana **greutate** – **variabila independentă  $x$**  - care reprezintă greutatea în kg a vehiculelor

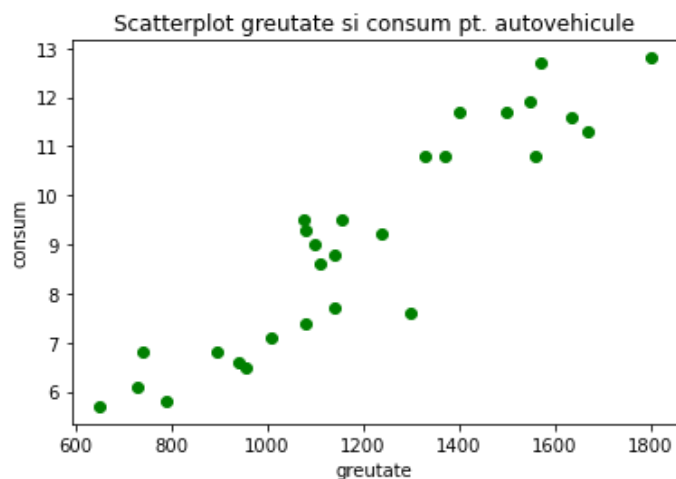
- coloana *consum* – *variabila dependentă y* - care reprezintă consumul de combustibil, în l/100km

n	greutate (x)	consum (y)
1	650	5.70
2	790	5.80
3	730	6.10
4	955	6.50
5	895	6.80
6	740	6.80
7	1,010	7.10
8	1,080	7.40
9	1,100	9.00
10	1,500	11.70
11	1,075	9.50
12	1,155	9.50
13	1,140	8.80
14	1,080	9.30
15	1,110	8.60
16	1,140	7.70
17	1,370	10.80
18	940	6.60
19	1,400	11.70
20	1,550	11.90
21	1,330	10.80
22	1,300	7.60
23	1,670	11.30
24	1,560	10.80
25	1,240	9.20
26	1,635	11.60
27	1,800	12.80
28	1,570	12.70

**Figura 1.1: Tabel greutate-consum**

Pe baza relației existente între cele două variabile/atribute, care au fost măsurate în prealabil, astfel că datele sunt cunoscute, putem estima variabila  $y$ .

Graficul scatter plot din Figura 1.2 (norul de puncte), arată existența unei relații liniare între cele două variabile. Relația liniară este pozitivă, întrucât creșterea valorilor variabilei  $x$  (axa  $ox$ ), se face concomitent cu creșterea valorilor variabilei  $y$  (axa  $oy$ ).



**Figura 1.2: Scatterplot greutate si consum pt. autovehicule**

### **Coeficientul de corelație:**

Pentru determinarea corelației între cele două variabile s-a calculat coeficientul de corelație Pearson, a cărui valoare poate fi cuprinsă între -1 și 1.

**Coeficientul de corelație este de 0.9263264**, deci este pozitiv, ceea ce înseamnă că avem o corelație directă, respectiv cele două variabile variază în același sens. Deoarece corelația se încadrează în segmentul [0.8-1] înseamnă că avem o corelație foarte înaltă, ceea ce arată că poate exista o legătură statistică între variabile.

### **Parametrii ecuației de regresie:**

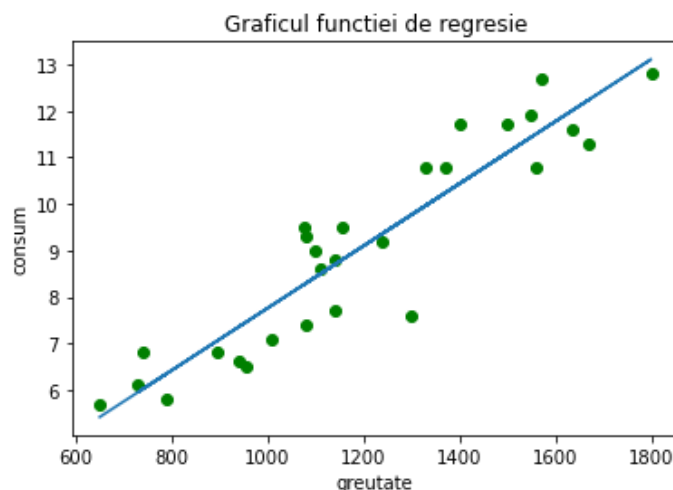
Utilizând metoda celor mai mici pătrate, calculăm estimatorii parametrilor din ecuația de regresie. Astfel, obținem următoarele valori:

- constanta/intercept este locul de pe ordonată unde dreapta de regresie se intersectează cu axa Oy și reprezintă o constantă care este valoarea medie a lui  $y$  atunci când  $x = 0$ :  **$a = 1.06269123$**
- panta de regresie/slope, care arată cu cât crește variabila  $y$ , atunci când  $x$  crește cu o unitate:  **$b = 0.00669386$ ,**

**Ecuația funcției liniare de regresie** este următoarea:

$$\hat{y}_i = [1.06269123] + [0.00669386]x_i$$

Vizualizarea grafică a funcției de regresie s-a făcut prin graficul scatter plot din Figura 1.3:



**Figura 1.3: Graficul funcției de regresie**

Se poate observa că dreapta de regresie trece aproximativ prin mijlocul norului de puncte. Evaluarea vizuală nu este însă suficientă pentru validarea modelului, astfel că avem nevoie de anumite criterii cantitative pentru evaluarea modelului.

### **Analiza varianței și coeficientul de determinare**

Obiectivul regresiei liniare este acela de a minimiza suma pătratelor erorilor, respectiv suma pătratelor diferențelor între valorile variabilei  $y$  și valorile estimate prin ecuația de regresie:

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

Variația erorilor/reziduurilor asociată dreptei de regresie arată abaterile punctelor digramei de dispersie de la dreapta de regresie, respectiv măsura dispersiei norului de puncte față de dreapta de regresie.

Variația totală  $SST$  este exprimată ca suma între variația explicată și variația neexplicată/variația erorilor:

$$SST = SSR + SSE$$

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

Pornind de la ecuația varianței, se poate determina coeficientul de determinare  $R^2$ , care descrie proporția varianței variabilei  $y$  explicată de modelul regresiei:

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

Cu cât valoarea lui  $R^2$  este mai aproape de 1, cu atât modelul este mai reprezentativ, iar variabila  $x$  permite determinarea valorilor estimate ale variabilei  $y$ , întrucât valoarea  $SSE$  este mai aproape de 0, deci obiectivul minimizării sumei pătratelor erorilor este atins. Dacă valoarea lui  $R^2$  este mai aproape de 0, atunci variația erorilor nu este minimizată, iar valoarea  $SSE$  se apropie de 1, ceea ce arată faptul că variabila  $y$  nu poate fi explicată prin variabila  $x$ .

În cazul modelului prezentat, **valoarea coeficientului de determinare este de 0.85808059**, ceea ce arată că modelul este reprezentativ pentru explicarea variabilei dependente.

### Ipotezele asupra proprietăților estimatorilor

Pentru verificarea ipotezelor, au fost efectuate mai multe teste și au fost calculați mai mulți indicatori, prezentați în Figura 1.4:

OLS Regression Results						
Dep. Variable:	y	R-squared:	0.858			
Model:	OLS	Adj. R-squared:	0.853			
Method:	Least Squares	F-statistic:	157.2			
Date:	Sun, 18 Apr 2021	Prob (F-statistic):	1.58e-12			
Time:	16:35:46	Log-Likelihood:	-34.378			
No. Observations:	28	AIC:	72.76			
Df Residuals:	26	BIC:	75.42			
Df Model:	1					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Intercept	1.0627	0.659	1.612	0.119	-0.292	2.418
x	0.0067	0.001	12.538	0.000	0.006	0.008
Omnibus:	1.054	Durbin-Watson:	1.799			
Prob(Omnibus):	0.590	Jarque-Bera (JB):	0.883			
Skew:	-0.409	Prob(JB):	0.643			
Kurtosis:	2.705	Cond. No.	5.03e+03			

**Figura 1.4: Rezultate regresie liniara simpla OLS**

Pentru validarea modelului au fost verificate următoarele ipoteze de lucru:

#### ***Liniaritatea modelului***

Verificarea liniarității s-a făcut grafic, prin intermediul graficului scatter plot prezentat în Figura 1.3. Întrucât norul de puncte care arată relația dintre variabila predictor și cea estimată este dispus în model liniar, se poate spune că există o relație liniară între cele două variabile. Coeficientul de corelație Pearson care are valoarea 0.9263264, confirmă ipoteza existenței unei relații liniare între cele două variabile.

**Media erorilor este aproape de 0:**  $\bar{\varepsilon} = 0$

Media erorilor este 1.9032394, fiind foarte apropiată de 0, ceea ce confirmă ipoteza.

**Normalitatea erorilor** – variabila  $\varepsilon$  urmează o lege normală de medie zero și variație  $\sigma^2$ :  
 $\varepsilon_i \sim N(0, \sigma^2)$ ;

Normalitatea erorilor a fost verificată prin testul Jarque-Bera. A fost calculat coeficientul de asimetrie (Skewness) și coeficientul de boltire (Kurtosis) a distribuției erorilor reziduale. O distribuție perfect normală are un coeficient de asimetrie  $S = 0$  și un coeficient de boltire  $K = 3$ .

- $S > 0$  arată o repartiție asimetrică la dreapta,  $S < 0$  arată o repartiție asimetrică la stânga
- $K > 3$  arată o repartiție afectată de boltire, iar  $K < 3$  arată o repartiție aplatizată

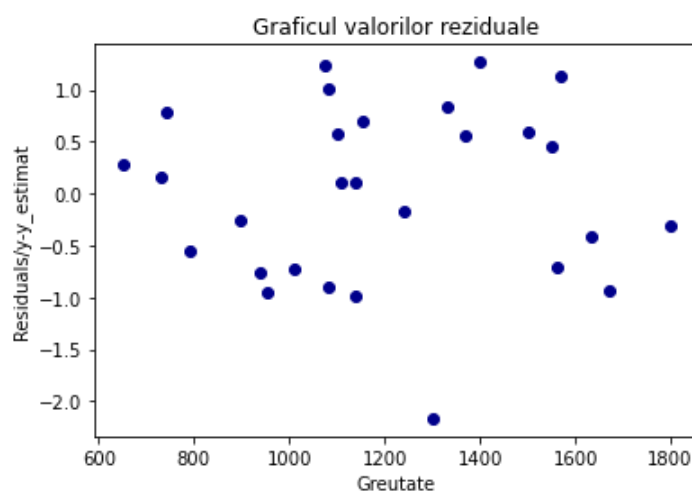
După cum se observă în Figura 1.4:

- coeficientul de asimetrie (Skew) = -0.409, deci este foarte apropiat de 0
- coeficientul de boltire (Kurtosis) = 2.705, deci este foarte apropiat de 3
- valoarea calculată a testului Jarque-Bera = 0.883, depășește foarte puțin valoarea teoretică de 0.643

Conform teoremei limitei centrale, cu cât eșantionul de date crește, cu atât distribuția tinde să fie normală. Dacă mărimea eșantionului depășește 30 de observații, distribuția mediei aritmetice a eșantioanelor va fi o distribuție normală. Având în vedere faptul că setul de date conține doar 28 de observații, putem considera că erorile reziduale sunt distribuite normal.

**Homoscedasticitatea (omogenitatea varianței erorilor)** – varianțele erorilor sunt constante oricare sunt valorile variabilei predictor  $x$ :  $v(\varepsilon_i) = \sigma^2$

Verificarea ipotezei homoscedasticității s-a făcut prin reprezentarea grafică a relației dintre variabila  $x$  (axa ox) și reziduuri (axa oy), conform Figura 1.5.



**Figura 1.5: Graficul valorilor reziduale**



Se poate observa că norul de puncte se află dispus normal, fără a forma un anumit model sau tipar, astfel că ipoteza homoscedasticității se confirmă.

**Necorelarea erorilor** – erorile sunt necorelate între ele:  $cov(\varepsilon_i, \varepsilon_j) = 0$ ; erorile asociate unor valori ale variabilei  $y$  nu sunt influențate de erorile asociate altor valori ale variabilei  $y$ .

În analiza de regresie autocorelarea erorilor intervine dacă se încearcă modelarea unei relații liniare asupra unor date neliniare, caz în care reziduurile vor fi autocorelate.

Ipoteza necorelării erorilor a fost verificată prin testul Durbin Watson. Valoarea testului Durbin Watson este de 1.799 și se compară cu limita inferioară -  $d_L$  și cu limita superioară -  $d_U$  citite din tabela Durbin și Watson. Întrucât valoarea calculată obținută se află în intervalul 1.325 – 1.964 obținut din tabela cu valorile teoretice, putem considera că nu există autocorelare între valorile reziduale.

### **Testarea semnificativității pantei de regresie**

Demersul pornește de la formularea ipotezei nule  $H_0$ , conform căreia variabila  $y$  nu este influențată de variația variabilei predictor  $x$  și deci coeficientul  $b$  din ecuația de regresie nu este semnificativ diferit de zero.

$$H_0: b = 0$$

$$H_1: b \neq 0$$

Pentru verificarea ipotezei nule s-a folosit testul  $t$  Student pentru coeficientul de corelație simplă, valoarea  $t$  calculată fiind de 12.538. Valoarea calculată se compară cu valoarea teoretică din tabela  $t$  Student pentru  $n-2$  grade de libertate și un nivel de semnificație de 5%.

Întrucât  $t_{calculat} = 12.538 > t_{teoretic} = 1.58$  se respinge ipoteza nulă. Coeficientul de corelație este semnificativ diferit de zero, iar între variabila independentă și variabila dependentă există o legătură semnificativă.

### **Testarea semnificativității globale a modelului de regresie**

Evaluarea globală a modelului s-a făcut pe baza statisticii test  $F$ . Valoarea calculată a testului este de 157.2.

Întrucât  $F_{calculat} = 157.2 > F_{teoretic} = 1.58$ , pentru  $n-2$  grade de libertate și un nivel de semnificație de 5%, se consideră că modelul este semnificativ.

**În concluzie**, întrucât modelul implementat a trecut testele de evaluare, se poate considera reprezentativ și poate fi utilizat pentru estimarea variabilei dependente - consumul de combustibil, în funcție de variabila independentă - greutatea autovehiculului.

### Predicția punctuală

Pe baza modelului construit se pot face previziuni punctuale ale comportamentului variabilei dependente  $y$ , în funcție de valorile fixe ale variabilei independente  $x$ .

Pentru exemplificare s-a construit o metodă de predicție pe baza ecuației de regresie. Dorim să determinăm consumul unui autovehicul cu o greutate de  $x = 1140$  kg. Calculăm predicția punctuală astfel:

$$y = \hat{a} + \hat{b}x = 1.062691 + 0.006694 * 1140 = 8.69368917$$

Se estimează astfel că un vehicul cu o greutate de 1140 kg va consuma 8.69 litri la 100 km.

### Predicția pe interval de predicție

Pentru construirea unui interval de predicție avem nevoie de un interval de încredere în cadrul căruia se găsește cu o probabilitate semnificativă parametrul estimat, astfel că estimatorul să se găsească între valoarea inferioară și valoarea superioară a intervalului.

Pentru exemplificare s-a construit o metodă de predicție, pe baza intervalului de predicție, cu un nivel de încredere de  $(1-\alpha)$  după următoarea formulă:

$$\hat{y}_i \pm t_{1-\frac{\alpha}{2}} * \sigma_{\varepsilon} \sqrt{1 + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

Pentru calcularea intervalului de încredere luăm un risc  $\alpha$  de 5% și  $n-2$  grade de libertate, și citim din tabelul  $t$  Student valoarea 1.705618.

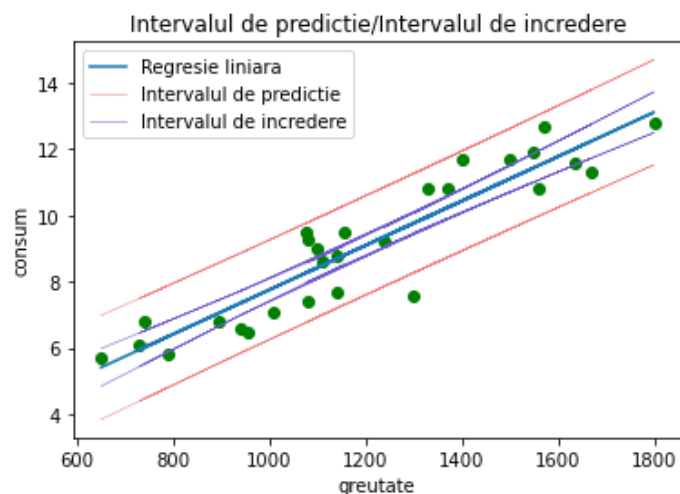
Pentru  $x = 1140$  kg se calculează intervalul de predicție după cum urmează:

$$[8.69368917 - 1.705618 * 0.872893 ; 8.69368917 + 1.705618 * 0.872893]$$

$$[7.204868 ; 10.182510]$$

Se estimează astfel că un vehicul în greutate de 1140 kg are 90% șanse să consume între 7.20 și 10.18 litri la 100 km.

Pentru vizualizare s-a construit graficul din Figura 1.6, care arată intervalul de încredere și intervalul de predicție.

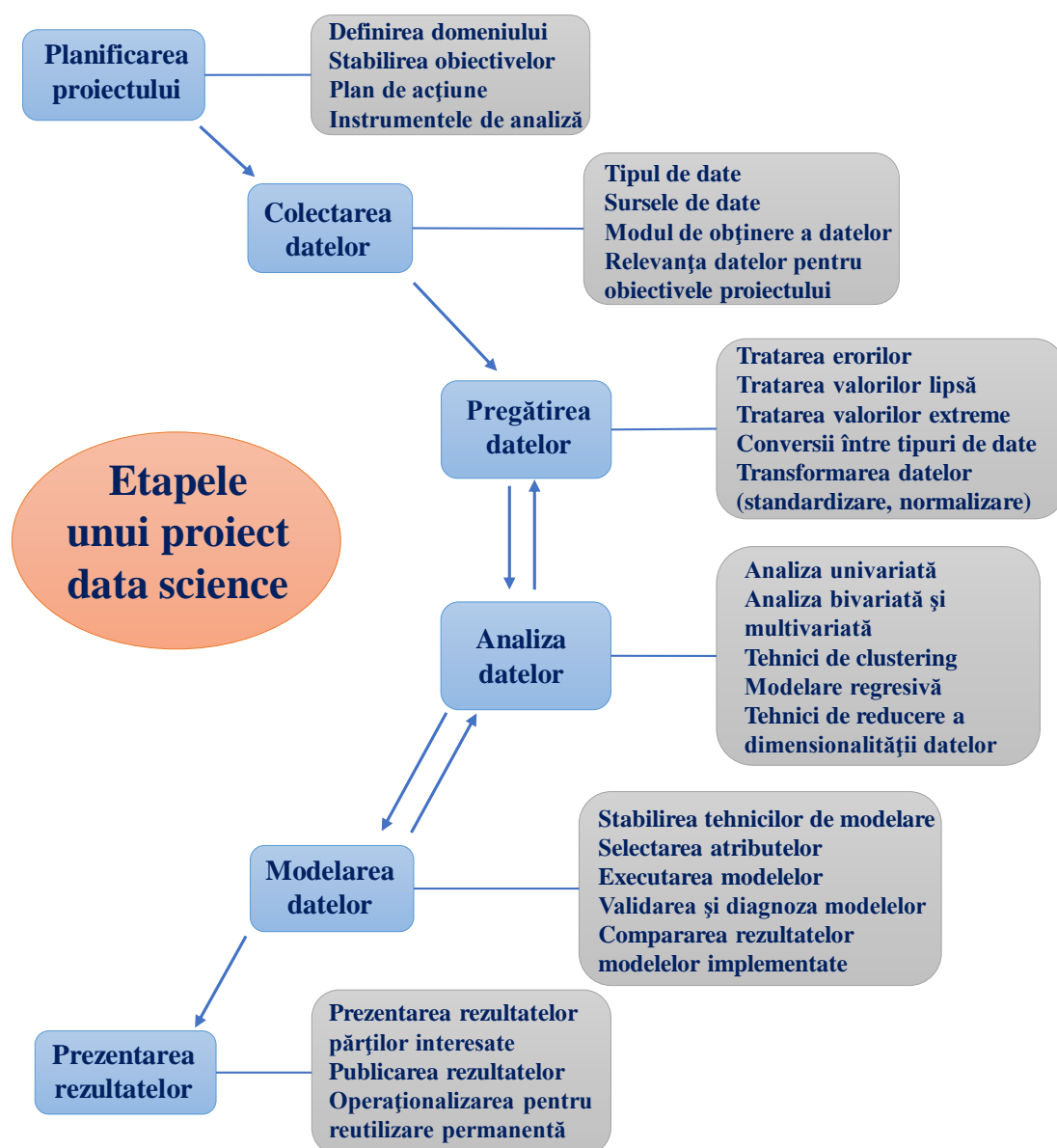


***Figura 1.6: Intervalul de predicție/Intervalul de încredere***

Modelul de regresie liniară simplă pe baza metodei celor mai mici pătrate este un model puternic, care poate fi utilizat pe multiple seturi de date din lumea reală. Modelul este consolidat pe fundamente teoretice solide, iar predicțiile elaborate pe baza lui sunt susținute de premise testate și verificate.

## 2. Etapele și conținutul unui proiect Data Science

Un proiect Data Science este un demers complex, care trebuie să parcurgă o serie de etape și să îndeplinească anumite condiții pentru a fi un proiect viabil și util. Etapele nu urmează o ordine strictă, de multe ori fiind necesară revenirea la etape anterioare, în funcție de rezultatele obținute pe parcursul demersului [3], care este prezentat în Figura 2.1.



*Figura 2.1: Etapele unui proiect data science*

## 2.1. Planificarea proiectului

Un proiect de data science trebuie să înceapă prin definirea domeniului studiat, stabilirea obiectivelor proiectului, datele și resursele necesare pentru colectarea lor, stabilirea planului de acțiune și a etapelor de urmat.

Proiectul trebuie să aibe un scop și obiective de cercetare foarte clar stabilite, un plan de acțiune, termene de finalizare pentru etapele proiectului, pentru a putea obține rezultate bine justificate. Înțelegerea obiectivelor studiului, a contextului în care se desfășoară cercetarea, stabilirea surselor de colectare a informațiilor și a instrumentelor de efectuare a analizei, identificarea resurselor necesare și disponibile, sunt condiții esențiale pentru reușita proiectului.

