



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

| | |
|--|-----------|
| Introducere | 1 |
| 1. Econometrie – studiu de caz..... | 5 |
| 2. Etapele și conținutul unui proiect Data Science..... | 14 |
| 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 |
| 3. Aspecte teoretice și practice ale unui proiect data science bazat pe analiza de regresie | 27 |
| 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 |
| 4. | Studiu de caz – predictia prețului autoturismelor prin analiza de regresie..... | 52 |
| 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 |
| | Bibliografie..... | 77 |

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)

ε = 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

σ = 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).

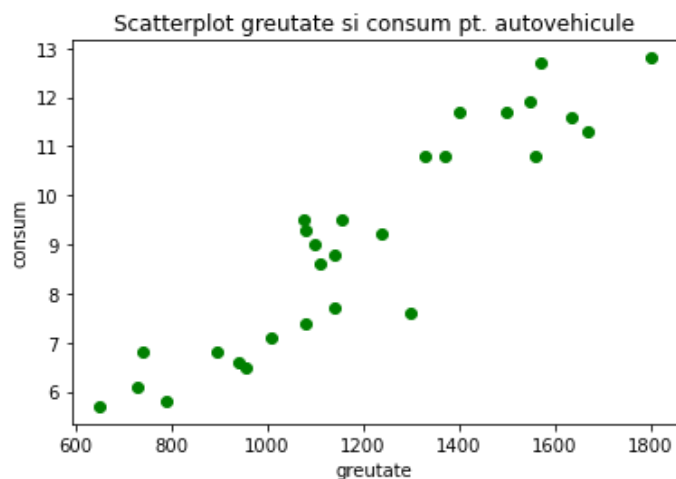


Figura 1.2: Scatterplot greutate si consum pt. autovehicule

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:

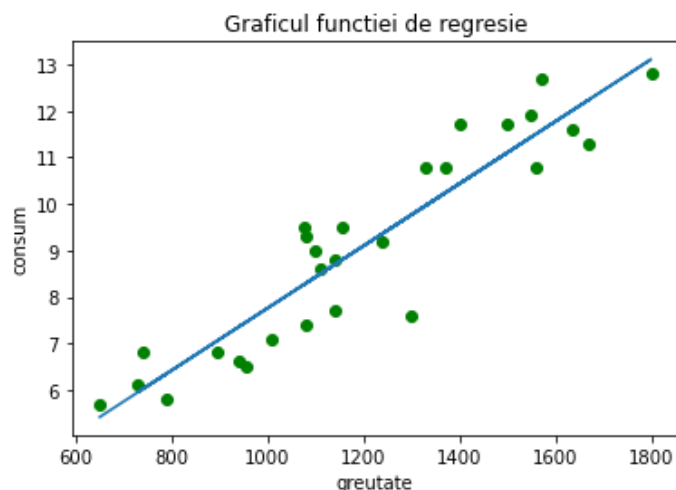


Figura 1.3: Graficul funcției de regresie

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 ε urmează o lege normală de medie zero și variație σ^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 aplătizată

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.

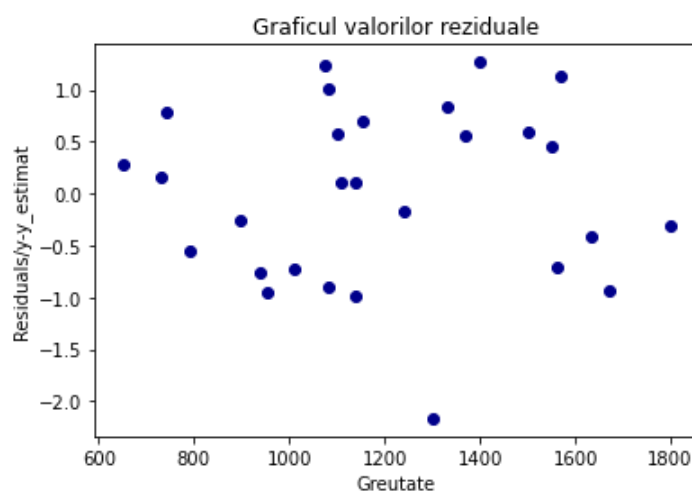


Figura 1.5: Graficul valorilor reziduale

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 α 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.

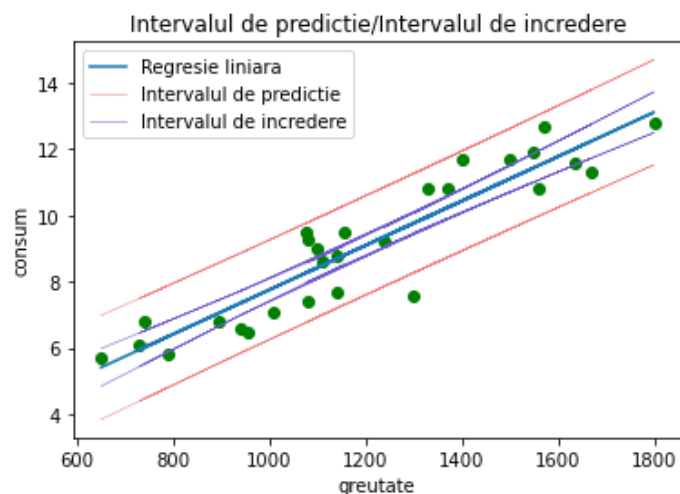


Figura 1.6: Intervalul de predicție/Intervalul de încredere

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.

