

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 α 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 maximumului 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 α 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 α stabilit. Pentru aflarea valorii teoretice se utilizează un risc de $\alpha/2$, întrucât distribuția Student este simetrică, iar suprafața de respingere α 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 ε :

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
Regresie	$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	$k - 1$	$MSR = \frac{SSR}{k - 1}$	$\frac{MSR}{MSE}$
Eroare	$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2$	$n - k$	$MSE = \frac{SSE}{n - k}$	
Total	$SST = \sum_{i=1}^n (y_i - \bar{y})^2$	$n - 1$	$\frac{SST}{n - 1}$	

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 α 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), $k-1$ subseturi 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

λ 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 λ , 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 λ și suma pătratelor parametrilor.

Valoarea lui λ se determină prin tehnica validării încrucișate. Astfel, se încearcă mai multe valori pentru λ 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 histograme (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:

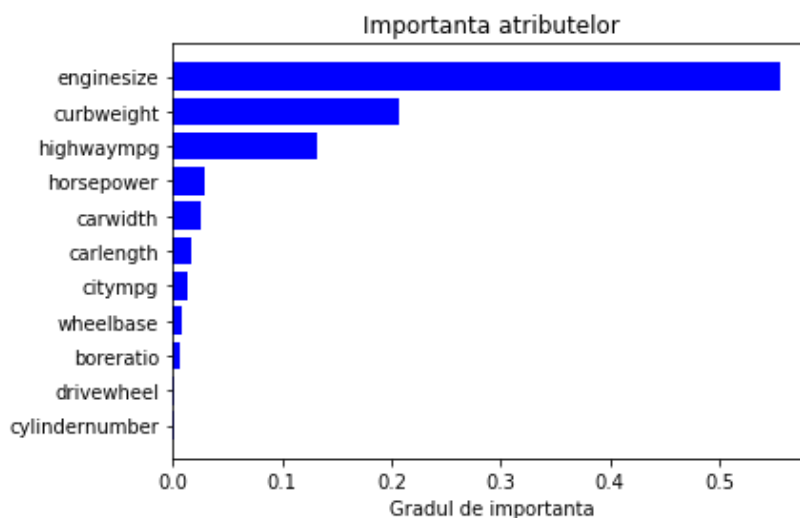


Figura 4.3: Random Forest – importanța atributelor

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')
```

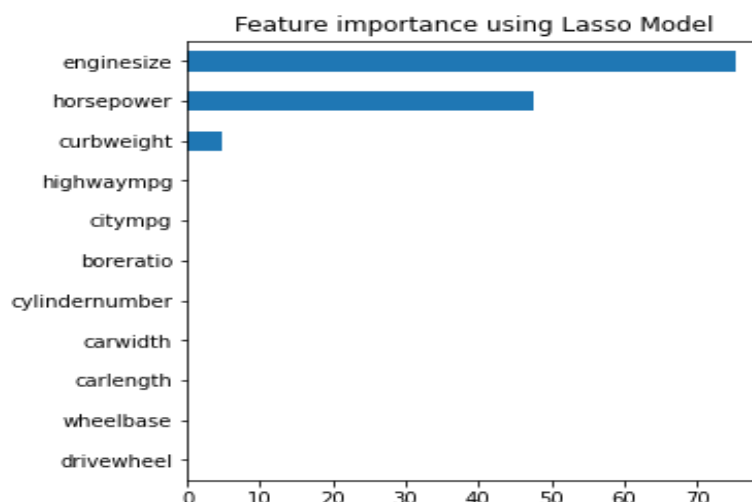


Figura 4.4: Lasso – importanța atributelor

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 CO₂.