## 2.2. Obținerea/colectarea datelor

Primele demersuri care trebuiesc făcute pentru colectarea datelor sunt stabilirea tipului de date necesare analizei, stabilirea surselor de date și a modului de obținere a datelor. Datele pot fi colectate din diverse surse și pot fi obținute în diverse formate.

În funcție de domeniul de acțiune și de obiectivele analizei, datele pot fi obținute din mai multe surse [4]:

- prin interogarea a diverse tipuri de baze de date, dacă proiectul se desfășoară în cadrul unei companii și avem acces la bazele de date sau la diverse informații
- prin web scraping, de pe diverse website-uri
- prin intermediul Web APIs de pe website-uri; Facebook sau Twitter de exemplu permit utilizatorilor să se conecteze la serverele lor web și să acceseze diverse date prin Web API-urile lor
- importante seturi de date pot fi descărcate direct de pe diverse site-uri (inclusiv site-uri guvernamentale), ca de exemplu Kaggle, UCI Machine Learning Repository, Data World, Dataset Search de la Google, AWS Public Data Sets de la Amazon și multe altele

Odată obținute datele, este necesar să se facă verificarea relevanței lor pentru scopul și obiectivele proiectului, respectiv trebuie verificat dacă colecția de date asigură resursele necesare pentru efectuarea analizei și modelării.

## **2.3. Pregătirea datelor**

Pregătirea datelor presupune organizarea și structurarea datelor brute pentru transformarea lor în date viabile pentru analiză și modelare. Este o etapă care se concentrează pe asigurarea calității și consistenței datelor, astfel ca datele să devină reprezentative pentru procesele de analiza și modelare.

Etapa pregătirii datelor are 3 subetape: curățarea datelor, integrarea datelor, transformarea datelor [3]. Pregătirea datelor are ca scop eliminarea datelor inconsistente sau completarea datelor lipsă acolo unde este posibil, astfel încât să nu se influențeze rezultatul.

### **2.3.1. Tratarea erorilor**

Există diverse erori în cadrul setului de date, atât pentru datele numerice, cât și pentru cele categoricale, care afectează analiza și modelarea datelor. Aceste erori trebuie identificate și eliminate, pentru a nu produce rezultate eronate.

În cadrul variabilelor categoricale trebuie găsite și eliminate spațiile libere prezente la începutul sau sfârșitul unei valori de tip string. De asemenea trebuie evitată capcana sensibilității la majuscule ale datelor de tip string. În ceea ce privește variabilele numerice, trebuie verificate valorile imposibile, de ex. vârsta de peste 100 de ani, înălțimea de peste 2 metri, etc.

### **2.3.2. Tratarea valorilor care lipsesc**

Unul dintre primele aspecte care trebuie analizate la un set de date, este prezența spațiilor libere, a elementelor lipsă. Există diverse metode pentru tratarea datelor lipsă, printre care, eliminarea instanțelor (linii) care conțin observații cu valoare zero sau a celor fără valoare (nan), eliminarea unor atribute (coloane) în cazul în care conțin prea multe valori lipsă și nu pot fi luate în calcul la analiza și modelarea datelor [5]. Datele lipsă pot fi înlocuite și cu indicatori statistici adecvați, de ex: media, mediana sau modulul coloanei pe care se găsesc datele, acolo unde este posibil și înlocuirea nu va influența rezultatul final.

În cazul în care există o relație liniară între 2 sau mai multe variabile, se poate utiliza și regresia liniară pentru identificarea și înlocuirea valorilor lipsă.

### 2.3.3. Conversii între tipuri de date

Datele care fac obiectul analizei pot fi numerice, categoriale, binare, imagini, etc. În scopul analizei și modelării datelor se poate impune codificarea sau convertirea lor dintr-un tip de date în alt tip de date.

Asupra datelor numerice (cantitative) reprezentate de numere întregi sau reale, se pot face operații aritmetice de adunare, scădere, înmulțire, împărțire, se pot calcula indicatori statistici de tendință centrală, de variație, etc. De asemenea se poate defini o ordine între valorile atributelor numerice, un minim, un maxim, se pot încadra în anumite intervale sau cuantile. Exemple de date numerice pot fi vârsta, greutatea, prețul, temperatura, etc.

Datele categoriale sau discrete sunt cele care au două sau mai multe categorii, și ele pot fi nominale sau ordinale.

Datele categoriale nominale sunt valori discrete care nu au o ordine intrinsecă și pot fi simboluri, caractere, șiruri de caractere, etc. Exemple de date nominale sunt genul (F/M), rasa, starea civilă, etc. Asupra datelor nominale se pot efectua operații de calcul al frecvențelor sau de verificare a egalității.

Datele categoriale ordinale pot fi sub formă de numere, simboluri, șiruri de caractere, și sunt date între care există o ordine clară în cadrul mulțimii de valori. Exemple de date ordinale sunt temperatura (ridică, medie, scăzută), nivele de calitate, etc. Asupra acestui tip de date nu are sens să se efectueze operații aritmetice.

Datele binare sunt date care acceptă doar două valori posibile (0 sau 1, adevărat sau fals) și se utilizează pentru a codifica prezența sau absența unor caracteristici, în general notându-se cu 1 prezența caracteristicii și cu 0 absența ei.

Datele categoriale pot fi convertite în echivalente numerice sau binare prin operații de codare (encoding), întrucât anumite tehnici de machine learning (de exemplu: regresia liniară, support vector machine, neural networks) pot fi aplicate doar pentru date numerice. De asemenea, operațiunile de normalizare, standardizare, analiza în componente principale pot fi aplicate doar atributelor numerice.

Conversia unui atribut numeric în unul categorial se numește discretizare (de exemplu vârsta poate fi împărțită pe categorii 20-35, 36-45, etc. și se asignează câte o valoare fiecărui subdomeniu) sau binarizare (binning) dacă este vorba despre conversia în attribute binare.

### **2.3.4. Tratarea valorilor extreme (outliers)**

Identificarea valorilor extreme (outliers) și studierea naturii acestora în vederea determinării modului de tratare a acestor valori, este un alt element important al pregătirii datelor. Valorile extreme pot fi identificate prin utilizarea histogramelor (cea mai întâlnită distribuție în științele naturii este distribuția Gaussiană) [3], a graficelor scatter plot, a tabelelor de frecvență (în cazul variabilelor categoriale), etc.

Valorile extreme pot fi rezultatul unor metode inadecvate de măsurare sau colectare a datelor, caz în care sunt anomalii care trebuie eliminate. De exemplu dacă într-un eșantion avem înălțimea unor persoane în funcție de sexul lor și observăm valori anormale, atunci aceste date trebuie tratate ca anomalii și eliminate.

Există și cazuri în care valorile extreme pot fi observații reale care ies din tiparul general al populației statistice măsurate, caz în care modul de tratare a acestor valori se face în funcție de situația dată. Dacă datele extreme nu modifică rezultatul, dar influențează ipotezele, atunci ele pot fi eliminate. Dacă datele extreme afectează și rezultatul și ipotezele, atunci nu este indicată eliminarea lor, ci este recomandabil să se facă modelarea datelor atât în varianta cu includerea valorilor extreme, cât și în varianta cu eliminarea lor, prezentându-se rezultatele comparativ. De asemenea, se pot testa comparativ mai multe modele pentru a vedea care este mai potrivit.

Identificarea și tratarea valorilor extreme se poate face și prin utilizarea modelelor statistice robuste. Media și abaterea standard sunt sensibile la outliers, de aceea în locul mediei poate fi utilizată mediana, iar în locul deviației standard se poate utiliza deviația cvartilă. Distribuția normală și abaterea standard pot fi utilizate în eliminarea datelor extreme, întrucât conform regulii celor 3 sigma pentru variabilele aleatoare normale, 68% din datele unei distribuții normale se situează într-o deviație standard, 95% din date se află în două deviații standard, iar 99.7% se situează în 3 deviații standard, probabilitatea ca abaterea în valoare absolută să depășească 3 sigma este de 0.026, adică practic 0.

În cazul modelului regresiei liniare, se poate utiliza comparativ regresia robustă pentru identificarea datelor extreme, de exemplu regresia RANSAC.

### **2.3.5. Transformarea datelor**

Datele trebuie să fie transformate într-o formă adecvată pentru a permite măsurarea adecvată a proceselor pe care le exprimă, astfel că ele trebuie normalizate sau standardizate



astfel ca datele pe coloană să fie de același tip, să nu existe unități de măsură diferite ale observațiilor unei variabile. Normalizarea și standardizarea datelor sunt metode de scalare a variabilelor.

**Standardizarea** presupune scalarea datelor pentru a avea media 0 și abaterea standard 1 (se utilizează de exemplu la analiza în componente principale pentru calcularea matricei de covarianță) și este utilă în cazul datelor care urmează o distribuție normală, Gaussiană [6]. Standardizarea se face după următoarea formulă:

$$x' = \frac{x - \bar{x}}{\sigma}$$

**Normalizarea** (Min-Max scalling) presupune scalarea datelor astfel ca valorile să se încadreze într-un interval prestabilit, între 0 și 1. Normalizarea este indicat să se folosească atunci când datele nu urmează o distribuție normală, Gaussiană, cum sunt algoritmi neural networks sau algoritmi de clasificare bazați pe calcularea distanțelor [6].

Normalizarea se face după următoarea formulă:

$$x' = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Algoritmi de machine learning bazați pe arbori (decision trees, random forest, gradient boosting) utilizați pentru probleme de clasificare sau regresie, precum și regresia liniară bazată pe modelul celor mai mici pătrate, nu necesită scalarea variabilelor.

Regresia liniară, regresia logistică, rețelele neurale, etc., care utilizează tehnica gradient descent (pasul descendent) ca și tehnică de optimizare, necesită scalarea datelor, deoarece diferența de mărime între valorile variabilelor independente va determina praguri diferite pentru fiecare variabilă [6]. Utilizarea variabilelor scalate similar va ajuta algoritmul să convergă mai repede spre valoarea minimă, în timp ce pașii vor fi actualizați la aceeași rată pentru toate variabilele independente.

Algoritmi de clasificare bazați pe calcularea distanțelor, cum sunt K-Nearest Neighbors (KNN), K-means sau Support Vector Machine (SVM) sunt de asemenea influențați de intervalul de mărime a variabilelor, deoarece acești algoritmi utilizează calcularea distanțelor între puncte pentru a determina similaritatea lor [6]. De aceea, pentru creșterea performanței algoritmului și pentru ca algoritmul să nu fie afectat de bias, se recomandă scalarea datelor.

## 2.4. Analiza datelor

Datele colectate și pregătite sunt supuse analizei exploratorii pentru identificarea caracteristicilor principale, obiectivul fiind înțelegerea datelor și a relațiilor dintre ele, identificarea unor rezultate statistice semnificative. Analiza exploratorie a datelor contribuie la identificarea anumitor tipare, la descoperirea eventualelor anomalii și date extreme care au trecut de etapa de pregătire a datelor, la extragerea variabilelor semnificative și reducerea dimensionalității setului de date, la formularea și testarea unor ipoteze de lucru, la alegerea celor mai adecvate tehnici de modelare a datelor.

În sensul analizei datelor exploratorii sunt utilizate diverse tehnici și instrumente statistice care includ [7]:

- analiza univariată, care presupune utilizarea indicatorilor statistici în vederea studierii datelor în funcție de o singură variabilă pentru a se vedea principalele caracteristici și tendințe, astfel că nu este nevoie de o analiză care să arate raporturi de cauzalitate sau corelație;
- analiza bivariată și multivariată care presupune utilizarea metodelor statistico-matematice pentru identificarea legăturilor de asociere existente între două sau mai multe variabile;
- tehnici de clustering pentru identificarea eventualelor tendințe de grupare a datelor;
- tehnici de reducere a dimensionalității datelor care contribuie la creșterea performanței tehnicilor de modelare și la reducerea considerabilă a sistemului;
- modelare regresivă.

Analiza statistică univariată și multivariată se poate face prin tehnici statistice non-grafice sau prin vizualizare grafică, întrucât tehnicile fără vizualizare grafică nu pot reda imaginea distribuirii datelor și a relațiilor dintre ele.

Analiza univariată și multivariată prin reprezentare grafică se poate face prin diverse grafice univariate sau multivariate ca de exemplu: histograme sau diagrame de frecvență, curba densității de probabilitate, diagrama boxplot care arată deviația quartilă, diagrama cu coloane, diagrame cu bare, diagrame cu linii, grafice scatter plot care arată relația dintre variabile, diagrama de procente (pie-chart), grafice time-series, grafice heat map, etc.

O etapă importantă a analizei exploratorii a datelor este reducerea dimensionalității, care presupune reducerea numărului de atribute prin care poate fi explicat modelul, în scopul eliminării informațiilor redundante și a eficientizării implementării tehnicilor algoritmice de modelare a datelor.

### 2.4.1. Reducerea dimensionalității

Reducerea dimensionalității este una dintre cele mai importante etape în analiza datelor, înainte de modelarea lor. În machine learning dimensionalitatea se referă la numărul variabilelor prezente în setul de date. Atunci când numărul variabilelor independente este mare, algoritmi pot avea probleme de performanță, iar aceasta se numește „blestemul dimensionalității” [8].

Reducerea dimensionalității înseamnă transformarea datelor dintr-un spațiu cu un număr mare de dimensiuni/atribute, într-un spațiu cu un număr redus de dimensiuni/atribute. Transformarea trebuie să se realizeze fără a pierde prea mult din informația inițială, astfel ca reprezentarea să elimine atributele irelevante, dar prin reducerea dimensiunilor să păstreze proprietățile datelor inițiale într-un procent semnificativ.

Reducerea dimensiunii poate fi abordată prin selectarea atributelor (feature selection) și extragerea atributelor (feature extraction).

*Modelele bazate pe extragerea atributelor* determină un set de date cu un număr mai mic de atribute decât setul original, dar nu păstrează caracteristicile originale și crează caracteristici noi din combinarea matematică a atributelor inițiale. Pentru acest motiv modelele feature extraction nu sunt fiabile în cazul în care asupra datelor urmează să fie aplicați algoritmi de machine learning care presupun predicții sau previziuni ale unor variabile pe baza altor variabile, cum sunt modelele de regresie. Unii algoritmi au încorporat modele de extragerea atributelor, cel mai bun exemplu fiind deep learning. De asemenea, extragerea atributelor poate fi nesupervizată prin analiza în componente principale (PCA) sau supervizată prin analiza discriminantă liniară (LCA).

*Modelele de reducere a dimensionalității prin selectarea atributelor* păstrează un subset al caracteristicilor inițiale, motiv pentru care sunt fiabile pentru datele asupra cărora se vor aplica modele de regresie. Unii algoritmi de învățare supervizată cum sunt modelele de regresie regularizate sau random forests, au încorporat modele de reducere a dimensionalității prin selectarea atributelor. Astfel regresia ridge regularizează coeficienții dar păstrează toate atributele, regresia lasso regularizează coeficienții pentru a efectua selecția atributelor, iar elastic nets combină aceste două tipuri de regresie regularizată.

Reducerea dimensionalității are mai multe avantaje [9]:

- reducerea numărului atributelor înseamnă reducerea spațiului necesar pentru stocarea datelor;

- evitarea multicoliniarității prin eliminarea atributelor redundante;
- algoritmi de machine learning lucrează mai eficient pe seturi de date cu dimensiuni reduse, iar modelele rezultate sunt mai simple și datele rezultate sunt mai ușor de vizualizat.

Există mai multe metode de reducere a dimensionalității datelor, unele mai simple, altele mai complexe, printre care:

#### ***Filtrul corelației scăzute între variabilele independente și variabila dependentă***

Corelația exprimă interdependența sau legătura dintre două variabile, fără a exista o relație de cauzalitate între ele. Astfel, o variabilă dependentă care nu este corelată cu variabila target, nu poate avea o pondere importantă în explicarea comportamentului variabilei target, motiv pentru care poate fi eliminată din model [9].

#### ***Filtrul corelației ridicate între variabilele independente***

Un indicator mare al corelației între două variabile dependente arată faptul că cele două variabile au tendințe similare și este posibil să aducă în model informații similare, redundante, determinând multicoliniaritatea modelului.

Pentru evitarea acestor probleme se calculează corelația între variabilele independente (se poate utiliza matricea de corelație), iar dacă două variabile sunt foarte corelate (coeficientul de corelație depășește 0,5-0,6), atunci una poate fi eliminată fără a pierde foarte mult din informație [10]. Vom păstra variabila care este mai corelată cu variabila dependentă.

#### ***Eliminarea progresivă (backward feature elimination)***

Eliminarea progresivă se face în mai mulți pași. Se începe prin antrenarea modelului utilizând toate variabilele prezente în setul de date, iar apoi se calculează performanța modelului (pe baza coeficientului de determinare ajustat). Ulterior se elimină pe rând câte o variabilă și se aplică modelul la toate seturile de variabile  $n-1$ , calculându-se performanța pentru fiecare variantă. Se identifică variabila a cărei eliminare a produs cea mai mică variație a performanței modelului și se elimină din setul de date. Se repetă procedura, până când nu se mai pot elimina variabile [9].

#### ***Selecția progresivă (forward feature selection/construction)***

Este inversa procedurii eliminării progresive, întrucât în loc să se elimine pe rând variabilele, se încearcă găsirea celei mai utile variabile care îmbunătățește performanța modelului.

Se începe prin aplicarea modelului pe rând, la câte o singură variabilă independentă, până epuizăm toate cele  $n$  variabile independente. Se selectează ca și variabilă de pornire, variabila care determină cea mai mare performanță a algoritmului. Se repetă procesul prin adăugarea a câte o variabilă, pe rând, din cele  $n-1$  variabile independente rămase, iar variabila care determină cea mai bună performanță se reține pentru iterațiile viitoare. Se repetă procesul până când nu se mai identifică o îmbunătățire a performanței modelului [9].

Eliminarea progresivă și selecția progresivă sunt tehnici care se aplică modelelor de regresie, dar întrucât sunt metode consumatoare de timp și resurse, aceste tehnici sunt utilizate atunci când avem un set de date cu un număr redus de atribute. Variabilele se pot elimina prin compararea performanțelor modelului aplicat la diversele variante ale setului de date pe baza coeficientului de determinare ajustat sau selecția variabilelor se oprește când valorile  $t$  calculate devin mai mici decât valoarea critică citită din tabela Student.

### ***Random forest regressor – importanța atributelor***

Random forest este unul dintre cei mai utilizați algoritmi pentru feature selection, întrucât are încorporat un model de determinare a semnificativității atributelor în explicarea modelului [9].

Prin random forest, care este un model de tip ansamblu (ensemble learning), se generează un set de arbori de decizie cu privire la o variabilă target și apoi se utilizează statistica de performanță a fiecărui atribut predictor pentru a identifica subsetul de atribute care furnizează cea mai mare cantitate de informație cu privire la model.

Algoritmul generează un set de arbori, fiecare dintre aceștia fiind antrenat cu o parte mică din totalul atributelor predictor, luate în diverse combinații. Dacă un atribut este selectat de multe ori ca și criteriu de ramificare (best split), atunci este foarte probabil că acesta este un atribut semnificativ al modelului.

După generarea ansamblului de arbori, se calculează scorul fiecărui atribut, în funcție de câte ori a fost ales ca și criteriu de ramificare și la ce nivel, în cadrul atributelor candidate.

$$\text{Scor} = \text{ramificații}(\text{niv.0})/\text{candidați}(\text{niv.0}) + \text{ramificații}(\text{niv.1})/\text{candidați}(\text{niv.1}) + \dots$$

Scorul arată care sunt atributele cele mai importante în explicarea modelului. Se vor reține doar atributele care înregistrează un scor mai mare decât un prag stabilit.

### ***Regresia lasso pentru selectarea atributelor***

Regresia Lasso este un algoritm care crează un model regularizat al regresiei liniare, prin adăugarea unui parametru de regularizare. Prin modificarea acestui parametru în sensul creșterii

sau descreșterii lui, coeficienții ecuației de regresie liniară scad către valoarea 0. Variabilele a căror coeficient atinge valoarea 0 sunt variabilele care pot fi eliminate, întrucât contribuția lor în cadrul modelului este nesemnificativă [10].

## 2.5. Modelarea datelor

După finalizarea operațiunilor de pregătire și analiză a datelor, se poate trece la modelarea lor prin tehnici de statistică, data mining sau machine learning. Tehnicile de modelare trebuiesc alese în funcție de tipul de date care sunt supuse analizei și de scopul analizei.

Etapele procesului de modelare cuprind [3]:

- stabilirea tehnicilor de modelare și selectarea atributelor/variabilelor care se modelează
- executarea și validarea modelelor
- diagnoza modelelor și compararea rezultatelor obținute în urma modelării

În urma rezultatelor analizei exploratorii a datelor se pot stabili variabilele care explică cel mai bine modelul și care vor constitui baza pentru construirea unor modele performante. În alegerea modelelor trebuie să se țină cont de scopul analizei, de resursele necesare pentru implementare, pe ce perioadă de timp vor rămâne relevante rezultatele analizei și dacă va fi nevoie de întreținerea ulterioară a modelului pentru păstrarea relevanței, care este performanța modelului.

Baza procesului de modelare a datelor o constituie tehnicile de machine learning (învățarea automată), care sunt algoritmi care se pot autoperfecționa prin experiență. Algoritmii de machine learning construiesc modele bazate pe serii de date cu scopul de a elabora predicții pe baza cărora să se poată lua decizii, prin descoperirea relațiilor de cauzalitate sau de corelație între diverse variabile. Tehnicile de machine learning pot contribui la structurarea și identificarea relațiilor existente între anumite elemente ale unui sistem de informații și la descoperirea de noi structuri, noi relații și semnificații.

Alegerea și implementarea algoritmilor de machine learning utilizați pentru modelarea unui set de date trebuie să se facă în funcție de sursa, dimensiunea și tipul datelor utilizate, de sarcinile pe care trebuie să le îndeplinească și de performanța avută în vedere în scopul autoperfecționării prin experiență.

Algoritmii de machine learning utilizați în Data Science pot fi împărțiți în două tipuri, în funcție de modul de învățare, în algoritmi supervizați și nesupervizați.