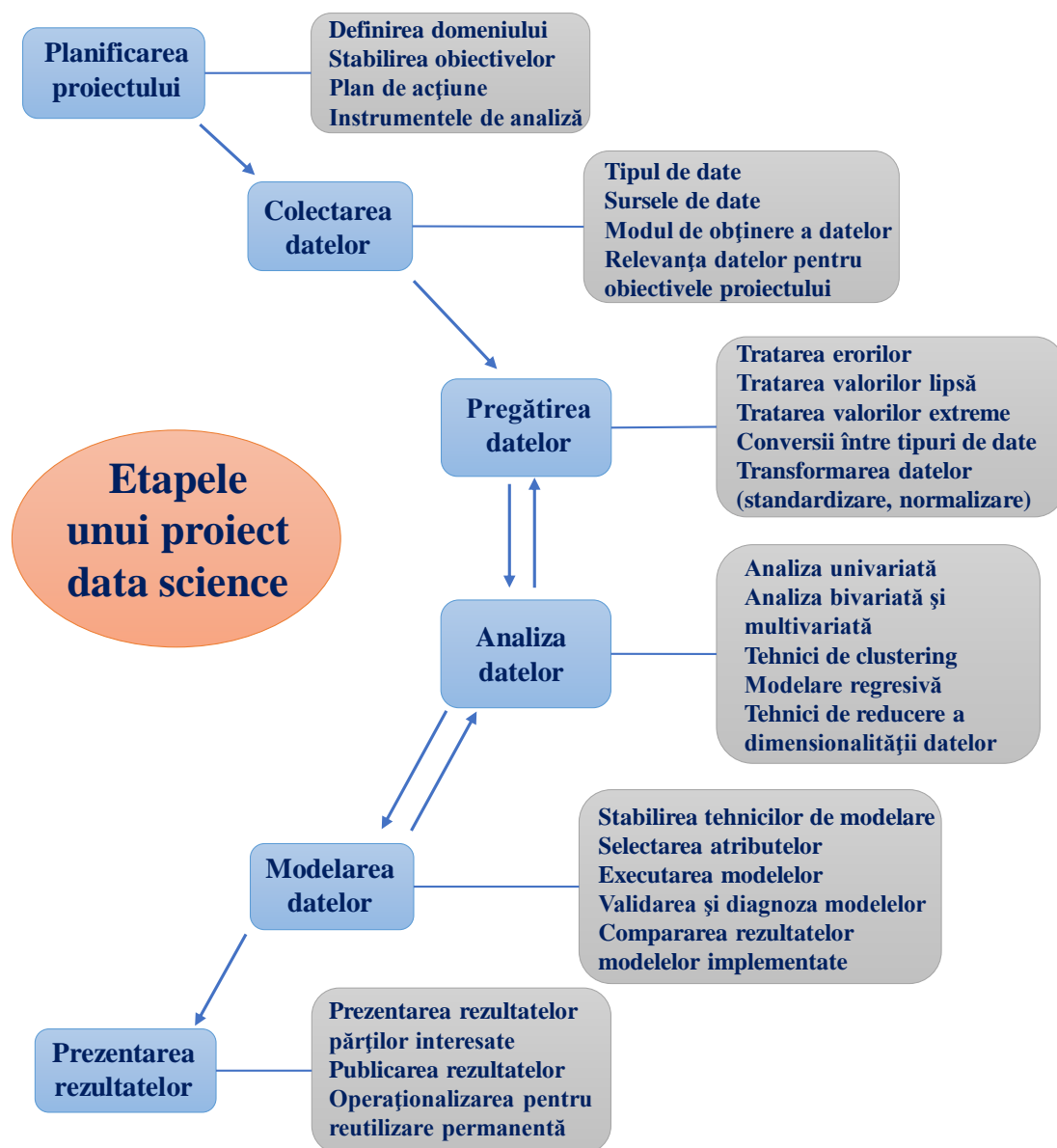


Figura 2.1: Etapele unui proiect data science

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 attribute (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 categoricale), 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, algoritmiile 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/tribute. 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 tehnici 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 ε , 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 - ε 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 ε urmează o lege normală de medie 0 și variație σ^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: $\text{cov}(\varepsilon_i, \varepsilon_j) = 0$; 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]