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:

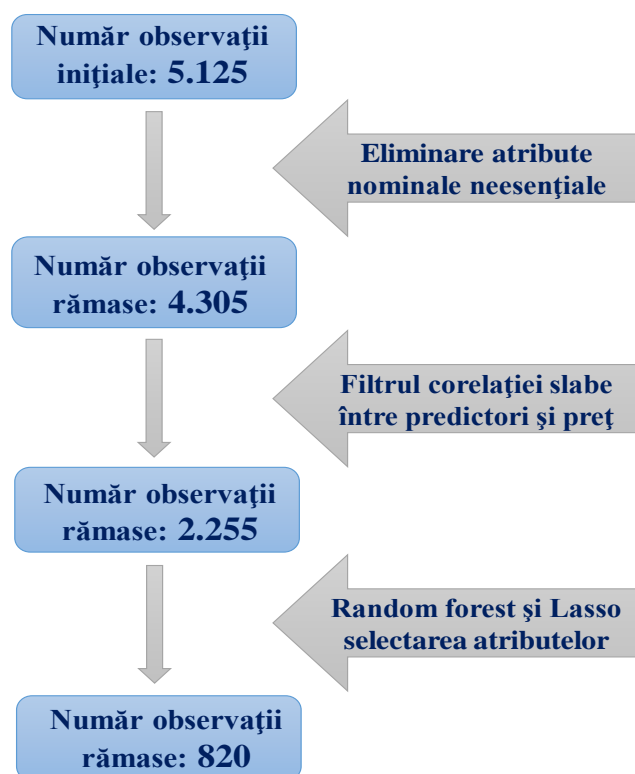


Figura 4.5: Reducerea nr. de observații în urma reducerii dimensionalității datelor

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:

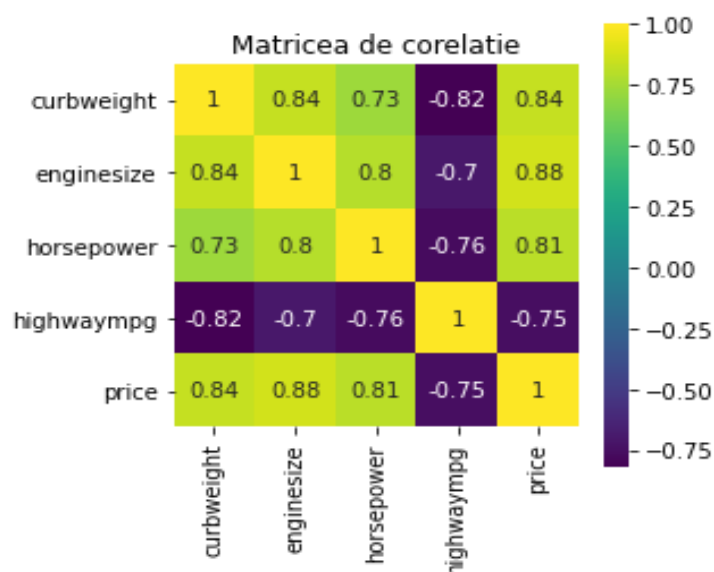


Figura 4.6: Matricea de corelație