**Algoritmii de învățare supervizată (*supervised learning*)** sunt algoritmi care construiesc modele matematice pe seturi de date în care rezultatele sunt cunoscute, iar prin instruirea pe

datele de antrenare, algoritmi învață să facă predicții pe seturi noi de date, care nu au făcut parte din datele de antrenament [11].

Exemple: algoritmi de regresie, de clasificare, K-nearest neighbors, decision tree, random forest, support vector machine, naïve Bayes, neural networks (multilayer perceptron), etc.

**Algoritmii de învățare nesupervizată (*unsupervised learning*)** sunt algoritmi care au la dispoziție doar date neetichetate, iar rezultatele nu sunt cunoscute [11]. Acești algoritmi încearcă în mod iterativ să extragă anumite reguli din datele disponibile și să organizeze datele pe baza identificării unor similitudini între date și pot avea ca obiectiv gruparea unor eșantioane cu comportament similar sau identificarea unor caracteristici similare în cadrul datelor.

Exemple: algoritmi de grupare (clustering), reducerea dimensionalității datelor (analiza în componente principale), învățarea regulilor de asociere, algoritmi de detectare a anomaliilor, etc.

Algoritmii de machine learning pot fi utilizați și în alte etape ale unui proiect, nu doar în partea de modelare. Tehnicile de modelare pot fi utilizate în pregătirea și analiza datelor, în completarea valorilor lipsă, în detectarea valorilor extreme, reducerea dimensionalității datelor, etc.

În procesul de modelare, tehnicile de machine learning pot fi utilizate pentru explicarea cauzelor unui fenomen (root cause analysis), pentru clasificarea unor fenomene sau pentru predicția unor fenomene.

Data science are la dispoziție diverse tehnici de modelare a datelor, de aceea este foarte important să se stabilească care sunt modelele adecvate pentru analiză, în funcție de tipul datelor și de scopul analizei. Asupra unui set de date pot fi aplicate în paralel mai multe tehnici de modelare, a căror performanțe trebuie evaluate și comparate, pentru a determina care dintre modele explică cel mai bine fenomenele studiate.

***Un model eficient trebuie să îndeplinească două condiții de bază:***

- să aibe putere predictivă asupra fenomenelor studiate
- să poată fi generalizat și extrapolat pe seturi noi de date

Pentru îndeplinirea acestor condiții este necesară implementarea următoarelor strategii:

- ***evaluarea modelelor*** prin definirea unor instrumente de măsurare a erorilor pe care pot să le genereze modelele, ca de exemplu rata erorii de clasificare, matricea de confuzie pentru clasificatori, sau coeficientul de determinare  $R^2$ , MSE, RMSE, MAE, pentru modelele de regresie;

- stabilirea unor *strategii de validare a modelului*, ca de exemplu divizarea setului de date în date de antrenare și date de test, utilizarea validării încrucișate, implementarea unor metode de regularizare a modelelor pentru tehnicile de regresie, în scopul obținerii unui compromis între bias și varianță, etc.

Validarea modelului este extrem de importantă, deoarece extrapolarea rezultatelor pe seturi noi de date determină dacă modelul poate funcționa în condiții reale și dacă este o reprezentare reală a fenomenelor studiate.

## **2.6. Prezentarea rezultatelor**

În urma modelării datelor și a obținerii rezultatelor cu privire la fenomenele analizate, rezultatele trebuiesc prezentate părților interesate și/sau pot constitui baza implementării unor aplicații.

Rezultatele și concluziile studiului pot fi prezentate ca atare, prin utilizarea unor instrumente de prezentare sau pot fi publicate în medii științifice de specialitate. Rezultatele pot fi automatizate și utilizate prin implementarea unor aplicații, iar întreg demersul de analiză poate fi operaționalizat pentru reutilizare permanentă și pentru integrarea cu alte instrumente [3].



### 3. Aspecte teoretice și practice ale unui proiect data science bazat pe analiza de regresie

Un proiect data science este un demers complex prin care date reale culese din diverse domenii sunt modelate prin tehnici de statistică, matematică, machine learning sau data mining. În funcție de tipul de date și de obiectivul studiului, datele pot fi analizate și modelate pentru regresie, clasificare, clustering, etc.

Analiza de regresie este un instrument prin care se pot analiza legăturile, conexiunile sau corelațiile existente între două sau mai multe fenomene. Pe baza rezultatelor analizei de regresie se pot lua decizii obiective, bazate pe date reale, observate prin măsurarea fenomenelor sociale, economice și din oricare alte domenii de activitate.

Analiza de regresie cuprinde tehinci statistice și matematice, precum și strategii de evaluare și validare a modelelor. Principalele concepte care stau la baza analizei de regresie și a strategiilor de validare sunt prezentate în prezentul capitol.

#### 3.1. Consideratii statistice

Statistica „*este știința care studiază fenomenele și procesele de tip colectiv din societate, natură, etc., din punct de vedere cantitativ, în scopul descrierii acestora și al descoperirii legilor care guvernează manifestarea lor*” [12].

Statistica este utilizată în data science în scopul pregătirii și analizării datelor, pentru a identifica anumite tendințe sau relații între datele analizate, pentru a înțelege datele analizate (statistica descriptivă) sau pentru a descoperi noi structuri, noi relații și semnificații despre problemele analizate (statistica inferențială). Datele sunt analizate prin intermediul unor indicatori statistici care arată tendința centrală sau variabilitatea lor.

**Indicatorii tendinței centrale** arată care este tendința evenimentelor analizate de a se grupa în jurul unui punct central. Principalii indicatori statistici care arată tendința centrală a unui set de date sunt media, mediana, modul [13].

**Media** unei serii de date statistice este raportul dintre suma valorilor individuale și numărul lor, și reprezintă măsura valorii centrale a setului de date în jurul căreia fluctuează datele. Se utilizează în general atunci când fenomenele cercetate prezintă o tendință liniară.

Acest indicator nu este reprezentativ dacă termenii seriei de date sunt prea dispersați, întrucât este sensibil față de valorile extreme [14].

$$\bar{x} = \frac{\sum x_i}{n}$$

Suma diferențelor dintre toate valorile individuale ale seriei și media setului de date este 0.

$$\sum_{i=1}^n (x_i - \bar{x}) = 0$$

**Mediana** este valoarea centrală din setul de date ordonat crescător sau descrescător, respectiv valoarea care împarte seria în două părți egale. Mediana nu este afectată de valorile extreme, dar poate fi nerepresentativă dacă valorile individuale nu se grupează în jurul valorii centrale [13] [14].

- pentru o serie cu număr impar de termeni mediana se calculează astfel:

$$Me = \frac{x_{n+1}}{2}$$

- pentru o serie cu număr par de termeni mediana se calculează astfel:

$$Me = \frac{1}{2}(x_{\frac{n}{2}} + x_{\frac{n}{2}+1})$$

**Modul** este valoarea care are cea mai mare frecvență în seria de date, fiind un indicator specific seriilor de distribuție unimodală (cu un singur maxim) și multimodală (cu mai multe maxime) [14].

**Indicatorii variației sau variabilității** arată pe ce domeniu se întinde setul de date, respectiv cât de dispersate sau de extinse sunt datele, având rol în inferența statistică (estimarea parametrilor pentru întreaga populație pornind de la rezultatele obținute la nivelul eșantionului). Se utilizează deoarece indicatorii tendinței centrale nu dau nicio indicație asupra împrăștierei datelor. Principalii indicatori ai variației sunt coeficientul de variație, deviația quartilă, deviația standard (abaterea standard), dispersia (varianța), covarianța [13].

**Coeficientul de variație** propus de Pearson este raportul dintre deviația standard și media seriei de date, fiind un indicator al împrăștierei datelor în raport cu tendința centrală.

$$Cv = \frac{\sigma}{\bar{x}} * 100$$

Acest indicator statistic se utilizează în special în studiul omogenității unei populații și pentru compararea a două serii de date atunci când au mediile mult diferite și abaterea standard nu lămurește diferențele. Coeficientul de variație se exprimă în procente și are valorile cuprinse între 0 și 100 %. Cu cât valoarea lui este mai apropiată de 0 cu atât seria statistică este mai omogenă și media este mai reprezentativă, iar cu cât valoarea se apropie de 100 cu atât seria este mai eterogenă iar media calculată este mai puțin reprezentativă. Pragul de la care se consideră că seria de date nu mai este omogenă este de 30-35% [13].

**Deviația quartilă** măsoară dispersia datelor aflate în zona de mijloc a distribuției și determină un interval centrat pe mediană.

$$DQ = \frac{(Q3 - Q1)}{2}$$

Se ordonează datele, după care se calculează valoarea mediană a jumătății superioare  $Q3 = UQ$  (upper quartile) și valoarea mediană a jumătății inferioare  $Q1 = LQ$  (lower quartile) [14].

Indicatorii de localizare se prezintă astfel:

- $Q1 \geq 25\%$  din datele seriei
- $Q2 =$  mediana (50% din date)
- $Q3 \geq 75\%$  din datele seriei

**Deviația standard (abaterea standard)** este un indicator al împrăștierei seriei de date în jurul valorii medii și se calculează ca medie pătratică sau ponderată a abaterilor valorilor seriei față de media lor, fiind de fapt radicalul dispersiei. O serie de date este omogenă dacă abaterea standard are o valoare mică [13].

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$$

**Dispersia (varianța)** este tot un indicator al gradului de împrăștiere a datelor și se calculează ca medie a pătratelor abaterilor termenilor seriei de la tendința lor centrală. Dispersia este abaterea standard ridicată la pătrat. O valoare mare a dispersiei arată o împrăștiere mare a valorilor seriei și o serie eterogenă, iar o valoare mică arată o împrăștiere mică a datelor, respectiv o serie omogenă.

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

Dacă comparăm două sau mai multe serii de date, la medii aproximativ egale, este mai împrăștiată seria cu dispersia mai mare, iar la dispersii aproximativ egale, este mai împrăștiată seria cu media mai mică [15].

**Covarianța** este măsura de variație comună a două variabile aleatorii și se calculează după următoarea formulă:

$$cov(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})$$

Dispersia (varianța) este un caz special al covarianței, în care cele două variabile sunt identice.

$$cov(x, x) = var(x)$$

Dacă valorile mari ale unei variabile corespund cu valorile mari ale altei variabile și în mod corespunzător acest fapt este valabil și pentru valorile mici, atunci avem două variabile cu comportamente similare, ceea ce înseamnă că avem covarianță pozitivă. Dacă acest lucru nu este valabil și avem valori mari la o variabilă și valori corespunzătoare mici la cealaltă variabilă, atunci cele două variabile au comportamente opuse și covarianța este negativă. Semnul covarianței arată practic direcția relației liniare care există între două variabile [16].

Totuși covarianța este greu de interpretat, deoarece nu este normalizată și depinde de interpretarea celor două variabile. Din acest motiv, se utilizează covarianța aplicată la date normalizate în prealabil sau varianta normalizată a covarianței, respectiv corelația, care poate arăta puterea relației liniare între două variabile.

**Corelația (coeficientul de corelație)** se utilizează pentru a defini interdependența sau legătura dintre două variabile. Acest indicator nu descrie o relație cauzală între două variabile și nu permite predicția unei variabile pe baza altei variabile. Coeficientul de corelație propus de Pearson se calculează astfel:

$$r_{x,y} = \frac{cov(x, y)}{\sqrt{var(x)} * \sqrt{var(y)}}$$

$$r_{x,y} = \frac{\sum((x_i - \bar{x}) * (y_i - \bar{y}))}{\sqrt{\sum(x_i - \bar{x})^2} * \sqrt{\sum(y_i - \bar{y})^2}}$$

Semnul valorii corelației indică direcția relației între două variabile:

- coeficient de corelație pozitiv = corelație directă, înseamnă că cele două variabile variază în același sens (dacă  $x$  crește - crește și  $y$ , dacă  $x$  scade - scade și  $y$ )

- coeficient de corelație negativ = corelație inversă, înseamnă că cele două variabile variază în sens contrar (daca  $x$  crește –  $y$  scade, daca  $x$  scade –  $y$  crește) [15]

Coeficientul de corelație  $r_{x,y}$  ia valori între  $[-1,1]$ , iar valoarea absolută a coeficientului constituie un indiciu al intensității corelației sau legăturii între cele două variabile  $x$  și  $y$ . Corelația este foarte stransă când  $r_{x,y} \rightarrow 1$  și foarte redusă când  $r_{x,y} \rightarrow 0$ .

Astfel, atunci când coeficientul de corelație ia valoarea -1 avem o corelație perfectă, inversă sau negativă, când ia valoarea 0 avem o corelație nulă sau inexistentă (asociere aleatoare) și când ia valoarea 1 avem o corelație perfectă, directă sau pozitivă.

### 3.2. Regresia liniară

Este un instrument de modelare statistică, utilizat pentru a determina existența unor legături sau corelații între date numerice. Modelul este reprezentat grafic printr-o linie care unește cel mai bine datele, care sunt reprezentate prin puncte [2]. Analiza de regresie se face pe baza unei ecuații de regresie, care exprimă evoluția variabilelor analizate și printr-un calcul al semnificativității utilizării acestei tehnici, care arată cât de adecvat este modelul pentru analizarea datelor. Pe baza datelor de antrenare, reprezentate de variabile numerice, se pot face predicții sau previziuni ale unor variabile pe baza altor variabile.

Regresia liniară identifică linia dreaptă care reprezintă cea mai bună aproximație dintr-o diagramă de dispersie a datelor [2]. Regresia liniară simplă utilizează două variabile pentru analiză, iar regresia multiplă utilizează mai multe variabile independente care afectează o variabilă dependentă.

Regresia liniară simplă face predicții ale unei variabile dependente  $y$ , pe baza unei variabile independente care este cunoscută  $x$ , prin utilizarea metodei celor mai mici pătrate.

**Ecuația funcției liniare de regresie** este următoarea:

$$y = a + bx + \varepsilon$$

Pentru regresia liniară multiplă ecuația funcției de regresie este următoarea:

$$y = a + b_1 * x_1 + b_2 * x_2 + \dots + b_i * x_i + \varepsilon$$

Interceptul și coeficientul de regresie sunt estimați prin rezolvarea următoarei probleme de optimizare [17] [18]:

$$\min_{\{a,b\}} F(a,b), \quad \text{unde} \quad F(a,b) = \sum_{i=1}^n \{y_i - (a + bx_i)\}^2$$

Minimul se obține prin calcularea sistemului de derivate parțiale:

$$\begin{cases} \frac{\partial F(a, b)}{\partial a} \big|_{a=\hat{a}, b=\hat{b}} = 0, \\ \frac{\partial F(a, b)}{\partial b} \big|_{a=\hat{a}, b=\hat{b}} = 0, \end{cases}$$

care se transformă în:

$$\begin{cases} -2 \sum_{i=1}^n (y_i - \hat{a} - \hat{b}x_i) = 0 \\ -2 \sum_{i=1}^n (y_i - \hat{a} - \hat{b}x_i)x_i = 0 \end{cases}$$

Ca urmare a rezolvării sistemului de ecuații, se obține formula de calcul a estimatorului interceptului și formula de calcul a estimatorului coeficientului de regresie.

**Interceptul** este locul de pe ordonată unde dreapta de regresie se intersectează cu axa Oy și reprezintă o constantă care este valoarea medie a lui  $y$  pentru  $x = 0$ , respectiv nivelul variabilei dependente  $y$  care nu este determinat de variabila independentă  $x$ , ci de alte variabile sau alți factori. Interceptul se calculează după următoarea formulă:

$$\hat{a} = \bar{y} - \hat{b} * \bar{x}$$

**Panta de regresie** arată variația medie a variabilei  $y$ , atunci când  $x$  crește cu o unitate, indicând faptul că între cele două variabile  $x$  și  $y$  există o legătură directă. Astfel la creșterea cu o unitate a variabilei independente  $x$ , variabila dependentă  $y$  crește cu valoarea reprezentată de coeficientul  $b$ . Dacă valoarea coeficientului de regresie este 0, înseamnă că nu există nicio relație între cele două variabile.

Întrucât regresia liniară este strâns legată de corelația dintre variabilele analizate, panta de regresie se calculează după formula:

$$\hat{b} = \frac{\sum_{i=1}^n ((x_i - \bar{x}) * (y_i - \bar{y}))}{\sum_{i=1}^n (x_i - \bar{x})^2} = r_{x,y} * \frac{\sigma_y}{\sigma_x}$$

În realitate variabila  $y$  nu depinde de o singură variabilă, ci de un ansamblu de variabile, exprimate în cadrul modelului prin variabila eroare sau reziduu  $\varepsilon$ , care însumează influența altor variabile, pe care modelul regresiei liniare nu le ia în calcul, asupra variabilei dependente  $y$ . Eroarea exprimă abaterile între valorile observate și valorile estimate prin aplicarea modelului.

Pentru estimarea celor doi parametri, constanta și panta de regresie, s-a utilizat metoda celor mai mici pătrate [19] [17] [20], care minimizează pătratele abaterilor dintre valorile date (datele observate) –  $y$  și cele calculate –  $\hat{y}$ , și care presupune parcurgerea mai multor etape:

- calcularea tuturor erorilor (diferențelor) între datele de antrenare (datele observate) și valorile rezultate (estimate) în urma aplicării modelului, care sunt reprezentate liniar

- valoarea fiecărei erori se ridică la pătrat și se calculează suma
- se găsește linia unde suma pătratelor erorilor sau reziduurilor ia valoarea cea mai mică

$$SSE = \sum (y_i - \hat{a} - \hat{b}x_i)^2 = \sum (y_i - \hat{y}_i)^2 = \text{minim}$$

Cu cât coeficientul de corelație între variabilele analizate este mai mare sau mai apropiat de valoarea absolută 1, cu atât valoarea erorii de estimare -  $\varepsilon$  va fi mai mică.

### 3.2.1. Ipoteze statistice asupra modelului regresiei liniare

Modelul regresiei liniare prin metoda celor mai mici pătrate trebuie să respecte anumite ipoteze de lucru cu privire la proprietățile estimatorilor. Prin validarea sau invalidarea ipotezelor, se poate determina dacă modelul liniar este adecvat pentru explicarea fenomenelor studiate și se poate evalua precizia estimatorilor.

**Liniaritatea modelului** – între variabila dependentă și variabila independentă există o relație liniară.

- este esențială pentru validarea coeficienților estimați ai ecuației de regresie
- verificarea liniarității se poate face grafic prin intermediul scatterplots, respectiv dacă norul de puncte care arată relația dintre variabila predictor și cea estimată este dispus în model liniar, atunci există o relație liniară între cele două variabile
- dacă relația între variabile nu este liniară, atunci se pot face transformări de liniarizare asupra variabilelor (modelul exponențial, modelul de tip putere) [20]
- în cazul în care relația dintre variabile nu se poate liniariza, se folosesc modele de regresie neliniare, de exemplu modelele polinomiale, care au la baza funcții neliniare

**Inexistența multicolinearității** – variabilele predictor din model nu sunt corelate între ele (testarea se face doar pentru regresia multiplă)

- pentru identificarea existenței corelațiilor între predictor/variabilele independente se afișează matricea de corelație, calculată pe baza coeficientului de corelație Pearson; dacă între două variabile predictor există o corelație mare, atunci una dintre ele trebuie eliminată din model; se păstrează variabila care are o corelație mai mare cu variabila dependentă  $y$  [21]

**Normalitatea erorilor** – variabila  $\varepsilon$  urmează o lege normală de medie 0 și variație  $\sigma^2$ :  $\varepsilon_i \sim N(0, \sigma^2)$ ; implică și verificarea ipotezei că **media erorilor este 0**:  $\bar{\varepsilon} = 0$

- dacă ipoteza e confirmată, atunci înseamnă că și estimatorii parametrilor modelului de regresie urmează o lege normală
- dacă ipoteza e încălcată, înseamnă că estimatorii parametrilor nu sunt eficienți
- testarea ipotezei de normalitate a erorilor se poate face grafic (histograma/diagrama reziduurilor) sau prin procedee numerice (testul Kolmogorov-Smirnov, testul Jarque-Bera); prin intermediul diagramei reziduurilor se poate observa foarte ușor dacă erorile au o distribuție normală [20] [22]

**Homoscedasticitatea (omogenitatea varianței erorilor)** – varianțele erorilor sunt constante oricare sunt valorile variabilei predictor  $x$ :  $v(\varepsilon_i) = \sigma^2$ ;

- ipoteza se poate confirma prin metode grafice sau non-grafice
- **metode grafice**: prin reprezentarea grafică (scatterplot) a relației dintre reziduuri și variabila  $x$ , sau reziduri și variabila estimată  $\hat{y}$ , în cazul regresiei multiple, se poate observa dacă norul de puncte se află dispus normal în jurul liniei de regresie, fără a forma un anumit model; dacă se identifică un model definit (liniar sau pătratic sau în formă de pâlnie), atunci avem heteroscedasticitate
- **metode non-grafice**: testul Park, testul Glejser, testul White, testul Goldfield-Quandt, testul t Student pentru coeficientul de corelație neparametrică Spearman [20] [22]

**Necorelarea erorilor** – erorile sunt necorelate între ele:  $cov(\varepsilon_i, \varepsilon_j) = 0$ ; erorile asociate unor valori ale variabilei  $y$  nu sunt influențate de erorile asociate altor valori ale variabilei  $y$

- este aplicabilă în special în **analiza time-series**, între valorile aceleiași variabile observată la diferite momente în timp; cea mai simplă metodă de verificare a ipotezei este prin intermediul unui scatterplot a valorilor reziduale vs. evoluția în timp a variabilei observate (mai întâi se aranjează datele în ordine temporală) [20]
- în analiza de regresie autocorelarea erorilor intervine dacă modelul e incorect specificat, respectiv dacă se încearcă modelarea unei relații liniare asupra unor date neliniare, atunci reziduurile vor fi autocorelate
- pentru testarea ipotezei se pot utiliza metode grafice (scatterplot) sau non-grafice: testul Durbin Watson, Runs test, testul Ljungbox [20] [23]



### 3.2.2. Predicția punctuală

Pe baza modelului de regresie liniară aplicat se poate face o previziune punctuală a comportamentului variabilei  $y$ , în funcție de valorile fixe pe care le ia variabila independentă  $x$ . Rezultatele obținute se referă la un comportament mediu al variabilei  $y$  [19].

### 3.2.3. Estimarea parametrilor pe bază de interval de încredere

Predicția punctuală a unui parametru poate fi la o anumită distanță față de valoarea reală a parametrului estimat. Prin urmare estimarea parametrului se poate face și pe baza unui interval de predicție, între limitele căruia poate să ia valori parametrul estimat, în cazul nostru variabila estimată  $\hat{y}$ .

Estimarea pe interval de predicție se face pe baza unui **interval de încredere** în cadrul căruia se găsește cu probabilitate mare parametrul estimat, modul de calcul fiind de forma: estimare  $\pm$  eroare, astfel că estimatorul să se situeze între valoarea inferioară a intervalului și valoarea superioară a acestuia.

$$[\text{valoarea inferioară} < \text{estimator} < \text{valoarea superioară}] = 1-\alpha,$$

unde  $\alpha$  este nivelul de semnificație, iar intervalul definit de cele două valori, cea inferioară și cea superioară va cuprinde estimatorul populației cu o probabilitate de  $1-\alpha$  [19].

$$[\hat{y} - t_{1-\frac{\alpha}{2}} * \sigma_{\hat{y}} ; \hat{y} + t_{1-\frac{\alpha}{2}} * \sigma_{\hat{y}}]$$

Calculul intervalului de predicție se face pe baza următorilor indicatori [19]:

- variația erorilor de predicție/reziduurilor, care se calculează ca și rădăcina pătrată a diferențelor între valorile lui  $y$  și valorile estimate a lui  $y$  prin aplicarea modelului (după legea Student s-au luat  $n-2$  grade de libertate)

$$\sigma_{\varepsilon} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n - 2}}$$

- variația estimatorului  $\hat{y}$  calculat în funcție de  $x$  după formula:

$$\sigma_{\hat{y}} = \sigma_{\varepsilon} \sqrt{1 + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

- gradul de încredere al intervalului specificat – scorul  $z$ , care corespunde nivelului de încredere; valoarea critică  $z$  se găsește din tabelul de scoruri  $z$ , întrucât se cunoaște valoarea abaterii standard a populației și se poate presupune că valorile sunt distribuite

normal; (valorile z comune sunt 1,645 pentru un nivel de încredere de 90%, 1,960 pentru un nivel de încredere de 95% și 2,576 pentru un nivel de încredere de 99%).

În funcție de indicatorii prezentați se calculează **intervalul de predicție** cu un nivel de încredere de  $(1-\alpha)$  [19] după următoarea formulă:

$$\hat{y}_i \pm t_{1-\frac{\alpha}{2}} * \sigma_{\varepsilon} \sqrt{1 + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

Spre deosebire de intervalul de predicție, **intervalul de încredere** pentru valorile estimate pe baza modelului, se calculează după formula [19]:

$$\hat{y}_i \pm t_{1-\frac{\alpha}{2}} * \sigma_{\varepsilon} \sqrt{\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

### 3.3. Regresia RANSAC (random simple consensus)

Întrucât regresia liniară este un model conceput pentru a se adapta la toate punctele care reprezintă setul de date, este puternic influențată de elementele extreme prezente în setul de date (erori de măsurare), astfel că acestea trebuie identificate și eliminate prin metode statistice de analiză a datelor. Metoda celor mai mici pătrate este un model adecvat dacă ipotezele de la care pornește analiza sunt adevărate și setul de date nu este afectat de valori extreme (observații care nu urmează tiparul observațiilor majoritare din setul de date).

Statistica robustă (solidă) este implementată pentru a depăși limitările metodelor tradiționale de analiză statistică, regresia robustă fiind o metodă statistică concepută pentru modelări care să nu fie afectate de valori extreme sau aberante [24].

Estimatorul RANSAC este un model statistic robust, care rămâne relevant în reprezentarea unui set de date, întrucât nu este influențat de anomalii (outliers), deoarece utilizează în modelarea liniară doar elementele din mulțimea de consens (inliers), care sunt determinate în mod iterativ prin eșantionări aleatorii de date. RANSAC este o metodă iterativă de detectare și de eliminare a anomaliilor dintr-un set de date, care nu încearcă să adapteze elementele extreme, ci să elimine valorile care nu sunt distribuite normal [25].

Pachetul sklearn pentru Python facilitează implementarea estimatorului RANSAC prin RANSACRegressor, care are ca estimator de bază pentru antrenarea datelor regresia liniară. RANSACRegressor utilizează abaterea absolută mediană (MAD) a variabilei dependente  $y$  pentru detectarea anomaliilor, întrucât față de medie, mediana este un indicator statistic robust [26]. Abaterea absolută mediană este un indicator al deviației de la valoarea mediană a setului

de date și este utilizat pentru clasificarea datelor în elemente interioare care satisfac modelul (mulțimea de consens/inliers) și date extreme (outliers).

Se alege pragul de distanță (residual\_threshold), care va determina seturi diferite ale mulțimii de consens (inliers), în funcție de mărimea pragului aleasă în mod iterativ. Rezultatele diverselor iterații sunt vizualizate într-un grafic scatter, care adaptează dreapta de regresie în funcție de valoarea pragului de distanță.

### 3.4. Regresia Decision Tree și Regresia Random Forest

Random forest este un model de învățare de ansamblu (ensemble learning). Învățarea de ansamblu presupune aplicarea mai multor algoritmi de machine learning sau a aceluiași algoritm de mai multe ori, pentru a face predicții mai precise asupra unui set de date. Algoritmul random forest este un model de învățare supervizată, care poate fi utilizat pentru clasificare sau regresie.

Random forest construiește în paralel o multitudine de arbori de decizie (decision trees) pe datele de antrenare și combină rezultatul predicțiilor tuturor arborilor de decizie pentru a elabora rezultatul final [27] [28].

*Arborele de decizie*, care este un model non-liniar, construiește modele de clasificare pentru date discrete, sau modele de regresie pentru valori numerice continue, în forma unui arbore.

Un arbore de decizie este un discriminator de categorii, care divide/ramifică în mod recursiv datele de antrenament prin intermediul unor noduri interne, până se obțin nodurile de decizie/frunzele, care reprezintă rezultatele finale:

- *nodul rădăcină* este nodul origine de unde se ramifică toate ramurile și care nu are intrare, doar ieșiri;
- *nodurile de decizie* sunt noduri interne, care au o intrare de la nodul rădăcină sau de la alte noduri interne și au câte două ramificații; fiecare nod reprezintă un punct de partiționare în funcție de un atribut (testarea după un anumit atribut) care determină modul de divizare al nodului respectiv;
- fiecare *ramură* reprezintă rezultatul testului;
- *frunzele*, care au câte o intrare de la nodul rădăcină sau de la noduri interne și nu au nicio ramificație, sunt noduri finale și reprezintă clase/categorii obținute după

calcularea tuturor atributelor; conținutul nodurilor frunze reprezintă rezultatele modelului.

În cazul în care avem *o singură variabilă predictor*  $x$ , algoritmul urmează următoarele etape:

Intervalul de valori al variabilei predictor  $x$  (axa  $Ox$ ) este împărțit în subseturi mai mici (regiuni sau segmente), delimitate de praguri, notate cu  $R_i$ . Pentru fiecare regiune se calculează media valorilor lui  $y$  cuprinse în regiunea respectivă, care reprezintă valoarea estimată a variabilei  $y$ , respectiv  $\hat{y}$ .

$$\hat{y} = \text{media valorilor variabilei } y$$

Astfel, pentru fiecare segment în parte, variația reziduală a erorilor este considerată a fi suma pătratelor diferenței dintre valorile lui  $y$  și media lui  $y$ . În cadrul fiecărui segment se urmărește minimizarea sumei pătratelor erorilor  $SSE$ .

Se calculează suma pătratelor erorilor pentru fiecare segment/regiune și se găsește segmentul care are cea mai mică valoare a  $SSE$ . Acest segment devine nodul rădăcină din care se ramifică toate celelalte noduri de decizie (ramificare binară < decât valoarea prag și respectiv  $\geq$  decât valoarea prag) [29].

$$\text{nodul rădăcina} = \min\{SSE_i \mid i \in R_i\} = \min\left\{\sum (y_j - \hat{y}_j)^2 \mid i \in R_i\right\}$$

Algoritmul se repetă recursiv până la nivelul nodurilor finale (frunzelor), care nu mai pot fi divizate binar în grupuri/categorii mai mici. Fiecare frunză reprezintă media variabilei  $y$ , a unui segment (cluster) din totalul observațiilor.

Algoritmul se va opri atunci când în segmentul  $R$  rămâne un singur punct care nu mai poate fi divizat. Se poate stabili iterativ numărul maxim de divizări sau ca nodurile să se dividă până când în fiecare nod sau regiune se regăsește un număr minim de observații, sub care nodurile nu se mai pot diviza. Numărul de niveluri ale nodurilor de decizie și ale nodurilor frunze reprezintă adâncimea arborelui.

Dacă adâncimea arborelui e mare sau dacă pragul de divizare stabilit este mic (numărul minim de observații dintr-o regiune până la care se poate diviza nodul), atunci algoritmul se ajustează mai bine pe datele de antrenare, dar generează overfitting și o variație mare pe datele de test. Reducerea overfittingului se poate face prin stabilirea unei adâncimi mai mici a arborelui sau a unor valori mai mari pentru pragul de divizare.

În cazul în care avem *mai multe variabile predictor*, pașii de mai sus se vor relua pentru fiecare variabilă. Se va calcula pentru fiecare variabilă predictor segmentul cu cea mai mică valoare a  $SSE$ , și pentru fiecare variabilă atribut acest segment va deveni candidatul pentru a fi nodul rădăcină al arborelui [29].

Se compară candidații, iar atributul care vine cu cea mai mică valoare a  $SSE$  este atributul câștigător. Segmentul prezentat de atributul câștigător, devine nodul rădăcină al arborelui, care va fi construit în continuare prin compararea celor mai mici valori ale  $SSE$  pentru fiecare predictor și repetând procedura pentru observațiile rămase, până când acestea nu mai pot fi divizate în grupuri mai mici [29].

Principalele dezavantaje ale modelului arborelui de decizie sunt următoarele:

- modelul este sensibil la valorile extreme (outliers)
- se ajustează foarte bine pe datele de antrenament, dar nu performează la fel de bine pe datele de test, având un bias mic pe datele de antrenare, dar o varianță mare pe datele de test
- modelul rezultat nu poate fi extrapolat pentru valori aflate în afara minimului sau maximului valorilor observate

Pentru evitarea overfitting-ului pe datele de antrenare se poate utiliza modelul *random forest*, care construiește în paralel o multitudine de arbori de decizie pe datele de antrenament, reducând proprietatea de overfitting a algoritmului arborilor de decizie. Random forest construiește mai mulți arbori, care sunt instruiți pe diverse subseturi ale datelor de antrenament, în scopul reducerii varianței pe datele de test. Cu cât crește numărul arborilor de decizie construiți, cu atât scade varianța. Astfel modelul pierde din interpretabilitatea pe datele de antrenare, dar are o performanță îmbunătățită pe datele de test.

Random forest folosește tehnica bagging (bootstrap aggregation), care stabilește în mod aleator, cu înlocuire, subseturi ale setului de date de antrenare, fiecare subset de date fiind utilizat pentru antrenarea unui arbore de decizie asociat [27] [29]. Numărul arborilor de decizie care vor fi construiți se stabilește în mod iterativ.

Pentru fiecare observație nouă dintr-un set de test, fiecare dintre arborii de decizie va face o predicție cu privire la valoarea lui  $y$ . Valoarea finală a lui  $y$  estimată de algoritm va fi reprezentată de media rezultatelor tuturor arborilor individuali [29].

Algoritmul random forest este foarte performant în analiza de regresie și are multiple avantaje întrucât este foarte stabil, lucrează atât cu variabile numerice cât și cu variabile categoricale, nu este afectat de variabile nescalate și nu este afectat de bias.

Principalul dezavantaj al algoritmului este dat de însăși complexitatea sa, întrucât necesită resurse computaționale și resurse de timp pentru antrenare mai mari față de alți algoritmi. Un alt dezavantaj este faptul că modelul rezultat nu poate fi extrapolat pentru valori aflate în afara minimului sau maximului valorilor observate [30].

### **3.5. Evaluarea performanței algoritmilor de machine learning**

În modelarea datelor sunt utilizate tehnici de învățare automată supervizată (supervised learning), întrucât setul de date conține date de intrare, iar rezultatele sunt cunoscute. Astfel sunt implementați algoritmi care printr-un proces de instruire pe baza datelor existente, să facă predicții, care să fie îmbunătățite ca urmare a analizei performanței aplicării modelelor respective. Procesul de instruire continuă până când se atinge un nivel de precizie optim. Datele de intrare sunt disponibile sub formă de vectori sau matrici, iar prin aplicarea unor funcții matematice asupra datelor de antrenare, algoritmi învață să facă predicții pe seturi noi de date, care nu au făcut parte din datele de antrenament.

În modelarea datelor însă nu este suficient să aplicăm diverse tehnici de machine learning, ci trebuie să și evaluăm performanța algoritmilor implementați, atât pentru a vedea care dintre modele este cel mai adecvat pentru prelucrarea datelor, cât și pentru a valida rezultatele obținute, prin compararea eficienței modelelor aplicate.

În scopul evaluării performanței modelelor, se pot utiliza tehnici de evaluare a modelelor și strategii de validare a modelelor.

Pentru atingerea acestor obiective, pot fi utilizate următoarele metode:

- determinarea semnificativității modelului în explicarea relației între variabile
- măsurarea erorilor prin diverse tehnici statistice
- implementarea mai multor algoritmi de machine learning și compararea lor performanțelor lor
- împărțirea datelor în seturi diferite pentru date de antrenare și date de testare
- validarea încrucișată (k-fold cross-validation)
- compensarea bias-varianță prin strategii de regularizare a modelelor de regresie liniară

## 3.6. Metode de evaluare a performanței modelelor

### 3.6.1. Coeficientul de determinare

Coeficientul de determinare arată care este proporția sau procentul din variația variabilei  $y$  care este explicată de variabila  $x$ , arătând care este proporția în care relația de variabilitate dintre  $x$  și  $y$  este explicată de modelul de regresie. Coeficientul de determinare este esențial în determinarea performanței modelelor liniare.

În cazul regresiei liniare simple, coeficientul de determinare este coeficientul de corelație Pearson la pătrat, prin urmare testarea coeficientului de determinare are sens doar dacă coeficientul de corelație este semnificativ:

$$R^2 = r_{x,y}^2$$

Pentru ca funcția de regresie aleasă să fie semnificativă, trebuie ca variația reziduală să fie minimă. Prin urmare, coeficientul de determinare se calculează pe baza următorilor indicatori [19] [17] [18]:

- **SSR** – variația totală a datelor  $\hat{y}$  estimate sau suma pătratelor abaterilor de regresie, unde  $\hat{y}$  sunt valorile estimate și  $\bar{y}$  este media datelor de intrare:

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

- **SSE** – variația între datele de intrare și cele estimate sau suma pătratelor valorilor reziduale, unde  $y$  sunt datele de intrare, iar  $\hat{y}$  sunt valorile estimate prin aplicarea modelului:

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

- **SST** – variația totală a datelor  $y$  de intrare sau suma pătratelor abaterilor individuale de la medie, unde  $y$  sunt datele de intrare, iar  $\bar{y}$  este media datelor de intrare:

$$SST(\text{variația totală}) = SSR(\text{variația explicată}) + SSE(\text{variația neexplicată})$$

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

Coeficientul de determinare se calculează după următoarea formulă:

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

Coeficientul de determinare ia valori între 0 și 1 și se interpretează procentual:

- Dacă  $SSE = 0$  și  $R^2 = 1$  înseamnă că modelul regresiei e perfect, variațiile lui  $y$  sunt explicate complet de  $x$  și deci între cele două variabile există o legătură liniară;
- Dacă  $SSE = 1$  și  $R^2 = 0$  înseamnă că modelul regresiei nu explică o relație între variabile, respectiv între cele două variabile nu există o legătură liniară.

Complementul lui  $R^2$ , respectiv  $(1 - R^2)$  se numește coeficient de nedeterminare și arată proporția în care variabila  $y$  nu este explicată de variabila  $x$ , ci de alți factori care nu sunt luați în considerare de model.

### 3.6.2. Coeficientul de determinare ajustat

Coeficientul de determinare ajustat este esențial în evaluarea modelelor de regresie liniară multiplă. Coeficientul de determinare  $R^2$  nu poate arăta dacă predicțiile sunt influențate de bias. Dacă un model are prea mulți predictor, este posibil să fie supraestimat (overfitting) și să își piardă abilitatea de a face predicții pe seturi noi de date.

În cazul regresiei multiple, întrucât adăugând modelului noi variabile predictor, valoarea  $SSE$  va scade, implicit valoarea  $R^2$  va crește. Coeficientul de determinare  $R^2$  crește cu fiecare variabilă predictor adăugată modelului, chiar dacă unele variabile nu sunt esențiale în explicarea modelului.

Coeficientul de determinare ajustat  $R_{adj}^2$  este corectat cu gradele de libertate, astfel ca adăugând noi variabile neesențiale în explicarea modelului, acest indicator va scade, aducând o penalizare pentru utilizarea de variabile inutile. Dacă se adaugă modelului variabile utile, atunci valoarea  $R_{adj}^2$  va crește [31].

Coeficientul  $R_{adj}^2$  va fi întotdeauna mai mic sau egal cu  $R^2$  și se calculează după următoarea formulă:

$$R_{adj}^2 = 1 - \left[ \frac{(1 - R^2) * (n - 1)}{n - k - 1} \right]$$

Coeficientul  $R_{adj}^2$  permite compararea puterii de predictibilitate a modelelor de regresie care conțin numere diferite de predictor, deoarece arată dacă un număr mai mare de predictor explică mai bine modelul sau nu. Valoarea lui e întotdeauna mai mică decât valoarea lui  $R^2$  [31].



### 3.6.3. Eroarea medie patrată (mean squared error)

- reprezintă media sumei pătratelor diferențelor între valorile originale ale variabilei  $y$  și valorile estimate  $\hat{y}$  [32] [33] [34]
- măsoară variația erorilor
- un scor bun al  $MSE$  tinde spre zero, dar evaluarea depinde de specificul setului de date
- este sensibilă la valori extreme (outliers)

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

### 3.6.4. Eroarea rădăcinii medie pătratică (root mean squared error)

- este rădăcina pătrată a erorii medii pătratice [32] [33] [34]
- măsoară abaterea/deviația standard a erorilor
- valorile sunt scalate la nivelul variabilei estimate  $\hat{y}$ , prin urmare vom avea erori cu aceeași unitate de măsură ca și datele originale
- un scor bun al  $RMSE$  tinde spre zero, dar evaluarea depinde de specificul setului de date

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

### 3.6.5. Media erorilor absolute (mean absolute error)

- reprezintă media în valori absolute a diferenței între valorile actuale a variabilei  $y$  și valorile estimate  $\hat{y}$  [32] [33] [34]
- măsoară media valorilor reziduale
- prin comparare cu media variabilei  $y$  putem afla cât la sută din medie este  $MAE$ ; cu cât procentul e mai mic, cu atât modelul e mai precis
- are unități de măsură similare cu datele originale, deci poate fi comparată doar între modele ale căror erori sunt măsurate în scală similară
- nu este sensibilă la valori extreme (outliers), deci este mai robustă decât  $MSE$

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

### 3.6.6. Graficul rezidurilor din regresie

Graficul valorilor reziduale este un instrument foarte util de vizualizare a tiparelor variației rezidurilor și se construiește luând pe axa  $ox$  variabila  $x$ , iar pe axa  $oy$  valorile reziduale, care se calculează ca diferență între valorile actuale ale variabilei  $y$  și valorile estimate ale acesteia prin aplicarea modelului [19] [20].

Din forma norului de puncte se poate observa distribuția valorilor reziduale care validează ipoteza modelului regresiei liniare, respectiv:

- dacă există un pattern al acestora înseamnă că există o corelație între variabila  $x$  și reziduri, și relația nu poate fi modelată prin regresia liniară
- dacă sunt dispuse la întâmplare fără să formeze un anumit model, înseamnă că variabila  $x$  și reziduurile sunt independente, și relația poate fi modelată prin regresia liniară

### 3.6.7. Testarea modelului de regresie pe baza statisticii $t$ Student

**Testul  $t$  Student** se folosește în regresia liniară pentru testarea semnificativității statistice a coeficientului de regresie/panta de regresie, în determinarea existenței unei legături liniare între variabila predictor  $x$  și variabila dependentă  $y$ .

În acest sens, se formulează **ipoteza nulă  $H_0$** , prin care se consideră că variabila  $y$  nu este influențată de variația variabilei  $x$  și deci coeficientul  $b$  din ecuația de regresie nu este semnificativ diferit de zero. **Ipoteza alternativă  $H_1$**  presupune că variabila  $y$  este influențată de variabila  $x$  și prin urmare valoarea coeficienților de regresie este semnificativ diferită de zero [19] [20].

$$H_0: b = 0$$

$$H_1: b \neq 0$$

În cazul regresiei multiple:

$$H_0: b_1, b_2, \dots, b_m = 0$$

$$H_1: \exists j = \overline{1, i}, b_j \neq 0$$

Dacă se respinge ipoteza  $H_0$ , cu un prag de semnificație  $\alpha$  ales, înseamnă că între variabila/variabilele independente și variabila dependentă există o legătură semnificativă. Pragul de semnificație arată probabilitatea de a obține datele observate pornind de la premisa că ipoteza nulă este adevărată. În practică se consideră de regulă un prag de semnificație  $\alpha=0.05$ , adică un risc de 5% de a se respinge în mod eronat ipoteza  $H_0$ , atunci când este adevărată.

Pragul de semnificație de 0.05 în testarea ipotezei nule, a fost stipulat de Fisher, părintele statisticii moderne.