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.

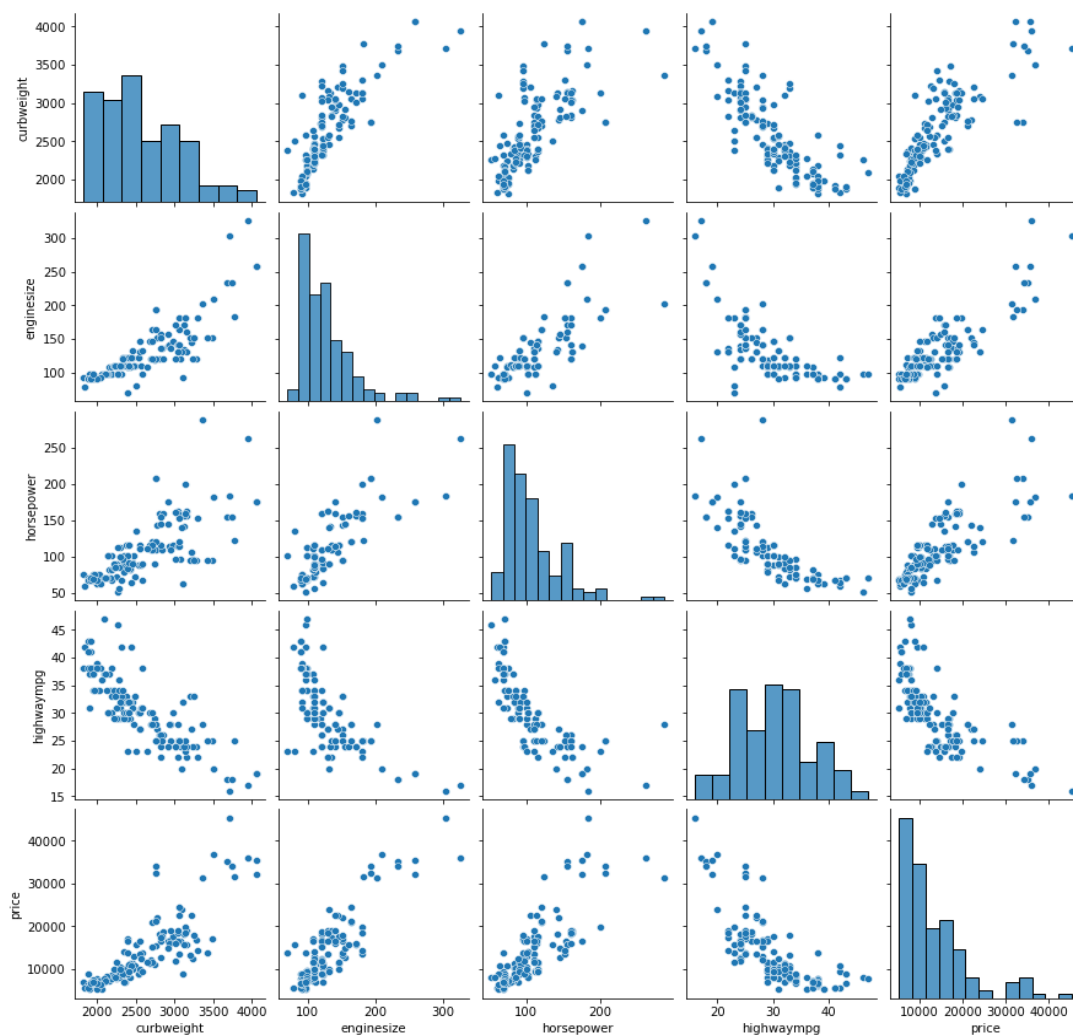


Figura 4.7: Graficul perechilor de variabile

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.

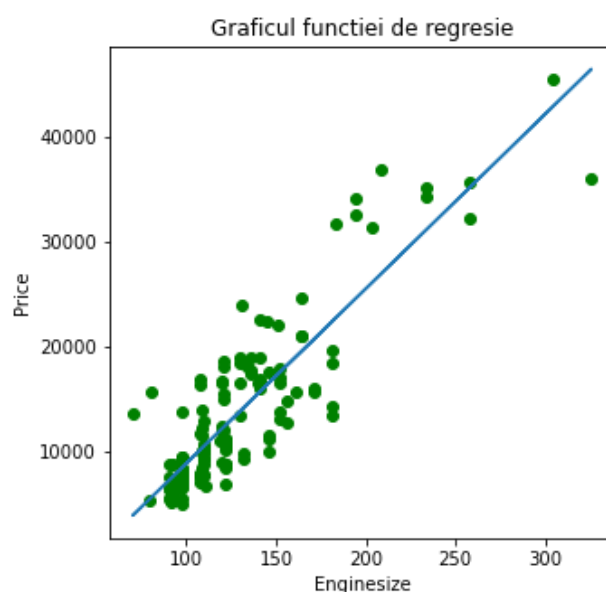


Figura 4.9: Graficul funcției de regresie liniară simplă

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.