Verificarea ipotezei  $H_0$  se face cu ajutorul testului  $t$  Student [19] [18] [20], după formula:

$$t = \frac{\hat{b} - b_0}{SE_{\hat{b}}} \sim T_{n-2}, \quad t \text{ este o statistică Student cu } (n-2) \text{ grade de libertate}$$

unde:  $\hat{b}$  este coeficientul de regresie

$b_0$  este valoarea asumată 0 conform ipotezei  $H_0$

$SE_{\hat{b}}$  este abaterea standard a coeficientului de regresie

iar abaterea standard a coeficientului de regresie este:

$$SE_{\hat{b}} = \frac{\sqrt{\frac{1}{n-2} * \sum_{i=1}^n (y_i - \hat{y}_i)^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

astfel:

$$t = \frac{(\hat{b} - b_0) * \sqrt{n-2}}{\sqrt{\frac{SSE}{\sum_{i=1}^n (x_i - \bar{x})^2}}} = \frac{(\hat{b} - b_0) * \sqrt{n-2}}{\sqrt{1 - \hat{b}^2}}$$

O altă metodă de determinare a lui  $t$  este pe baza coeficientului de corelație a lui Pearson [20], după formula:

$$t = \frac{r_{x,y} * \sqrt{n-2}}{\sqrt{1 - R^2}}$$

Valoarea calculată a lui  $t$  se compară cu valoarea teoretică a lui  $t$ , obținută din tabelul Student, pentru  $(n-2)$  grade de libertate și pentru nivelul de semnificație  $\alpha$  stabilit. Pentru aflarea valorii teoretice se utilizează un risc de  $\alpha/2$ , întrucât distribuția Student este simetrică, iar suprafața de respingere  $\alpha$  este împărțită în două părți egale [19] [20].

Dacă valoarea absolută a testului  $t$  este mai mare decât valoarea teoretică a testului, se respinge ipoteza nulă, coeficientul de regresie fiind considerat semnificativ pentru explicarea modelului.

### 3.6.8. Testarea modelului de regresie pe baza statisticii test $F$

Testul  $F$  Fisher este un test global de semnificație a ansamblului coeficienților de regresie. Spre deosebire de testul  $t$  Student, care se poate aplica doar pentru un coeficient de regresie

odată, testul  $F$  se poate aplica pentru coeficienți multipli de regresie simultan, ceea ce permite compararea între modele liniare cu număr diferit de coeficienți de regresie.

Testul  $F$  compară modelul care se implementează, cu un model fără variabile independente, deci care are doar intercept, fără a avea coeficienți de regresie și verifică dacă coeficienții de regresie îmbunătățesc modelul.

Modelul redus bazat doar pe intercept, sugerează că fiecare valoare a variabilei  $y$  este o funcție compusă dintr-o medie generală reprezentată de intercept (media lui  $y$ ) și o eroare  $\varepsilon$ :

Prin urmare ipotezele sunt următoarele [35] [36]:

**$H_0: b = 0$  (ipoteza nulă)** - modelul fără variabile independente este la fel de semnificativ ca și modelul implementat cu variabile independente:

$$y_i = a + \varepsilon_i$$

**$H_1: b \neq 0$  (ipoteza alternativă)** - modelul implementat cu variabile independente este mai semnificativ decât modelul bazat doar pe intercept: (ipoteza alternativă)

$$y_i = a + b_1 * x_1 + \dots b_i * x_i + \varepsilon_i$$

Testul  $F$  se poate calcula pe baza tabelului ANOVA din Figura 3.1, ca raport între varianța explicată și varianța neexplicată [19] [20] [17]:

Sursa variației	Suma pătratelor	Grade de libertate	Media pătratelor	F
<b>Regresie</b>	$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	$k - 1$	$MSR = \frac{SSR}{k - 1}$	$\frac{MSR}{MSE}$
<b>Eroare</b>	$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2$	$n - k$	$MSE = \frac{SSE}{n - k}$	
<b>Total</b>	$SST = \sum_{i=1}^n (y_i - \bar{y})^2$	$n - 1$	$\frac{SST}{n - 1}$	

**Figura 3.1: Tabelul ANOVA**

$$F = \frac{\text{variata explicata}}{\text{variata neexplicata}} = \frac{SSR}{SSE} * \frac{n - k}{k - 1} \sim F_{k-1, n-k}$$

Valoarea  $F$  obținută se compară cu valoarea teoretică (valoarea critică) din tabelul F, care se citește în funcție de: nivelul de semnificație  $\alpha$  stabilit și de gradele de libertate ( $k-1$ ,  $n-k$ ). Rezultatul este semnificativ dacă valoarea testului  $F$  este mai mare decât valoarea teoretică citită din tabelul F, caz în care se respinge ipoteza nulă și se acceptă ipoteza alternativă, considerându-

se că modelul cu coeficient de regresie este mai semnificativ decât modelul bazat doar pe intercept.

### **3.7. Strategii de validare a modelelor**

#### **3.7.1. Validarea simplă (seturi de date disjuncte de antrenare și testare)**

În scopul eliminării supraestimării (overfitting), datele pot fi împărțite în date de antrenare și date de testare. Nu este eficient să antrenăm algoritmi pe un set de date și să utilizăm exact același set de date pentru a evalua algoritmi. Dacă utilizăm un algoritm pe un set de date pentru antrenare și utilizăm același set pentru testare, atunci algoritmul va avea un scor perfect pe setul de date de antrenament, dar nu știm cum se va comporta modelul aplicat altor seturi de date, care nu apar în setul de antrenament. Este foarte posibil ca performanța aplicării aceluiași model pe alte seturi de date să fie foarte slabă. Din acest motiv, este indicat să se facă testări ale modelului implementat pe un set de date de testare, diferit de setul utilizat pentru datele de antrenare [37].

Prin compararea rezultatelor și a scorului obținut prin aplicarea algoritmului asupra datelor de antrenament, cu cele obținute prin aplicarea modelului asupra datelor de testare, se poate estima dacă modelul nu este afectat de overfitting, dacă nu este afectat de variabile neluate în calcul (de ex. date extreme – outliers – care pot altera rezultatele) și deci dacă modelul este eficient pentru utilizare operațională.

Pe baza datelor de antrenare se estimează parametrii funcțiilor de regresie, iar pe baza datelor de testare se evaluează modelele implementate.

#### **3.7.2. Validarea încrucișată (k-fold cross validation)**

În scopul evaluării performanței modelelor implementate se poate folosi validarea încrucișată, care efectuează mai multe iterații pe același set de date, prin împărțirea setului în sub-seturi egale de date, după cum urmează [37] [38]:

- setul de date este împărțit în  $k$  părți sau grupuri egale (de aici denumirea de  $k$ -fold),  $k-1$  sub-seturi pentru antrenare și 1 set pentru testare (dimensiunea unui sub-set = dimensiunea setului /  $k$ )
- algoritmul este antrenat pe  $k-1$  părți și testat pe o parte din cele  $k$  stabilite
- se repetă antrenarea și testarea de  $k$  ori, pentru toate variantele de împărțire a grupurilor

Alegerea numărului  $k$  se face în funcție de mărimea setului de date, astfel că dimensiunea fiecărui grup sau partiție să fie suficient de mare pentru a fi reprezentativă.

Un exemplu de validare încrucișată pentru  $k = 5$  este prezentat în Figura 3.2:

Împărțirea 1	Grupul 1 testare	Grupul 2 antrenare	Grupul 3 antrenare	Grupul 4 antrenare	Grupul 5 antrenare	=>	Scorul 1
Împărțirea 2	Grupul 1 antrenare	Grupul 2 testare	Grupul 3 antrenare	Grupul 4 antrenare	Grupul 5 antrenare	=>	Scorul 2
Împărțirea 3	Grupul 1 antrenare	Grupul 2 antrenare	Grupul 3 testare	Grupul 4 antrenare	Grupul 5 antrenare	=>	Scorul 3
Împărțirea 4	Grupul 1 antrenare	Grupul 2 antrenare	Grupul 3 antrenare	Grupul 4 testare	Grupul 5 antrenare	=>	Scorul 4
Împărțirea 5	Grupul 1 antrenare	Grupul 2 antrenare	Grupul 3 antrenare	Grupul 4 antrenare	Grupul 5 testare	=>	Scorul 5

**Figura 3.2: Tabel validare încrucișată**

Ca urmare a executării tuturor iterațiilor, fiecare grup de date va fi utilizat de  $k-1$  ori pentru antrenare și o dată pentru testare.

Validarea încrucișată aplicată diversilor algoritmi de machine learning utilizați, permite compararea modelelor pentru a se determina eficiența aplicării acestora asupra setului de date, respectiv performanța algoritmilor.

### 3.7.3. Compensare bias - varianță

În evaluarea eficienței unui model de machine learning este importantă obținerea unui compromis între bias-ul și variația parametrilor estimați, compromis care contribuie la implementarea de modele adecvate pentru datele analizate, precum și la evitarea greșelilor de overfitting (supraestimare) și underfitting (subestimare). Bias-ul și variația parametrilor estimați sunt importante surse de erori, care trebuie minimizate pentru a putea implementa algoritmi de învățare automată supervizată care să poată fi aplicați unor seturi diferite de date și deci care să fie valabili dincolo de datele de antrenare [39] [40].

**Bias-ul** înseamnă incapacitatea unui model de machine learning de a descrie și reprezenta în mod fidel relația între date. Un model cu un bias mare este un model care întrucât nu se poate ajusta pe datele de antrenament, suprasimplifică descrierea modelului, ceea ce determină o valoare mare a variației erorilor, astfel că principiul minimizării sumei pătratelor erorilor nu poate fi respectat [39] [29] [41].

De exemplu dacă aplicăm un model liniar asupra unor date distribuite sub formă de curbă, vom avea underfitting, întrucât modelul liniar va simplifica reprezentarea printr-o dreaptă și va

genera o valoare mare a sumei pătratelor erorilor. Pentru a evita underfittingul vom crește complexitatea modelului. Astfel, un model polinomial de grad mai mare ar avea probabil un bias minim și o variație a erorilor apropiată de zero, reprezentând foarte fidel relația între date ( $SSE \rightarrow 0$  și  $R^2 \rightarrow 1$ , ceea ce arată un model perfect pentru explicarea datelor, în care suma pătratelor erorilor este 0).

**Variația** înseamnă diferența între rezultatele obținute prin aplicarea modelului pe datele de antrenare și cele obținute din aplicarea modelului pe datele de testare. Modelele cu variație mare reprezintă foarte bine datele de antrenament, cu media pătratelor erorilor aproape de zero, dar înregistrează erori mari la aplicarea pe datele de testare [39] [41] [29].

De exemplu un model polinomial de grad mai ridicat aplicat datelor distribuite sub formă de curbă, ar avea un bias foarte aproape de zero pe datele de antrenare, dar datorită overfittingului, ar genera o variație a erorilor foarte mare, mult mai mare decât modelul liniar, atunci când ar fi aplicat datelor de testare. ( $SSE \rightarrow 1$  și  $R^2 \rightarrow 0$ , față de valorile aplicate setului de antrenare, care arată un model perfect).

În concluzie:

- dacă modelul implementat este prea simplu și are prea puțini parametrii, va avea un bias mare la datele de antrenare, dar va avea o variație mai mică decât modelul complex la datele de testare
- dacă modelul implementat este prea complex și are un număr prea mare de parametrii va avea un bias mic la datele de antrenare, dar va avea o variație mare la datele de testare

Prin urmare, pentru a avea un algoritm eficient care să modeleze în mod fidel relația între date și care să producă predicții consistente prin aplicarea asupra unor seturi diferite de date, este nevoie să existe un compromis între bias-ul și variația parametrilor estimați.

Compromisul dintre bias-ul și variația parametrilor estimați constituie baza conceptuală a metodelor de regularizare (optimizare) a regresiei, cum sunt **regresia lasso** și **regresia ridge**. Metodele de regularizare introduc în soluțiile de regresie conceptul de bias, care poate reduce variația parametrilor în mod considerabil, față de metoda celor mai mici pătrate.

### 3.7.4. Regresia Ridge

Regresia ridge (regularizare L2) presupune calcularea unei noi linii de regresie, care nu se potrivește perfect datelor de antrenare, prin introducerea unei cantități mici de Bias, numită factor de regularizare, cu scopul de a evita supraînvățarea modelului (overfitting). Chiar dacă

linia de regresie ridge nu reprezintă perfect estimatorii ca în cazul metodei celor mai mici pătrate, regresia ridge asigură predicții mai generalizate și pentru seturi de date noi, nu doar pentru cele de antrenare, întrucât determină scăderea variației erorilor la un nivel semnificativ.

Regresia ridge adaugă o penalizare funcției de cost (media pătratelor erorilor) a modelului regresiei liniare bazat pe minimizarea pătratelor erorilor [41] [29]:

$$mse = \frac{1}{n} \sum (y_i - \hat{y}_i)^2 + \lambda * \sum b^2$$

unde:  $b^2$  este penalitatea L2 la metoda celor mai mici pătrate, reprezentată de pătratul coeficientului de regresie

$\lambda$  determină mărimea penalității și poate lua valori cuprinse între 0 și infinit

Atunci când  $\lambda = 0 \Rightarrow$  dreapta de regresie liniară este egală cu dreapta de regresie ridge. La fiecare creștere a lui  $\lambda$ , panta devine tot mai puțin abruptă și se apropie asimptotic de axa ox, astfel că variabila dependentă  $y$  devine tot mai puțin sensibilă la modificările variabilei independente  $x$ .

Regresia ridge îmbunătățește predicțiile pentru seturi noi de date întrucât reduce variația erorilor pentru datele de test, făcând predicția mai puțin sensibilă față de datele de antrenare.

În cazul regresiei liniare multiple toți parametrii sunt regularizați în mod egal, factorul de regularizare fiind produsul dintre  $\lambda$  și suma pătratelor parametrilor.

Valoarea lui  $\lambda$  se determină prin tehnica validării încrucișate. Astfel, se încearcă mai multe valori pentru  $\lambda$  și se folosește validarea încrucișată pentru a determina care variantă are cea mai mică variație a erorilor și cel mai mare scor al coeficientului de determinare  $R^2$ .

### 3.7.5. Regresia Lasso

Regresia lasso (regularizarea L1) este similară cu regresia ridge, prin introducerea unui factor de regularizare. În acest caz, penalitatea L1 este reprezentată de valoarea absolută a coeficientului de regresie [41] [29]:

$$mse = \frac{1}{n} \sum (y_i - \hat{y}_i)^2 + \lambda * \sum |b|$$

Diferența între regularizarea L2 și L1 este aceea că, în cazul regresiei ridge panta se apropie asimptotic de 0, fără să ajungă la valoarea 0, iar în cazul regresiei lasso panta se apropie de 0, până la valoarea 0.



În cazul regresiei liniare multiple, regresia lasso poate duce la excluderea anumitor variabile care nu sunt necesare ecuației de regresie, ducând la reducerea variației în modelele care conțin variabile inutile în explicarea modelului și la o ecuație finală mai simplă și mai ușor de interpretat. Astfel, regresia lasso contribuie la minimizarea numărului variabilelor independente în explicarea unui model de regresie liniară multiplă.

### **3.7.6. Implementarea comparativă a algoritmilor de machine learning**

În scopul validării unui model de regresie se pot implementa mai mulți algoritmi de regresie, asupra cărora au fost aplicate validarea simplă, prin împărțirea setului de date în date de antrenare și date de testare și/sau validarea încrucișată, în scopul autoevaluării performanțelor fiecărui algoritm în parte.

Scorurile obținute prin aplicarea diverselor modele sunt comparate, atât pentru validarea reprezentativității fiecărui algoritm în parte, cât și pentru a identifica care este modelul cel mai eficient în explicarea datelor și a relațiilor dintre ele.

## **4. Studiu de caz – predicția prețului autoturismelor prin analiza de regresie**

### **4.1. Tehnologii utilizate pentru analiza și modelarea datelor**

Pentru prezentul proiect de data science, a fost utilizat Python, împreună cu Jupyter Notebook și pachetele Numpy, Pandas, Matplotlib și Scikit-learn.

#### ***Python***

Python este un limbaj de programare dinamic, multifuncțional, de nivel înalt, care face parte din categoria limbajelor de programare interpretate și care permite atât programarea imperativă, funcțională sau procedurală, cât și programarea orientată obiect [42] [43].

Limbajele de programare interpretate funcționează pe baza unui interpretor, adică un program care va interpreta codul scris și îl va transforma în cod mașină. Codul scris de utilizator va rămâne în mod text, adică nu va fi compilat până la nivel de cod mașină, până la momentul rulării. Din punct de vedere al funcționalității, avantajul unui limbaj interpretat este faptul că poate rula același program pe toate platformele unde există un interpretor (dacă scriem un program Python pe Windows, putem să-l rulăm și pe Linux sau Mac și invers).

Un mare avantaj al utilizării Python decurge din faptul că are o gamă largă de biblioteci standard de metode sau funcții implementate care pot fi folosite (lucrurile cu fișiere externe, conectare la baze de date, module pentru implementarea de interfețe grafice, etc.) și permite extinderea funcționalității prin numeroase pachete adiționale programate și puse la dispoziție de terți, care sunt create pentru a îndeplini diverse funcții. În noiembrie 2019, PyPI - Python Package Index (Indexul pachetelor Python) [44], colecția oficială a pachetelor software create de terți, conținea peste 200.000 de pachete, cu o arie largă de funcționalități, printre care: baze de date, analiză de date, interfețe grafice, procesare de imagini, informatică științifică, web scraping, machine learning, etc.

Prin toate aceste facilități pe care le pune la dispoziție, Python este un sistem multifuncțional foarte util în data science, deoarece facilitează: colectarea datelor din diverse surse (fișiere externe, baze de date, web scraping), pregătirea datelor (analizarea datelor brute, curățarea și prelucrarea lor prin transformarea din date brute în date viabile pentru prelucrare), analiza statistică descriptivă și inferențială a datelor, reprezentarea grafică a datelor pentru vizualizarea distribuirii datelor și a diverselor corelații între ele, modelarea datelor prin tehnici de statistică și machine learning

### ***Jupyter Notebook***

Proiectul Jupyter (dezvoltat din IPython începând cu anul 2014) a fost creat pentru dezvoltarea de software open-source și servicii pentru computing interactiv, și suportă medii de execuție pentru mai mult de 40 de limbaje de programare, printre care și Python.

Proiectul Jupyter a dezvoltat produsul pentru computing interactiv Jupyter Notebook (cunoscut inițial ca IPython Notebooks), care este o aplicație web open-source care permite crearea și partajarea de documente care conțin cod Python [45]. Un document scris în Jupyter Notebook este un document JSON, care conține o listă ordonată de celule de input și output, care conțin cod, text, formule matematice, grafice, etc. Jupyter Notebook permite vizualizare grafică de date, ecuații, analiza datelor și modelare statistică, machine learning, etc. și este foarte adecvat pentru proiecte de data science datorită interfeței grafice implicite și a modului de afișare a datelor, sub formă de jurnal, ca o succesiune de celule [45] [46].

### ***NumPy***

NumPy este o bibliotecă pentru Python, care permite lucrul eficient cu serii foarte mari de date prezentate sub formă de vectori și matrici unidimensionale sau multidimensionale. Pachetul NumPy pune la dispoziție o colecție mare de operatori și funcții matematice pentru manipularea rapidă a datelor de tip arrays. NumPy permite și procesarea imaginilor.

Pachetul NumPy, care suportă și aborarea orientată obiect, permite efectuarea rapidă de operații cu matrici, inclusiv manipulare matematică, logică, manipularea formei matricilor prin modificarea liniilor și coloanelor, sortarea și selectarea datelor după anumite criterii, operații de algebră liniară, operații statistice simple, etc [47].

Întrucat seriile de date care se analizează în proiectele de data science se prezintă sub formă de tabele și sunt procesate sub formă de matrici, pachetul NumPy oferă suportul pentru manipularea eficientă a acestor date, prin faptul că permite transformarea matricilor, selectarea anumitor coloane și linii, compunerea a două sau mai multe matrici, eliminarea anumitor elemente ale matricilor, etc. astfel ca datele brute analizate să fie transformate după cum este necesar [48].

### ***Pandas***

Este o bibliotecă Python, un instrument pentru manipularea și analiza seriilor de date. Pachetul pune la dispoziția utilizatorilor structuri de date și operații pentru manipulare de tabele și analiza time-series, care presupune observarea evoluției unor date furnizate de același eșantion pe diverse perioade de timp [49].

Datele sunt disponibile sub formă de tabele, desfășurate sub formă de linii și coloane, iar Pandas le stochează într-un DataFrame, o structură de date 2-dimensională care poate stoca date numerice, caractere, date categoriale, etc. Pachetul conține instrumente pentru citirea și scrierea datelor, analizarea datelor lipsă sau a anomaliilor, modificarea tabelelor, obținerea de sub-seturi de date din seturi mari de date, inserarea sau ștergerea de coloane și linii în tabele, compunerea unui tabel din mai multe tabele, selectare și sortare după diverse criterii, filtrare de date, împărțirea datelor în grupuri de date, etc. Totodată, Pandas permite diverse analize statistice pentru date numerice, facilitând calcularea diversilor indicatori statistici, necesari în analiza datelor [50].

### ***Matplotlib***

Matplotlib este o bibliotecă open-source pentru Python pentru vizualizare grafică a datelor cantitative și statistică descriptivă. Pachetul furnizează o interfață API orientată obiect, pentru prezentarea de grafice, în scopul vizualizării statice, animate sau interactive a datelor în Python. Datele pot fi prezentate sub formă de histogramme (tabel de frecvențe), diagrame cu bare, diagrama de procente (pie-chart sau diagrama tip sector circular), diagrame de tip arie sau de tip suprafață, boxplots, scatterplot (grafic de dispersie), grafice time-series, etc [51].

Deoarece sunt analizate cantități mari de date (în general numerice), în scopul prelucrării, analizării și comunicării acestor date este necesară reprezentarea grafică a datelor pentru vizualizarea lor, în scopul observării modului de distribuire sau de împrăștiere a datelor, a identificării unor eventuale relații de cauzalitate sau de corelație între date, pentru a putea identifica diverse modele sau relații între date pe baza variabilelor analizate și pentru a descrie valori codificate ca obiecte vizuale [52].

### ***Scikit-learn (Sklearn)***

Scikit-learn este o bibliotecă Python pentru machine learning, concepută pentru a funcționa în relații de interoperabilitate cu bibliotecile NumPy și SciPy [53]. Scikit-learn utilizează pachetul NumPy pentru operații de algebră liniară și operații cu matrici. De asemenea, Scikit-learn se integrează foarte bine cu alte librării Python, printre care Pandas și Matplotlib.

Scikit-learn facilitează preprocesarea datelor sub formă de array, întrucât furnizează funcții și clase pentru transformarea datelor în reprezentări adecvate pentru a fi analizate, tehnici de factorizare a matricilor (analiza în componente principale, analiza factorială, etc). De

asemenea, pe baza calculelor statistice, algoritmi de machine learning pot identifica anomaliiile și extremele dintr-un set de date [54].