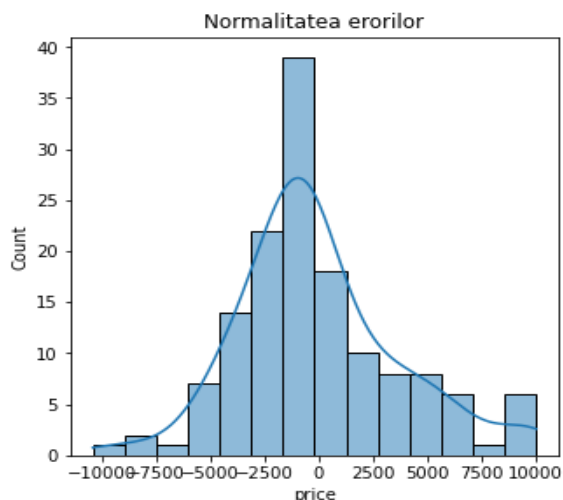


Figura 4.10: Normalitatea erorilor - regresia liniară simplă

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

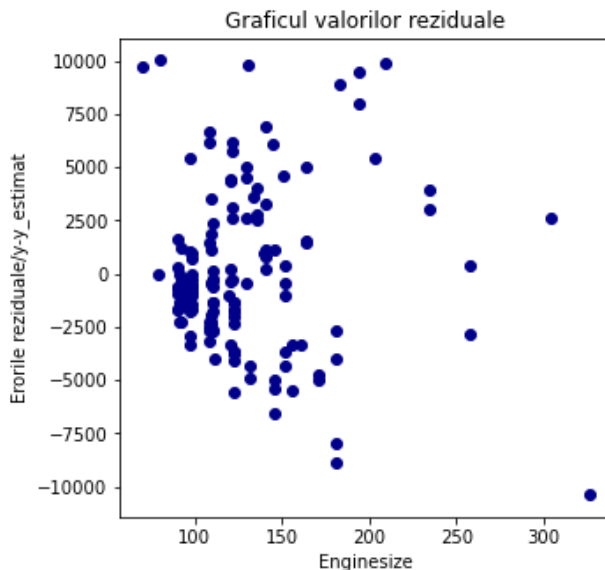


Figura 4.11: Graficul valorilor reziduale – regresia liniară simplă

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:

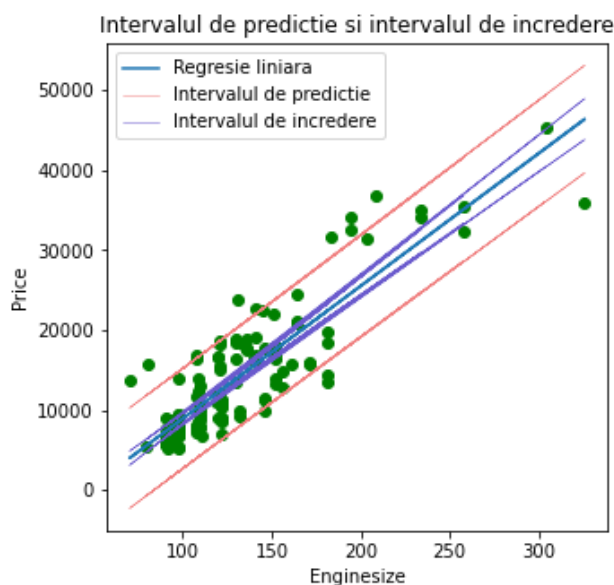


Figura 4.12: Intervalul de predicție și intervalul de încredere

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:

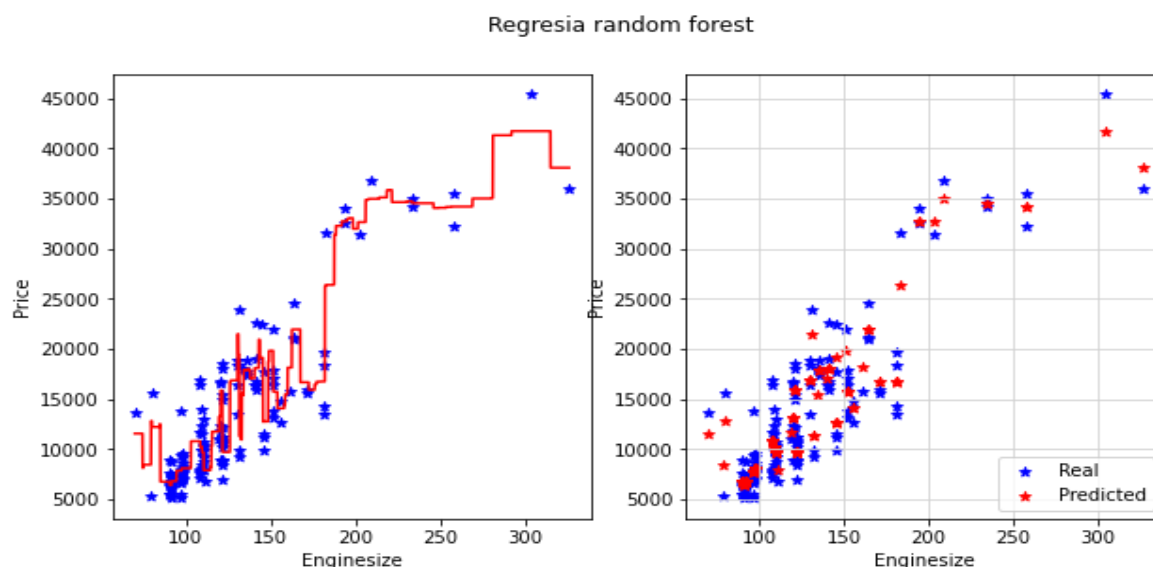