Pachetul implementează numeroși algoritmi pentru învățare automată supervizată, semi-supervizată și nesupervizată, respectiv algoritmi de regresie, analiza discriminantă liniară și quadratică, support vector machines, nearest neighbors, naive Bayes, arbori de decizie, rețele neuronale, clustering (grupare) și mulți alți algoritmi, fiind un suport important pentru proiectele de data science [54].

## 4.2. Analiza și modelarea datelor

### 4.2.1. Analiza preliminară

Setul de date care stă la baza prezentului studiu conține 205 observații și are 26 de attribute/variabile, de tip numeric și categorial. După eliminarea unor attribute care nu sunt necesare pentru analiza datelor, respectiv id, simbol, denumire auto, au rămas 22 de attribute, dintre care 14 de tip numeric și 8 de tip categorial, 21 fiind attributele explicative, conform Figura 4.1:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 205 entries, 0 to 204
Data columns (total 22 columns):
#   Column                Non-Null Count  Dtype
---  -
0   fueltype              205 non-null   object
1   aspiration             205 non-null   object
2   doornumber            205 non-null   object
3   carbody               205 non-null   object
4   drivewheel           205 non-null   object
5   wheelbase             205 non-null   float64
6   carlength             205 non-null   float64
7   carwidth              205 non-null   float64
8   carheight            205 non-null   float64
9   curbweight            205 non-null   int64
10  enginetype            205 non-null   object
11  cylindernumber        205 non-null   object
12  enginesize            205 non-null   int64
13  fuelsystem            205 non-null   object
14  boreratio             205 non-null   float64
15  stroke                205 non-null   float64
16  compressionratio      205 non-null   float64
17  horsepower            205 non-null   int64
18  peakrpm               205 non-null   int64
19  citympg               205 non-null   int64
20  highwaympg            205 non-null   int64
21  price                 205 non-null   float64
dtypes: float64(8), int64(6), object(8)
memory usage: 35.4+ KB
```

**Figura 4.1: Variabile inițiale car-price**

În scopul efectuării analizei datelor, observațiile de tip categorial au fost convertite în echivalente numerice prin operații de codare. Pentru fiecare atribut au fost calculați anumiți indicatori, respectiv media, abaterea standard, valorile minime și maxime, și nu au fost identificate date extreme sau erori de măsurare.

Pentru verificarea și validarea rezultatelor studiului, datele au fost împărțite în date de antrenare/train data și date de test/test data. Astfel, din cele 205 observații, 70% au fost selectate în mod aleator pentru antrenare și 30% pentru testare.

Variabilele independente sunt cele 21 de variabile numerotate de la 0-20, iar variabila price este variabila dependentă.

#### 4.2.2. Reducerea dimensionalității datelor

Pentru identificarea principalelor atribute/variabile independente care pot explica prețul, au fost utilizate mai multe metode.

##### *Filtrul corelației slabe între predictorii și variabila target*

Pentru a explica o variabilă pe baza altor variabile este necesar ca între variabilele independente și variabila dependentă să existe un anumit nivel de corelație. A fost utilizat coeficientul de corelație Pearson pentru determinarea variabilelor care au un coeficient de corelație cu variabila preț mai mare decât 0,5, iar variabilele cu un nivel al coeficientului de corelație sub 0,5 au fost eliminate.

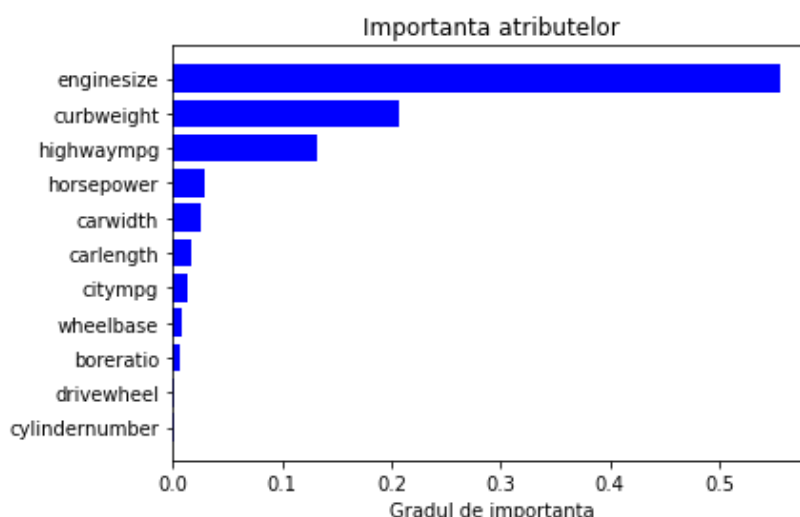
Numărul variabilelor esențiale pentru explicarea prețului s-a redus de la 21 la 11, conform Figura 4.2:

```
drivewheel      0.503460
wheelbase       0.550372
carlength       0.673316
carwidth        0.771034
curbweight      0.838979
cylindernumber  0.719352
enginesize      0.876751
boreatio        0.543367
horsepower      0.811184
citympg         0.736946
highwaympg      0.749342
price           1.000000
Name: price, dtype: float64
```

**Figura 4.2: Variabile esențiale car-price**

### ***Random Forest pentru determinarea importanței atributelor***

În urma aplicării Random Forest feature importances, cu un scor semnificativ de 0,9850 s-au obținut următoarele rezultate, conform Figura 4.3:



***Figura 4.3: Random Forest – importanța atributelor***

Scorul celor mai importante attribute se prezintă astfel:

```
1 0.5550654500067229
2 0.20750191106669158
3 0.1324827167558145
4 0.02972178482112183
```

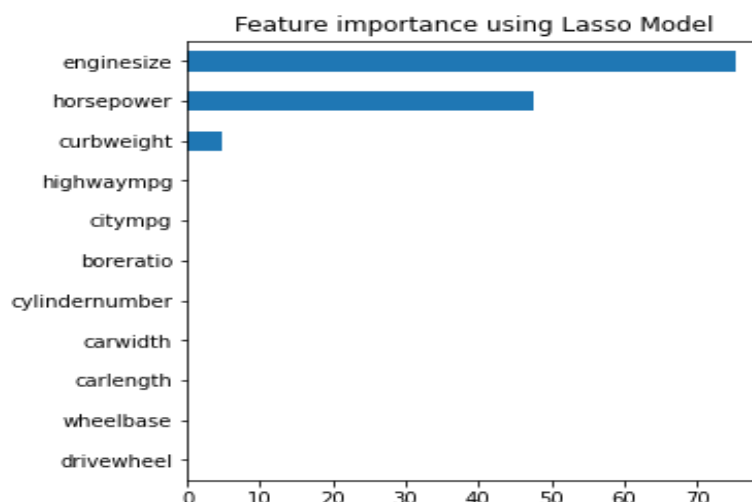
Este de reținut faptul că Random Forest va determina cele mai semnificative attribute, dar va subestima attributele care sunt foarte corelate între ele. Prin urmare a fost utilizată încă o metodă de selecție a atributelor și au fost comparate rezultatele.

### ***Regresia Lasso pentru selectarea atributelor***

Prin aplicarea regresiei lasso cu validare încrucișată a fost determinat cel mai bun coeficient de penalizare având valoarea de 3409,49, iar scorul (coeficientul de determinare) obținut de model la valoarea de 0,8064 este semnificativ. Rezultatele se prezintă conform Figura 4.4:

Scorul celor mai importante attribute se prezintă astfel:

```
(4.8572, 'curbweight')
(47.5523, 'horsepower')
(75.2475, 'enginesize')
```



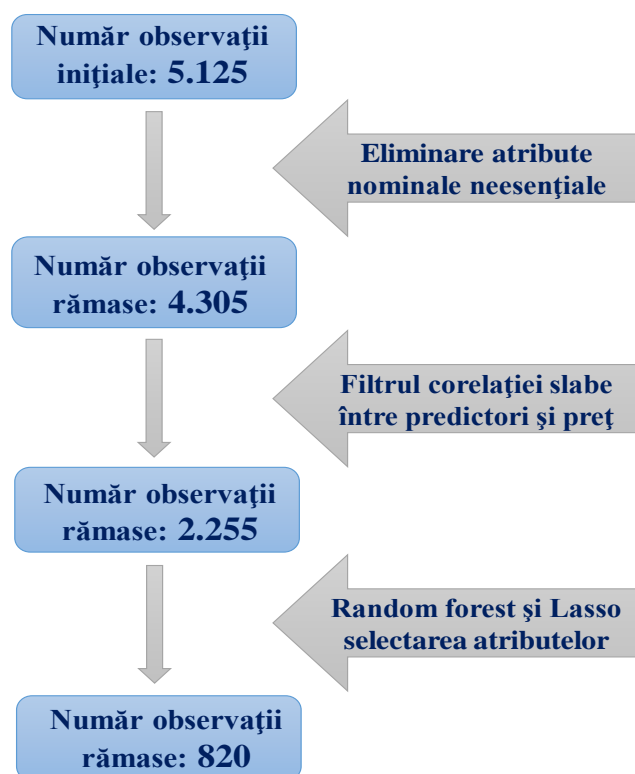
**Figura 4.4: Lasso – importanța atributelor**

După cum se poate observa, cele două modele de selecție a atributelor au rezultate asemănătoare, astfel că au fost reținute în explicarea modelului variabilele **enginesize**, **curbweight**, **horsepower** și **highwaympg**.

- **enginesize** reprezintă capacitatea cilindrică a motorului și determină puterea motorului și consumul de combustibil; se poate exprima în litri, centimetri cubi sau inci cubi;
- **curbweight** reprezintă greutatea totală a autovehiculului cu tot echipamentul standard;
- **horsepower** sunt caii putere care arată randamentul unui motor;
- **highwaympg** reprezintă un indicator al consumului de combustibil pe autostradă; nivelul consumului este determinat de tehnologii complexe și este dezirabil un consum cât mai mic, atât pentru economisirea resurselor cât și pentru scăderea emisiilor de CO<sub>2</sub>.

Prin utilizarea metodelor de reducere a dimensionalității datelor, setul de observații a fost redus substanțial, prin eliminarea atributelor explicative care nu sunt esențiale pentru explicarea variabilei target. Astfel, din 25 de attribute explicative aflate în setul inițial de date au rămas 4 attribute. De la 5.125 observații aflate în setul inițial de date, au rămas pentru modelare 820 de observații. Reducerea numărului de observații odată cu parcurgerea etapelor de reducere a dimensionalității este prezentată în Figura 4.5:

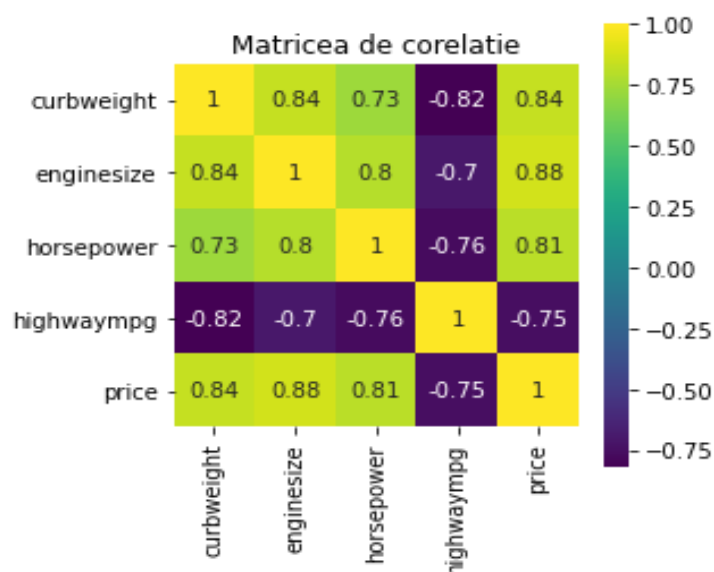




**Figura 4.5:** Reducerea nr. de observații în urma reducerii dimensionalității datelor

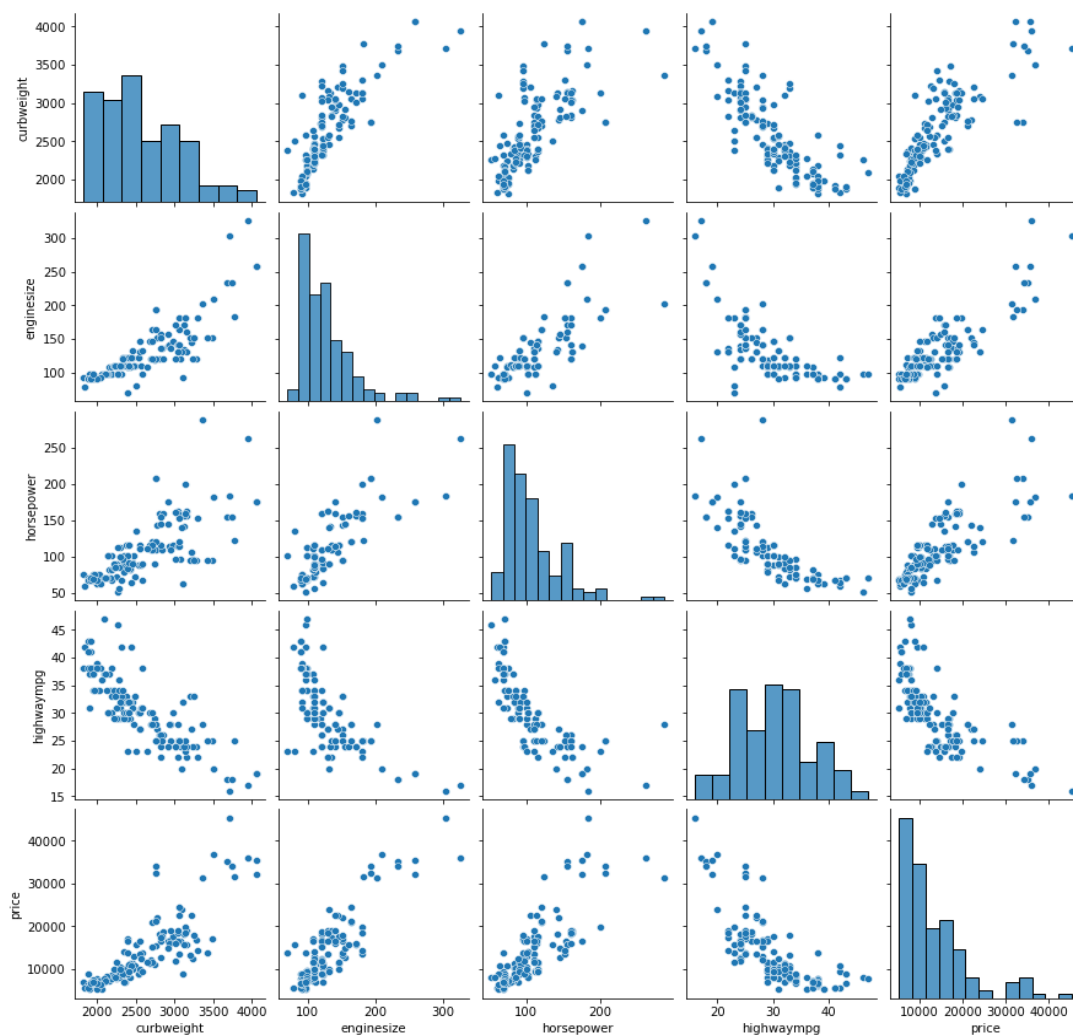
### 4.2.3. Determinarea strategiei și a tehnicilor de modelare

Pentru stabilirea gradului de corelație între toate variabilele implicate în model, a fost construită matricea de corelație, care poate fi vizualizată în Figura 4.6:



**Figura 4.6:** Matricea de corelație

Se poate observa faptul că între toate variabilele există un nivel de corelație ridicat, că toate variabilele au corelație pozitivă, în afară de variabila highwaympg, care este corelată negativ cu celelalte variabile, inclusiv cu variabila price. Pentru vizualizarea relațiilor dintre variabilele implicate în model s-a generat graficul perechilor de variabile, conform Figura 4.7.



**Figura 4.7: Graficul perechilor de variabile**

Vizual se observă existența unei relații care poate fi modelată liniar între variabila independentă enginesize și preț, iar între celelalte variabile independente și variabila target se observă existența unor relații neliniare.

Având în vedere gradul ridicat de corelație între fiecare dintre variabilele independente și variabila target, și pentru a respecta principiul multicolarității, se păstrează variabila care este cea mai corelată cu variabila dependentă, respectiv enginesize și sunt eliminate celelalte variabile din modelul explicativ. Întrucât se poate observa vizual o relație liniară între

enginesize și preț, s-a ales modelul regresiei liniare simple, care va fi completat cu regresia lasso și ridge, dar se va utiliza comparativ și regresia random forest.

Având în vedere specificul și scopul prezentului studiu, respectiv acela de a facilita întocmirea unui business plan pentru investiții, considerăm că varianta teoretică prin utilizarea doar a regresiei simple nu este suficientă pentru explicarea modelului. În demersurile de întocmire a unui plan de business, o companie are nevoie de cât mai multe informații și nu de un model simplificat care este insuficient pentru stabilirea strategiei de producție.

Astfel, se pornește de la ipoteza ca un model cu mai multe variabile independente va explica mai bine variabila dependentă și va aduce mai multe informații utile. Datorită acestor considerente se va efectua și analiza de regresie multiplă, păstrându-se toate cele 4 variabile independente în modelul explicativ.

**În concluzie**, vor fi implementate modele de regresie liniară, regresia ridge și lasso, regresia random forest, după cum urmează:

- regresia simplă:  $x = \text{enginesize} \rightarrow y = \text{preț}$ ;
- regresia multiplă:  $x_i = \text{enginesize, curbweight, horsepower, highwaympg} \rightarrow y = \text{preț}$ ;
- un model alternativ bazat pe analiza de regresie random forest.

În scopul evaluării performanțelor și a identificării celui mai semnificativ model, se va face o analiză comparativă a tuturor tehnicilor de modelare eficiente în explicarea prețului.

#### 4.2.4. Regresia liniară simplă

În urma implementării modelului regresiei liniare simple s-a obținut **ecuația funcției liniare de regresie**:

$$y = -7646.435077807868 + [165.76866245] x$$

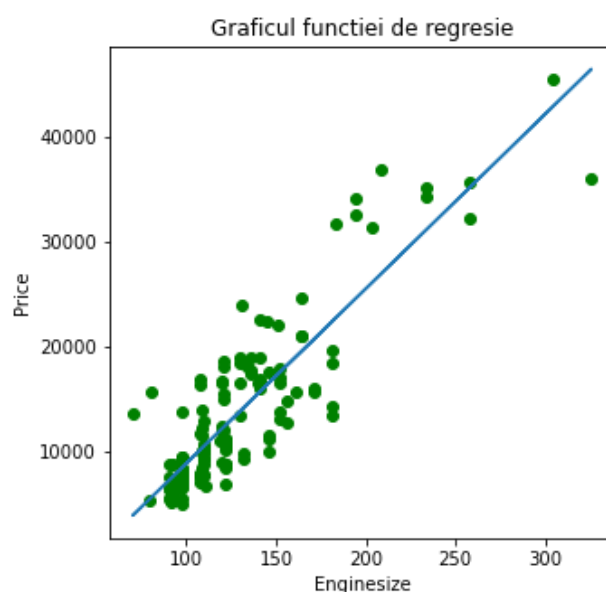
Rezultatele modelării prin regresia liniară simplă sunt prezentate în Figura 4.8:

OLS Regression Results						
Dep. Variable:	price	R-squared:	0.769			
Model:	OLS	Adj. R-squared:	0.767			
Method:	Least Squares	F-statistic:	468.6			
Date:	Thu, 27 May 2021	Prob (F-statistic):	1.14e-46			
Time:	16:57:50	Log-Likelihood:	-1379.9			
No. Observations:	143	AIC:	2764.			
Df Residuals:	141	BIC:	2770.			
Df Model:	1					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Intercept	-7646.4351	1022.771	-7.476	0.000	-9668.383	-5624.487
enginesize	165.7687	7.658	21.647	0.000	150.629	180.908
Omnibus:	9.817	Durbin-Watson:	2.111			
Prob(Omnibus):	0.007	Jarque-Bera (JB):	9.998			
Skew:	0.555	Prob(JB):	0.00674			
Kurtosis:	3.668	Cond. No.	432.			

**Figura 4.8: Rezultate regresia liniară simplă OLS**

### ***Ipotezele asupra proprietăților estimatorilor***

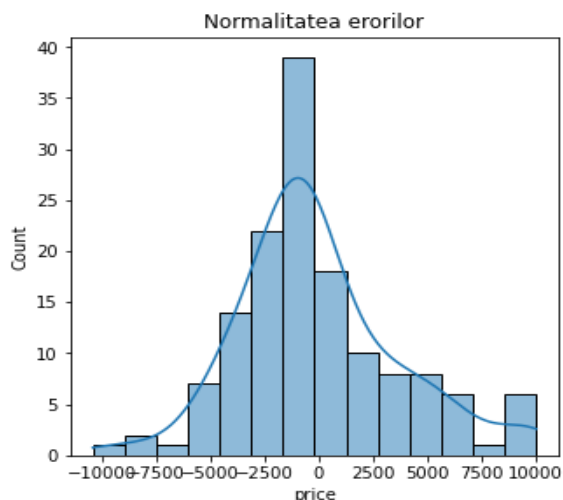
***Liniaritatea modelului:*** se poate observa vizual din Figura 4.9, unde se poate observa existența unei relații liniare între cele două variabile. Coeficientul de corelație Pearson cu valoarea de 0,88 confirmă ipoteza existenței unei relații liniare.



**Figura 4.9: Graficul funcției de regresie liniară simplă**

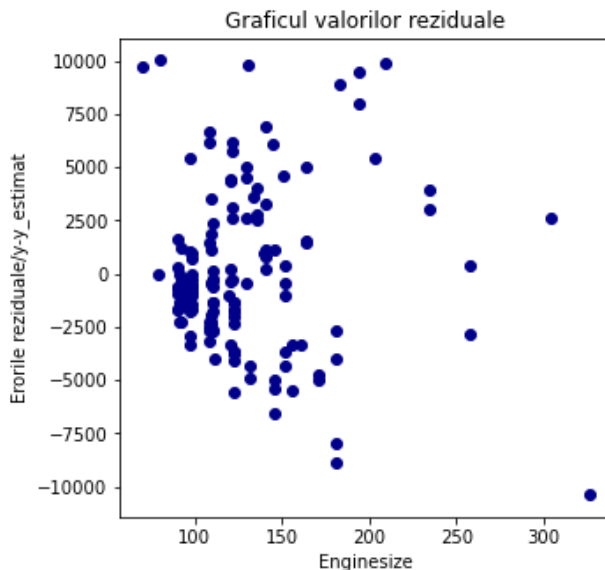
***Media erorilor*** este 8.1409, fiind foarte apropiată de 0, ceea ce confirmă ipoteza.

**Normalitatea erorilor** se poate confirma vizual prin histograma erorilor prezentată în Figura 4.10 și prin verificarea coeficientului de asimetrie (Skew) de 0,555 - apropiat de 0 și a coeficientului de boltire (Kurtosis) de 3,668 - apropiat de 3. Putem considera că erorile reziduale sunt distribuite normal.



**Figura 4.10: Normalitatea erorilor - regresia liniară simplă**

**Homoscedasticitatea** s-a verificat prin reprezentarea grafică a relației dintre variabila  $x$  și reziduuri, prezentată în Figura 4.11:



**Figura 4.11: Graficul valorilor reziduale – regresia liniară simplă**

Se poate observa faptul că punctele nu formează un anumit model sau tipar, astfel ca ipoteza homoscedasticității se confirmă.

**Necorelarea erorilor** a fost verificată prin testul Durbin Watson care are o valoare de 2,111. Pentru  $\alpha=0,5$  limita inferioară citită din tabelul Durbin Watson este de 1,694, iar limita

superioară este 2,346. Întrucât valoarea testului se află în intervalul valorilor teoretice, putem considera ca nu există autocorelare între valorile reziduale.

#### ***Testarea semnificativității pantei de regresie și a semnificativității globale***

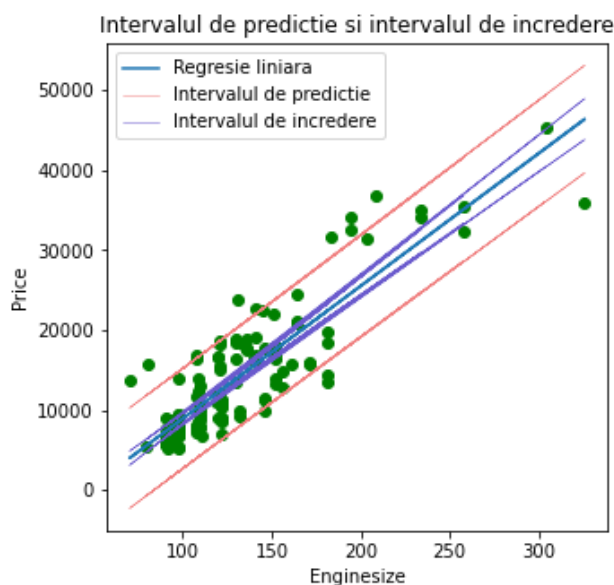
Testul  $t$  Student:  $t_{calculat} = 21,647 > t_{teoretic} = 1.14$ , se respinge ipoteza nulă.

Testul  $F$ :  $F_{calculat} = 468,6 > F_{teoretic} = 1.14$ , se respinge ipoteza nulă.

**Coefficientul de determinare** de 0,769 obținut pentru datele de antrenare este semnificativ. De asemenea pe datele de test s-a obținut un scor de 0,75.

**În concluzie**, întrucât modelul implementat a trecut testele de evaluare, se poate considera reprezentativ și poate fi utilizat pentru estimarea variabilei dependente price, în funcție de variabila independentă enginesize.

A fost dezvoltată o metodă de calcul a **intervalului de predicție** și a **intervalului de încredere** a pantei modelului de regresie, care sunt reprezentate grafic în Figura 4.12:



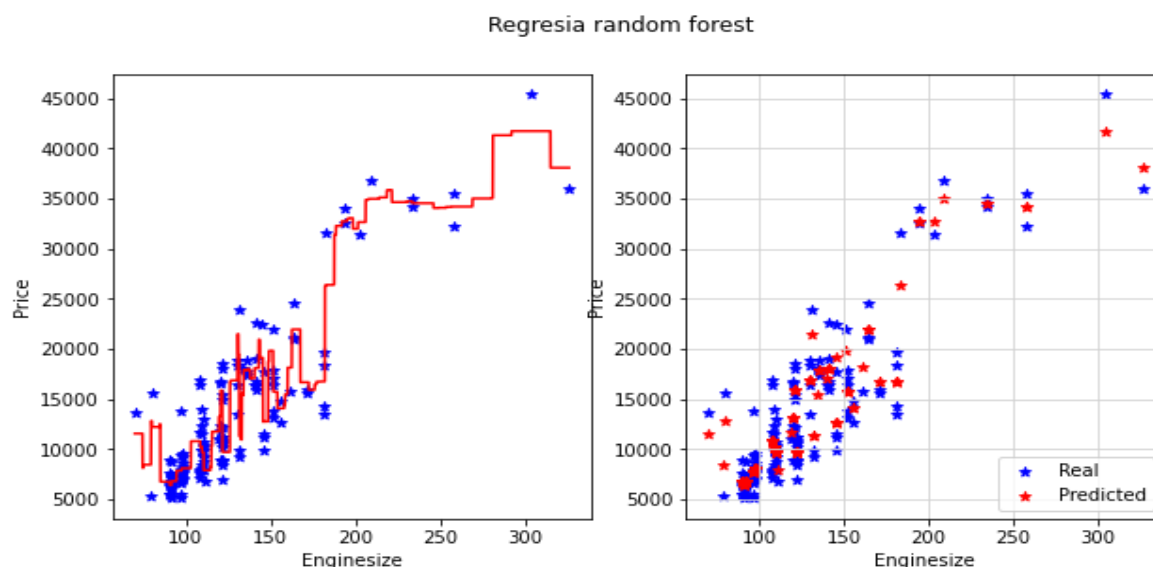
**Figura 4.12: Intervalul de predicție și intervalul de încredere**

Prin regresia lasso și ridge s-au obținut scoruri aproximativ egale cu scorul obținut de regresia liniară, astfel că a fost reținut pentru aplicare, modelul regresiei liniare.

#### **4.2.5. Regresia random forest cu o variabilă explicativă**

Prin aplicarea regresiei random forest la setul de date de antrenare a fost obținut un scor de 0,9315, iar pentru datele de test s-a obținut un scor de 0,9101, ambele superioare scorurilor

obținute de regresia liniară simplă. Grafic regresia random forest se prezintă conform Figura 4.13:



**Figura 4.13: Regresia random forest simplă: enginesize  $\rightarrow$  price**

Regresia random forest este un model neliniar, astfel că se valorile estimate se prezintă sub forma unui nor de puncte, foarte apropiate de valorile reale. Având în vedere performanțele obținute prin aplicarea modelului atât pe datele de antrenare cât și pe datele de test, modelul a fost reținut pentru aplicare.

#### 4.2.6. Regresia liniară multiplă

În urma implementării modelului regresiei liniare multiple s-a obținut **ecuația funcției liniare de regresie**:

$$y = -9524.689312612985 + [3.82247295 \ 82.78858615 \ 44.86893789 \ -71.26461143] * xi$$

Rezultatele modelării prin regresia liniară multiplă sunt prezentate în Figura 4.14:

OLS Regression Results						
Dep. Variable:	price	R-squared:	0.828			
Model:	OLS	Adj. R-squared:	0.823			
Method:	Least Squares	F-statistic:	165.7			
Date:	Sun, 23 May 2021	Prob (F-statistic):	1.20e-51			
Time:	11:08:39	Log-Likelihood:	-1358.9			
No. Observations:	143	AIC:	2728.			
Df Residuals:	138	BIC:	2743.			
Df Model:	4					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Intercept	-9524.6893	4973.971	-1.915	0.058	-1.94e+04	310.361
curbweight	3.8225	1.229	3.110	0.002	1.392	6.253
enginesize	82.7886	14.580	5.678	0.000	53.960	111.618
horsepower	44.8689	13.112	3.422	0.001	18.942	70.796
highwaympg	-71.2646	87.354	-0.816	0.416	-243.990	101.460
Omnibus:	5.568	Durbin-Watson:	1.988			
Prob(Omnibus):	0.062	Jarque-Bera (JB):	6.419			
Skew:	0.252	Prob(JB):	0.0404			
Kurtosis:	3.907	Cond. No.	4.74e+04			

**Figura 4.14: Rezultate regresia liniară multiplă OLS**

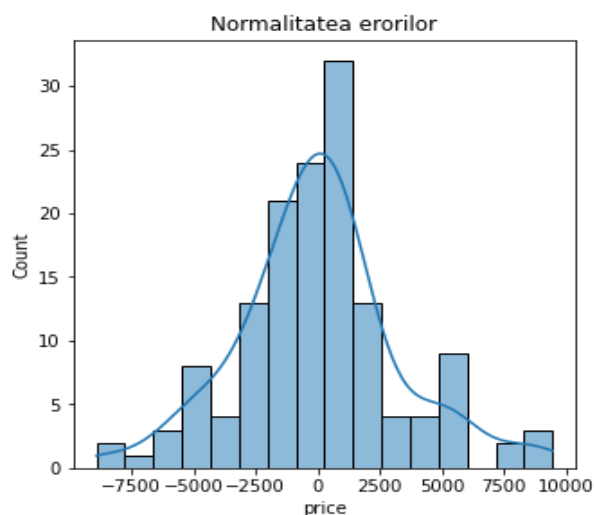
#### ***Ipotezele asupra proprietăților estimatorilor***

***Liniaritatea modelului:*** se poate observa vizual din graficul perechilor de variabile, existența unor relații dispuse liniar între fiecare variabilă independentă și variabila dependentă.

***Media erorilor*** este de -4,2231 fiind foarte apropiată de 0, ceea ce confirmă ipoteza.

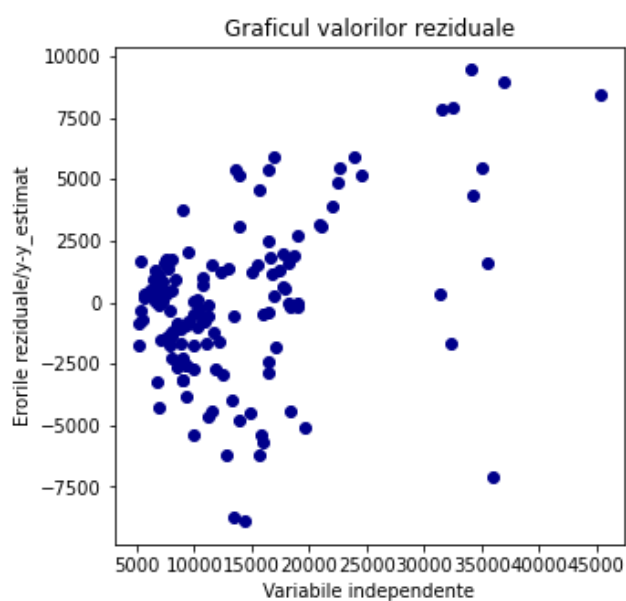
***Normalitatea erorilor*** se poate confirma vizual prin histograma erorilor prezentată în Figura 4.15 și prin verificarea coeficientului de asimetrie (Skew) de 0,252 - apropiat de 0 și a coeficientului de boltire (Kurtosis) de 3,907 - apropiat de 3. Putem considera că erorile reziduale sunt distribuite normal.





**Figura 4.15: Normalitatea erorilor – regresia liniară multiplă**

**Homoscedasticitatea** s-a verificat prin reprezentarea grafică a relației dintre variabila  $x$  și reziduuri, prezentată conform Figura 4.16:



**Figura 4.16: Graficul valorilor reziduale – regresia liniară multiplă**

Se poate observa faptul că punctele nu formează un anumit model sau tipar, astfel că ipoteza homoscedasticității se confirmă.

**Necorelarea erorilor** a fost verificată prin testul Durbin Watson care are o valoare de 1,988. Pentru  $\alpha=0,5$  limita inferioară citită din tabelul Durbin Watson este de 1,694, iar limita superioară este 2,346. Întrucât valoarea testului se află în intervalul valorilor teoretice, putem considera ca nu există autocorelare între valorile reziduale.

### ***Multicoliniaritatea***

Din matricea de corelație între variabilele independente reiese faptul că există un nivel ridicat de corelație între toate variabilele independente implicate, respectiv coeficienții de corelație în valoare absolută depășesc 0,7.

De asemenea, factorul de inflamare a dispersiei erorilor (VIF) prezintă valori foarte ridicate:

	VIF	variables
0	59.072436	curbweight
1	45.814650	enginesize
2	23.913537	horsepower
3	8.786194	highwaympg

Având în vedere indicatorii de mai sus, este evident faptul că modelul este afectat de multicoliniaritate, care are următoarele efecte: coeficienții de regresie vor fi instabili și vor avea o eroare standard mare, intervalul de încredere al coeficienților de regresie va fi mare, coeficienții de regresie vor fi greu de interpretat. Testarea semnificativității pantei de regresie și a semnificativității globale va determina valori foarte mici ale testelor, astfel că este foarte greu de demonstrat respingerea ipotezei nule. În cazul multicoliniarității se recomandă eliminarea variabilelor independente care sunt foarte corelate între ele.

Cu toate acestea, am considerat că un model cu mai multe variabile va aduce mai multă informație și este necesar pentru realizarea obiectivului studiului, acela de la întocmi un business plan pe baza rezultatelor obținute, întrucât este imposibil de a realiza o strategie eficientă de business, care presupune întocmirea unui buget de venituri și cheltuieli, și a unui plan de producție, pe baza informației oferite de o singură variabilă. Considerăm că principalele variabile care influențează în mod consistent prețul autovehiculelor sunt toate cele patru variabile independente: enginesize, curbweight, horsepower și highwaympg. Faptul că cele patru variabile sunt corelate între ele arată că setul de date este corect întocmit și informațiile sunt reale, întrucât în realitate cele patru componente auto (în special capacitatea motorului, caii putere și greutatea mașinii) sunt într-o anumită măsură corelate între ele și au într-adevar un rol esențial în determinarea prețului.

Prin urmare, dacă analizăm evoluția coeficientului de determinare și în paralel a coeficientului de determinare ajustat, putem observa următoarele:

- coeficientul de determinare  $R^2$  înregistrează o creștere de la valoarea de 0,769 în cazul regresiei simple, la 0,828 în cazul regresiei multiple;

- coeficientul de determinare ajustat  $R_{adj}^2$  care în cazul regresiei simple era de 0,767 înregistrează o creștere la valoarea de 0,823 în cazul regresiei multiple.

Evoluția lui  $R_{adj}^2$  pe măsură ce adăugăm variabile explicative este prezentată în Figura 4.17:

Variabile	Coeficientul de determinare ajustat
enginesize	0.769
enginesize + curbweight	0.801
enginesize + highwaympg	0.801
enginesize + horsepower	0.800
curbweight + horsepower + highwaympg	0.783
enginesize + curbweight + highwaympg	0.809
enginesize + horsepower + highwaympg	0.812
enginesize + curbweight + horsepower	0.823
enginesize + curbweight + horsepower + highwaympg	0.823

**Figura 4.17: Evoluția coeficientului de determinare ajustat**

Coeficientul de determinare ajustat  $R_{adj}^2$  este corectat cu gradele de libertate, astfel că prin adăugarea de noi variabile inutile în model acesta va scădea, iar prin adăugarea de variabile utile, acesta va crește.

După cum se poate observa, pe măsură ce se adaugă variabile explicative în model,  $R_{adj}^2$  crește, lucru care arată un câștig de informație ca urmare a adăugării de noi variabile.

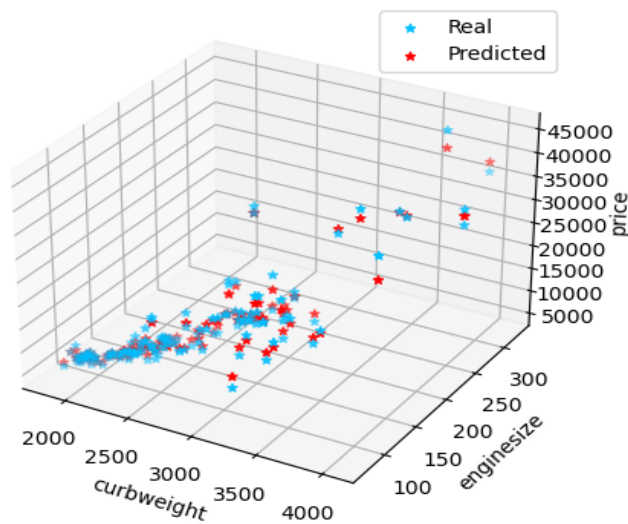
Într-un model afectat de multicolaritate este greu de explicat modul în care variabila  $y$  este afectată de către fiecare variabilă independentă în parte, datorită instabilității coeficienților de regresie, dar multicolaritatea nu scade capacitatea predictivă a modelului (lucru care se va vedea în continuare), astfel că modelul va fi reținut pentru aplicare.

#### 4.2.7. Regresia random forest cu patru variabile explicative

Prin aplicarea regresiei random forest cu patru variabile explicative la setul de date de antrenare a fost obținut un scor de 0,9852, iar pentru datele de test s-a obținut un scor de 0,94, ambele superioare scorurilor obținute de regresia liniară.

Pentru vizualizare grafică s-a modelat regresia random forest cu 2 variabile explicative, respectiv enginesize și curbweight, care a obținut un scor de 0,9798. Grafic regresia random forest cu doua variabile explicative se prezintă conform Figura 4.18:

Regresia random forest cu 2 variabile explicative



*Figura 4.18: Scatterplot regresia random forest cu 2 variabile explicative*

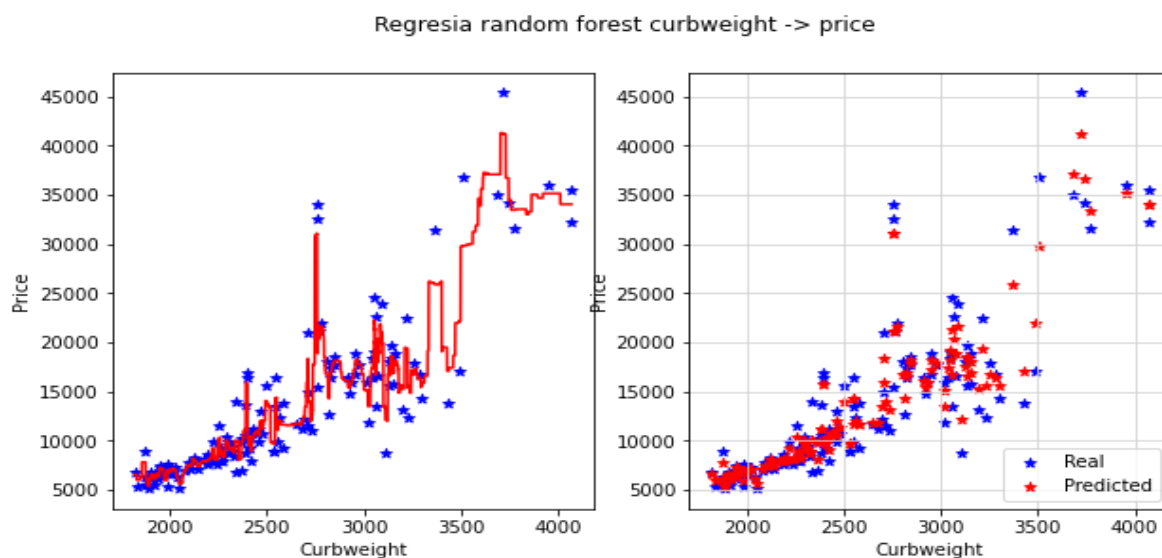
Regresia random forest este un model neliniar, astfel că se valorile estimate se prezintă sub forma unui nor de puncte, foarte apropiate de valorile reale. Având în vedere performanțele obținute prin aplicarea modelului atât pe datele de antrenare cât și pe datele de test, modelul a fost reținut pentru aplicare.

#### 4.2.8. Model alternativ bazat pe regresia random forest

Ca alternativă la modelele prezentate anterior a fost dezvoltat un model bazat pe regresiiile random forest simple. Astfel, a fost modelată regresia random forest separat pentru fiecare variabilă independentă:

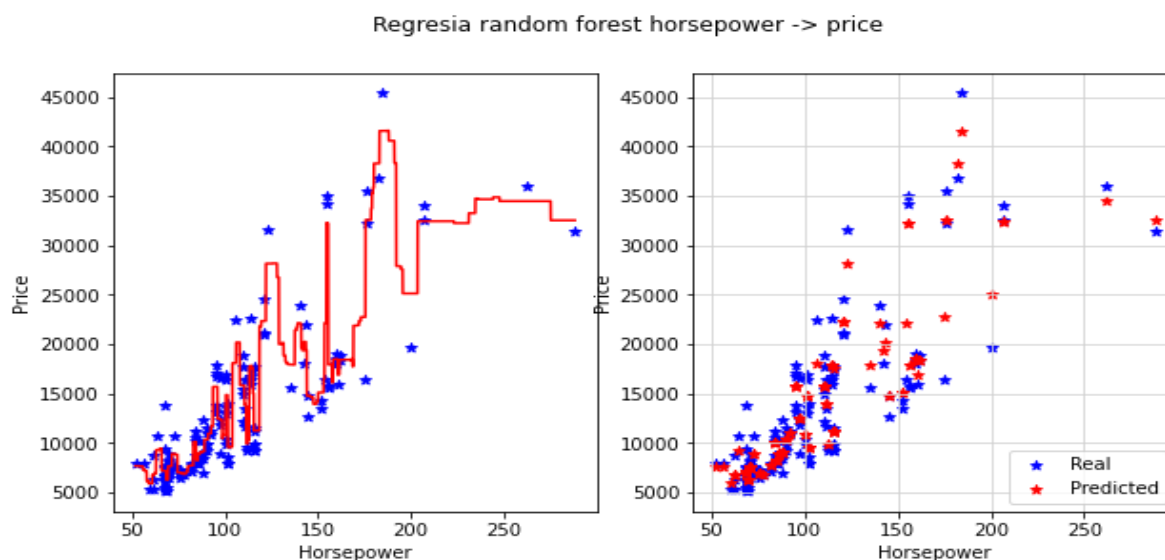