Figura 4.13: Regresia random forest simplă: enginesize → price

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.

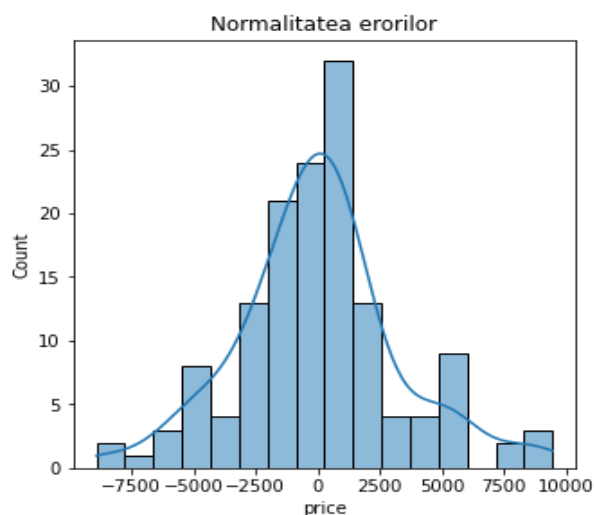


Figura 4.15: Normalitatea erorilor – regresia liniară multiplă

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

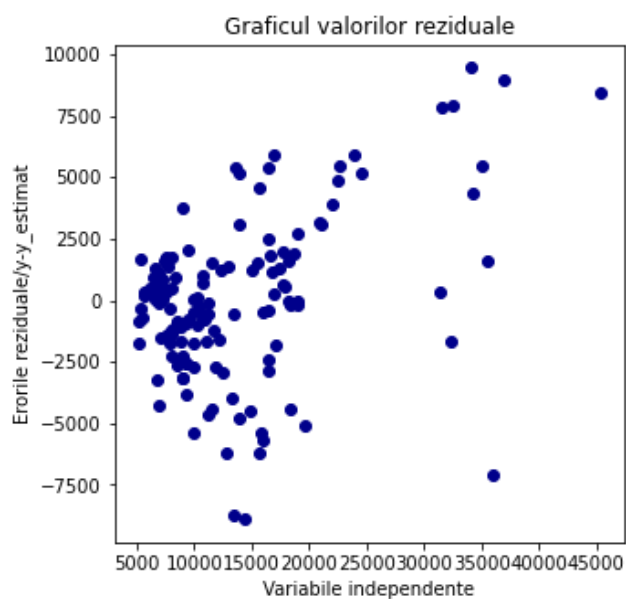


Figura 4.16: Graficul valorilor reziduale – regresia liniară multiplă

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 scade, 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

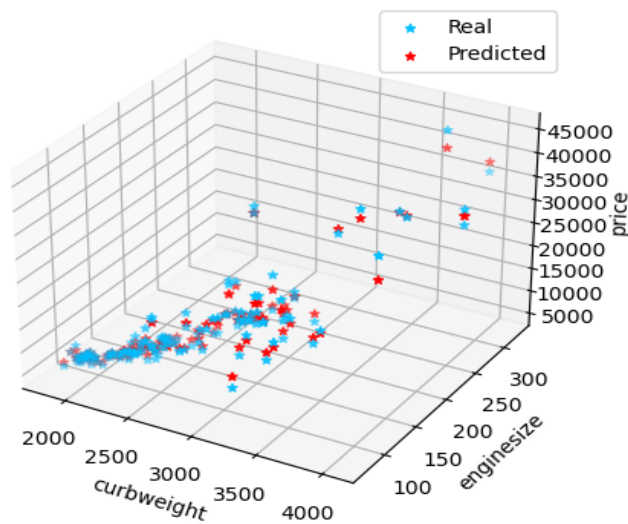


Figura 4.18: Scatterplot regresia random forest cu 2 variabile explicative