***Regresia random forest enginesize*** → ***price*** a fost prezentată anterior.

***Regresia random forest curbweight*** → ***price*** pentru setul de date de antrenare a obținut un scor de 0,9559, iar pentru setul de date de test a obținut un scor de 0,6999. Grafic modelul este prezentat în Figura 4.19:



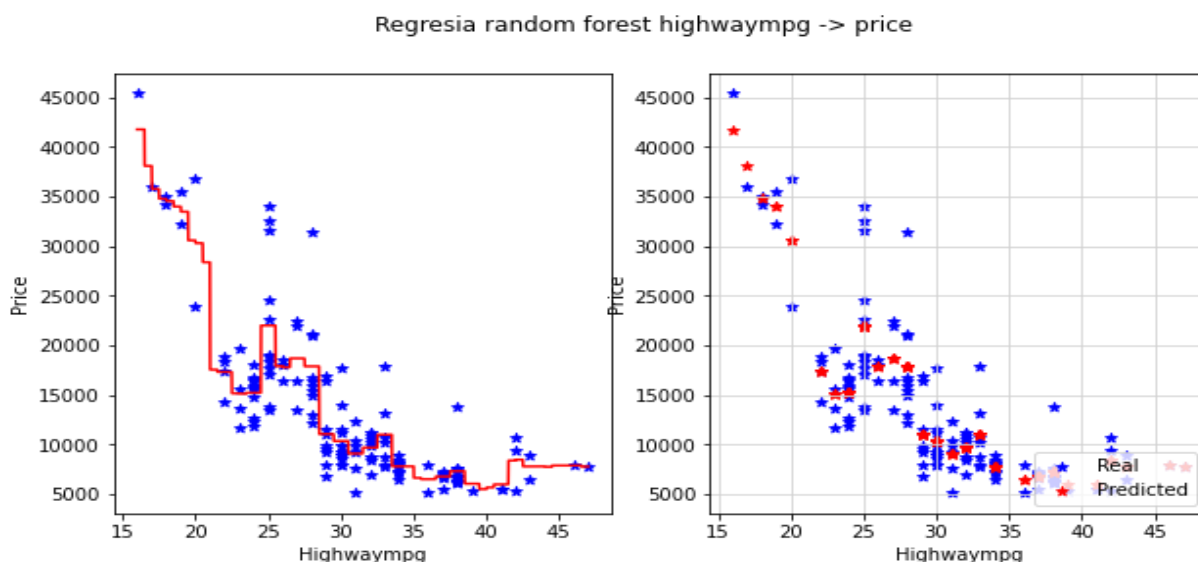
**Figura 4.19: Regresia random forest simplă: curbweight  $\rightarrow$  price**

**Regresia random forest horsepower  $\rightarrow$  price** pentru setul de date de antrenare a obținut un scor de 0,9390, iar pentru setul de date de test a obținut un scor de 0,8937. Grafic modelul este prezentat în Figura 4.20:



**Figura 4.20: Regresia random forest simplă: horsepower  $\rightarrow$  price**

**Regresia random forest highwaympg  $\rightarrow$  price** pentru setul de date de antrenare a obținut un scor de 0,8283, iar pentru setul de date de test a obținut un scor de 0,5966. Grafic modelul este prezentat în Figura 4.21:



**Figura 4.21: Regresia random forest simplă: highwaympg → price**

Au fost dezvoltate patru metode de predicție separate pentru fiecare model, iar pentru obținerea rezultatului final s-a calculat media celor 4 răspunsuri individuale, medie ponderată cu scorul obținut de fiecare metodă:

Pentru datele de antrenare s-a dezvoltat o metodă de predicție bazată pe următoarea ecuație:

$$y_{\text{estimat}} = (y_{\text{estimat\_enginesize}} * 0.9315 + y_{\text{estimat\_horsepower}} * 0.9390 + y_{\text{estimat\_curbweight}} * 0.9559 + y_{\text{estimat\_highwaympg}} * 0.8283) / (0.9315 + 0.9390 + 0.9559 + 0.8283)$$

Pentru datele de test s-a dezvoltat o metodă de predicție bazată pe următoarea ecuație:

$$y_{\text{estimat}} = (y_{\text{estimat\_enginesize}} * 0.9101 + y_{\text{estimat\_horsepower}} * 0.8937 + y_{\text{estimat\_curbweight}} * 0.6999 + y_{\text{estimat\_highwaympg}} * 0.5966) / (0.9101 + 0.8937 + 0.6999 + 0.5966)$$

#### 4.2.9. Compararea rezultatelor modelelor dezvoltate

Pentru evaluarea performanțelor modelelor implementate s-au comparat p-norme vectorilor rezultați din diferențele între variabila dependentă inițială  $y$  și  $\hat{y}$  estimat prin modelarea datelor, având în vedere faptul că variabila  $y$  și variabila  $\hat{y}$  estimat se prezintă sub forma unor vectori.

Dacă  $p$  este un număr real, iar  $p \geq 1$ , p-norma sau  $l_p$ -norma unui vector  $x_i = (x_1, \dots, x_n)$  este [55]:

$$||x||_p := (\sum_{i=1}^n |x_i|^p)^{1/p}$$

- dacă  $p = 1$ , avem 1-norma sau norma  $l_1$ , iar distanța care derivă din această normă se numește distanța  $l_1$ ; 1-norma este suma valorilor absolute ale coloanelor vectorului:

$$||x||_1 := \sum_{i=1}^n |x_i|$$

- dacă  $p = 2$ , avem norma euclidiană sau “radical din suma pătratelor”; astfel, pe un spațiu euclidian n-dimensional  $\mathbf{R}^n$ , lungimea vectorului este dată de formula:

$$||x||_2 := \sqrt{x_1^2 + \dots + x_n^2}$$

- dacă  $p \rightarrow \infty$ , p-norma se apropie de norma infinită sau norma maximă:

$$||x||_\infty := \max(|x_1|, \dots, |x_n|)$$

Toate aceste norme sunt echivalente, întrucât toate definesc aceeași topologie [55].

Pentru evaluarea performanțelor modelelor implementate s-a verificat diferența sau distanța între  $y$  inițial și  $\hat{y}$  estimat. Pentru ca un model să fie eficient, trebuie ca distanțele între valorile inițiale ale variabilei  $y$  și valorile lui  $\hat{y}$  estimat să fie cât mai mici.

Astfel, dacă avem:

$$d_{l_p}(y_i, \hat{y}_i) = \sqrt[p]{\sum |y_i - \hat{y}_i|^p}$$

- dacă  $p = 1$ , avem:

$$d_{l_1} = \sum |y_i - \hat{y}_i|, \quad \text{iar} \quad \frac{1}{n} * d_{l_1} = \frac{1}{n} * \sum |y_i - \hat{y}_i| = \text{MAE}$$

- dacă  $p = 2$ , avem:

$$d_{l_2} = \sqrt{\sum (y_i - \hat{y}_i)^2}, \quad \text{iar} \quad \sqrt{\frac{1}{n}} * d_{l_2} = \sqrt{\frac{1}{n}} * \sqrt{SSE} = \text{RMSE}$$

- dacă  $p \rightarrow \infty$ , avem valoarea maximă a vectorului, care arată care este valoarea lui  $\hat{y}$  estimat, cea mai depărtată de  $y$  inițial:

$$d_{l_\infty} = \max\{(y_i - \hat{y}_i), i = \overline{1, n}\}$$

Au fost calculați pentru fiecare dintre modele indicatorii, iar situația este prezentată în Figura 4.22:

Evaluarea performanțelor		Reg liniara simpla	Random forest 1 variabila explicativa	Reg liniara multipla	Random forest 4 variabile explicative	Model alternativ
$R^2$	train	0.77	0.93	0.83	0.99	-
	test	0.75	0.91	0.78	0.94	-
$d_{l_1}(MAE)$	train	2,814.75	1,594.25	2,352.26	639.18	1,201.17
	test	2,960.71	1,801.56	2,777.90	1,455.20	1,902.35
$d_{l_2}(RMSE)$	train	3,754.74	2,043.74	3,241.21	950.12	1,640.88
	test	4,132.71	2,496.18	3,922.16	2,008.81	2,792.89
$d_{l_\infty}(max)$	train	10,394.15	6,105.32	9,450.71	4,089.72	5,278.72
	test	14,315.78	6,435.13	14,018.59	6,706.10	10,807.45

**Figura 4.22: Tabel compararea rezultatelor modelelor**

Modelul random forest cu 4 variabile explicative a obținut cel mai mare coeficient de determinare  $R^2$ , respectiv 0.99 pentru datele de antrenare și 0.94 pentru datele de test.

Cele mai mici valori ale  $MAE$  și  $RMSE$  atât pentru datele de antrenare cât și pentru datele de test au fost obținute de modelul random forest cu 4 variabile explicative. După random forest cu 4 variabile explicative, cele mai mici valori ale  $MAE$  și  $RMSE$  pentru datele de antrenare au fost obținute de modelul alternativ, iar pentru datele de test, de către random forest cu o variabilă explicativă, diferența pentru datele de test între cele două fiind foarte mică.

În ceea ce privește valoarea maximă a vectorului rezultat din diferența între vectorul  $y$  initial și vectorul  $\hat{y}$  estimat, pe datele de antrenare cea mai mică valoare a fost obținută de către modelul random forest cu o variabilă explicativă, imediat urmând modelul alternativ, iar pe datele de test cea mai mică valoare a fost obținută de random forest cu o variabilă explicativă, la diferență foarte mică fiind modelul random forest cu 4 variabile explicative.

S-au acordat punctaje modelelor în funcție de performanța totală, performanța obținută pe datele de antrenare și performanța pe datele de test, iar rezultatele sunt prezentate în Figura 4.23:

Model	Punctaj pe datele de antrenare	Punctaj pe datele de test	Punctaj date antrenare și test
Random forest 4 variabile explicative	9	8	17
Random forest 1 variabilă explicativă	3	7	10
Model alternativ	6	3	9

**Figura 4.23: Punctaj comparativ modele implementate**



În concluzie, **cel mai performant model este regresia random forest cu 4 variabile** explicative, care a obținut cele mai bune rezultate atât pe datele de antrenare cât și pe datele de test, obținând astfel și un raport optim între bias și varianță. Chiar dacă modelul este afectat de multicolinearitate, aceasta nu scade capacitatea sa predictivă, care este cea mai ridicată dintre toate modelele aplicate.

Al doilea cel mai eficient model pe datele de antrenare este **modelul alternativ**, iar pe datele de test **modelul regresiei random forest simple** cu o variabilă explicativă. Aceasta arată faptul că modelul alternativ are un bias mai mic decât regresia random forest simplă și o variație mai mare decât aceasta, în timp ce regresia random forest simplă are un bias mai mare și o variație mai mică decât modelul alternativ.

Prin urmare, dacă nu ne interesează să extrapolăm modelul pentru alte date decât cele pe care le avem la dispoziție, modelul alternativ are performanțe mai ridicate decât modelul random forest simplu. Piața auto americană are un anumit specific care e posibil să nu mai fie întâlnit în altă parte, fiind caracterizată de preferința pentru autovehiule de dimensiuni mari, cu motoare de capacitate mare, existând o corelație ridicată între greutatea autovehiculului (curbweight), capacitatea cilindrică (enginesize) și randamentul motorului (horsepower), lucru care se observă din datele care stau la baza studiului. Piața auto din Europa de exemplu este diferită, fiind caracterizată de preferința pentru autovehicule de dimensiuni mai mici, care nu necesită neaparat motoare de capacitate mare, iar corelația între principalele variabile explicative din prezentul studiu este foarte posibil să nu fie așa de ridicată ca în cazul pieței auto americane.

Mai mult, întrucât modelul alternativ ține cont de toate cele patru variabile importante în explicarea modelului este mai adecvat pentru scopul studiului, acela de a determina care sunt variabilele importante care influențează prețul autovehiculelor și de a întocmi un business plan eficient.

Următorul model ca și eficiență este regresia liniară multiplă, care depășește performanțele regresiei liniare simple, arătând din nou faptul că multicolinearitatea nu reduce capacitatea predictivă a modelului.

Regresia random forest nu poate extrapola rezultatele pentru valori mai mari decât limita maximă sau mai mici decât limita minimă a valorilor variabilelor explicative.

$$x \geq x_{min} \text{ și } x \leq x_{max}$$

În afara intervalului de valori ale variabilelor explicative regresia random forest nu poate da rezultate. Spre deosebire de regresia random forest, modelul regresiei liniare funcționează și

pentru valori din afara intervalului de valori ale variabilelor independente. Prin urmare, dacă se dorește estimarea variabilei dependente și pentru valori dincolo de limita minimă sau maximă a variabilelor explicative, se poate utiliza modelul regresiei liniare.

Toate cele cinci modelele sunt valide și cu performanțe bune, fiecare putând fi utilizat pentru estimarea prețului unui autovehicul pe baza setului de variabile explicative. Cel mai performant model, care oferă cea mai multă informație despre variabila dependentă preț și care răspunde cel mai bine scopului prezentului studiu este regresia random forest cu patru variabile explicative (enginesize, curbweight, horsepower și highwaympg), întrucât vectorul de valori estimate este la cea mai mică distanță de vectorul de valori inițiale și reușește să obțină cel mai optim compromis între bias-ul pe datele de antrenare și varianța pe datele de test, putând fi extrapolat și la alte seturi de date.

## Bibliografie

- [1] Wikipedia, <https://ro.wikipedia.org/wiki/Econometrie> (accesare: 02.04.2021).
- [2] Wikipedia, [https://ro.wikipedia.org/wiki/Analiza\\_de\\_regresie](https://ro.wikipedia.org/wiki/Analiza_de_regresie) (accesare: 02.04.2021).
- [3] D. Cielen, A. D. B. Meysman and M. Ali, Introducing Data Science, Shelter Island: Manning Publications Co., 2016.
- [4] C. H. Lau, <https://towardsdatascience.com/5-steps-of-a-data-science-project-lifecycle-26c50372b492> (accesare: 05.04.2021).
- [5] M. Mayo, <https://www.kdnuggets.com/2019/06/7-steps-mastering-data-preparation-python.html> (accesare: 03.03.2021).
- [6] A. Bhandari, <https://www.analyticsvidhya.com/blog/2020/04/feature-scaling-machine-learning-normalization-standardization/> (accesare: 03.03.2021).
- [7] Official website IBM, <https://www.ibm.com/cloud/learn/exploratory-data-analysis> (accesare: 03.03.2021).
- [8] L. Sasu, <https://www.slideshare.net/lmsasu/curs-2-data-mining> (accesare: 04.03.2021).
- [9] P. Sharma, <https://www.analyticsvidhya.com/blog/2018/08/dimensionality-reduction-techniques-python/> (accesare: 04.03.2021).
- [10] Aman1608, <https://www.analyticsvidhya.com/blog/2020/10/feature-selection-techniques-in-machine-learning/> (accesare: 05.03.2021).
- [11] Wikipedia, [https://en.wikipedia.org/wiki/Machine\\_learning](https://en.wikipedia.org/wiki/Machine_learning) (accesare: 05.03.2021).
- [12] Wikipedia, <https://ro.wikipedia.org/wiki/Statistic%C4%83> (accesare: 07.03.2021).
- [13] D. Danciulescu, <http://inf.ucv.ro/documents/danciulescu/curs4-curs-5-curs6.pdf> (accesare: 20.11.2020).
- [14] D. Maniu, <http://www.phys.ubbcluj.ro/~dana.maniu/BIOSTAT/C2.pdf> (accesare: 20.11.2021).
- [15] UMFVC, <http://www.umfcv.ro/files/b/i/Biostatistica%20MG%20-%20Cursul%205%20-%20Corelatii.pdf> (accesare: 20.11.2020).
- [16] Wikipedia, <https://ro.wikipedia.org/wiki/Covarian%C8%9B%C4%83> (accesare: 20.11.2020).

- [17] Slideshare, <https://www.slideshare.net/Cattta89/regresie> (accesare: 01.12.2020).
- [18] M. Chavent, [http://www.math.u-bordeaux.fr/~mchave100p/wordpress/wp-content/uploads/2013/10/ModStat\\_C1\\_pres.pdf](http://www.math.u-bordeaux.fr/~mchave100p/wordpress/wp-content/uploads/2013/10/ModStat_C1_pres.pdf) (accesare: 20.02.2021).
- [19] R. Rakotomalala, Econometrie La regression lineaire simple et multiple, Lyon: Universite Lumiere Lyon 2, 2018.
- [20] M. T. Coadă, <https://www.slideshare.net/tiberiumarian92/49855810-capitolul2regresialiniarapp133slideej> (accesare: 25.11.2020).
- [21] A. Fahad, <https://machinelearningmind.com/2019/10/27/assumptions-of-linear-regression-how-to-validate-and-fix/> (accesare: 04.04.2021).
- [22] V. Cristescu și T. Sâia, [http://www.imst.pub.ro/Upload/Sesiune/ComunicariStiintifice/Lucrari\\_2015/06.16/16\\_L36.pdf](http://www.imst.pub.ro/Upload/Sesiune/ComunicariStiintifice/Lucrari_2015/06.16/16_L36.pdf) (accesare: 07.04.2021).
- [23] R. Atha, <https://medium.com/swlh/multi-linear-regression-using-python-44bd0d10082d> (accesare: 07.04.2021).
- [24] Wikipedia, [https://ro.xcv.wiki/wiki/Robust\\_regression](https://ro.xcv.wiki/wiki/Robust_regression) (accesare: 21.03.2021).
- [25] Wikipedia, [https://en.wikipedia.org/wiki/Random\\_sample\\_consensus#:~:text=Random%20sample%20consensus%20\(RANSAC\)%20is,as%20an%20outlier%20detection%20method.](https://en.wikipedia.org/wiki/Random_sample_consensus#:~:text=Random%20sample%20consensus%20(RANSAC)%20is,as%20an%20outlier%20detection%20method.) (accesare: 05.02.2021).
- [26] Official website scikit-learn, [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.RANSACRegressor.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.RANSACRegressor.html) (accesare: 15.03.2021).
- [27] D. Zaharie, [https://staff.fmi.uvt.ro/~daniela.zaharie/dm2017/RO/curs/dm2017\\_curs11.pdf](https://staff.fmi.uvt.ro/~daniela.zaharie/dm2017/RO/curs/dm2017_curs11.pdf) (accesare: 31.03.2021).
- [28] A. Chakure, <https://medium.com/swlh/random-forest-and-its-implementation-71824ced454f> (accesare: 01.04.2021).
- [29] J. Starmer, <https://www.youtube.com/user/joshstarmer> (accesare: 01.04.2021).
- [30] D. Mwit, <https://neptune.ai/blog/random-forest-regression-when-does-it-fail-and-why> (accesare: 03.04.2021).
- [31] S. Glen, <https://www.statisticshowto.com/probability-and-statistics/statistics-definitions/adjusted-r2/> (accesare: 31.03.2021).

- [32] A. Chugh, <https://medium.com/analytics-vidhya/mae-mse-rmse-coefficient-of-determination-adjusted-r-squared-which-metric-is-better-cd0326a5697e> (accesare: 04.09.2021).
- [33] J. Brownlee, <https://machinelearningmastery.com/regression-metrics-for-machine-learning/> (accesare: 08.04.2021).
- [34] D. N. Sadawi, <https://www.youtube.com/user/DrNoureddinSadawi/playlists> (accesare: 01.12.2020).
- [35] The Pennsylvania State University website, <https://online.stat.psu.edu/stat501/lesson/6/6.2> (accesare: 19.03.2021).
- [36] The Pennsylvania State University website, <https://online.stat.psu.edu/stat462/node/135/> (accesare: 19.03.2021).
- [37] A. Birlutiu, [http://adrianabirlutiu.uab.ro/cursuri/MIRF/note\\_curs\\_lab\\_5.pdf](http://adrianabirlutiu.uab.ro/cursuri/MIRF/note_curs_lab_5.pdf) (accesare: 30.03.2021).
- [38] A. Birlutiu, <http://adrianabirlutiu.uab.ro/cursuri/MIRF/2018curs5.pdf> (accesare: 15.03.2021).
- [39] S. Singh, <https://towardsdatascience.com/understanding-the-bias-variance-tradeoff-165e6942b229> (accesare: 28.02.2021).
- [40] Wikipedia, [https://ro.xcv.wiki/wiki/Bias%E2%80%93variance\\_tradeoff](https://ro.xcv.wiki/wiki/Bias%E2%80%93variance_tradeoff) (accesare: 28.02.2021).
- [41] D. Patel, <https://www.youtube.com/c/codebasics/playlists> (accesare: 15.11.2020).
- [42] Wikipedia, [https://en.wikipedia.org/wiki/Python\\_\(programming\\_language\)](https://en.wikipedia.org/wiki/Python_(programming_language)) (accesare: 01.02.2021).
- [43] Official website Python, <https://www.python.org/> (accesare: 01.02.2021).
- [44] Official website Python Package Index, <https://pypi.org/> (accesare: 01.02.2021).
- [45] Wikipedia, [https://en.wikipedia.org/wiki/Project\\_Jupyter](https://en.wikipedia.org/wiki/Project_Jupyter) (accesare: 01.02.2021).
- [46] Official website Project Jupyter, <https://jupyter.org/> (accesare: 01.02.2021).
- [47] Wikipedia, <https://en.wikipedia.org/wiki/NumPy> (accesare: 02.02.2021).
- [48] Official website Numpy, <https://numpy.org/> (accesare: 02.02.2021).
- [49] Wikipedia, [https://en.wikipedia.org/wiki/Pandas\\_\(software\)](https://en.wikipedia.org/wiki/Pandas_(software)) (accesare: 02.02.2021).
- [50] Official website Pandas, <https://pandas.pydata.org/> (accesare: 02.02.2021).

- [51] Wikipedia, <https://en.wikipedia.org/wiki/Matplotlib> (accesare: 03.02.2021).
- [52] Official website Matplotlib, <https://matplotlib.org/> (accesare: 03.02.2021).
- [53] Wikipedia, <https://en.wikipedia.org/wiki/Scikit-learn> (accesare: 04.02.2021).
- [54] Official website Scikit-learn, <https://scikit-learn.org/stable/> (accesare: 03.02.2021).
- [55] Wikipedia, [https://ro.wikipedia.org/wiki/Norm%C4%83\\_\(matematic%C4%83\)](https://ro.wikipedia.org/wiki/Norm%C4%83_(matematic%C4%83)) (accesare: 25.05.2021).