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:

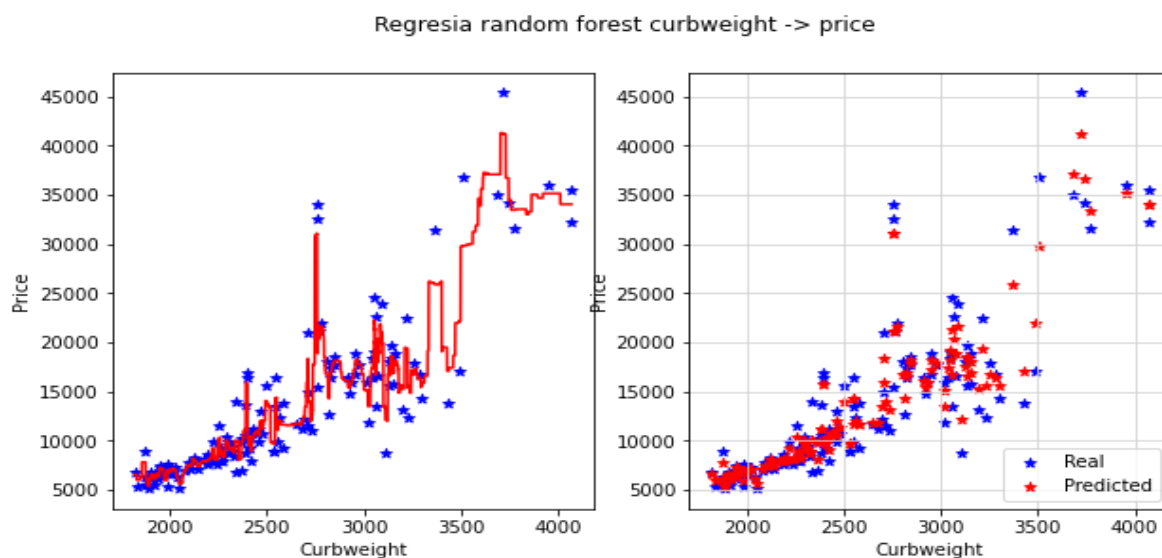


Figura 4.19: Regresia random forest simplă: curbweight \rightarrow price

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:

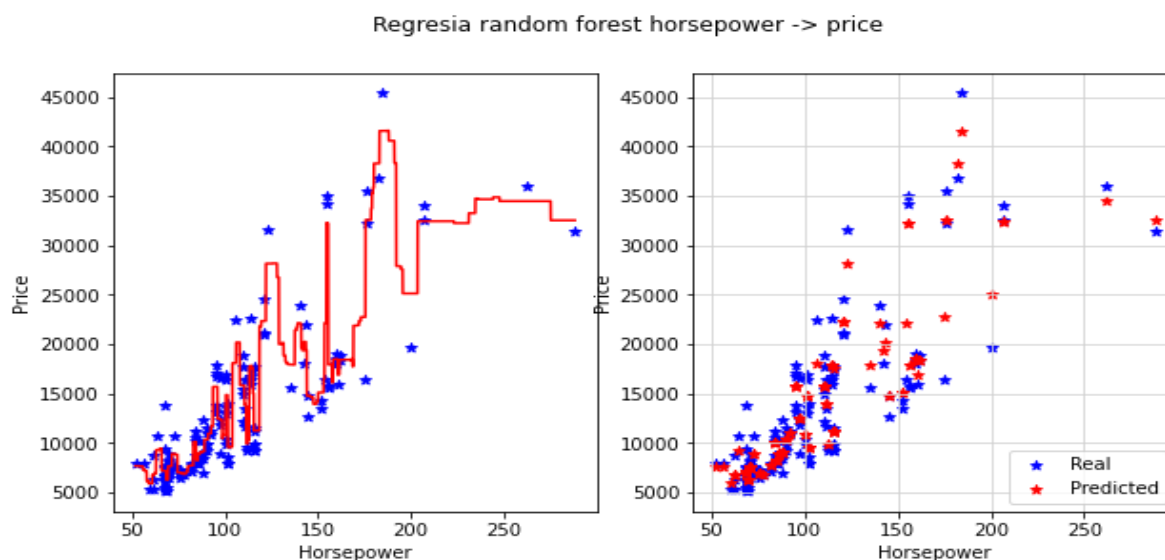


Figura 4.20: Regresia random forest simplă: horsepower \rightarrow price

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:

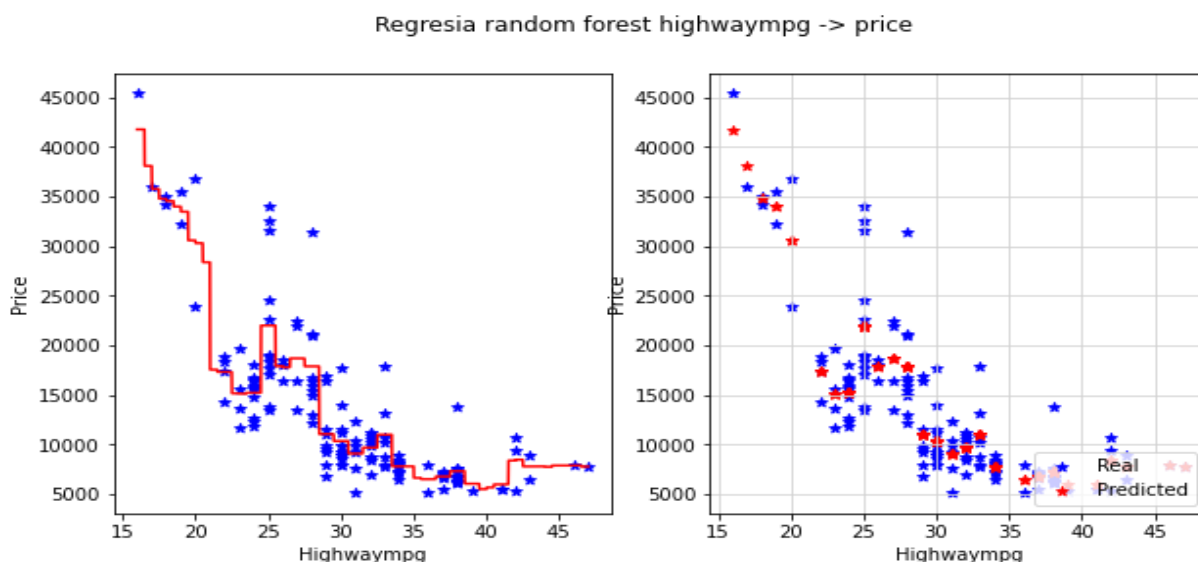


Figura 4.21: Regresia random forest simplă: highwaympg → price

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 (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 (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 (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 (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 (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 (accesare: 15.03.2021).
- [27] D. Zaharie, 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 (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 (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 (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).