

Predictia prețurilor mașinilor

Gabriel-Ioan Andreica

May 30, 2024

Contents

0.1	Motivația alegerii bazei de date	2
0.2	Contextul bazei de date și al proiectului	3
0.2.1	Cerinte	3
0.3	Obiective	3
0.4	Aspecte teoretice relevante	4
0.5	Implementarea aspectelor teoretice în cadrul proiectului	6
0.6	Testare și validare	8
0.6.1	Metode de testare	8
0.7	Rezultate ale testării	9
0.8	Rezultate	9
0.9	Concluzii	10

Motivația alegerii bazei de date

La fel ca și majoritatea domeniilor, cel al vanzarilor auto, a suferit schimbări majore în ultimii ani, determinate de evoluția tehnologică de la nivel global. Prin digitalizare, mulți dealeri și-au mutat activitatea în online, permițând clienților să vizioneze, compare sau chiar să achiziționeze un autovehicul, direct de pe internet. Această digitalizare, permite, printre altele, dealerilor auto, să colecteze date pe care să le poată prelucra, cu scopul de a îmbunătăți vânzările. Cu toate că la prima vedere, pare că tehnologia avansată lucrează în folosul firmelor și nu al clientului, aceasta poate fi folosită de ambele părți cu succes. În asemenea măsură în care inteligența artificială este folosită în analiza comportamentelor utilizatorilor, aceasta poate fi folosită și în analizarea unui set mare de mașini, pentru prezicerea unui preț corect. Am ales această baza de date, implicit acest subiect pentru proiectul meu, deoarece în perioada alegerii subiectelor, eram în căutarea unei mașini și am vrut să văd dacă mă pot ajuta de cunoștințele mele tehnice în achiziționarea acestora.

```
C:\Program Files\Facultate > MachineLearning > date > DATE.csv
1  car_id,symboling,fueltype,aspiration,doornumber,carbody,drivewheel,engineplacement,wheelbase,carlength,carwidth,carheight,curbweight,engine_type,cylindernumber,enginesize,
2  1,3,0,0,0,0,0,88.6,168.8,64.1,48.8,2548,0,0,130,0,3.47,2.68,9.0,111,5000,21,27,13495.0,0,0
3  2,3,0,0,0,0,0,88.6,168.8,64.1,48.8,2548,0,0,130,0,3.47,2.68,9.0,111,5000,21,27,13495.0,0,1,0
4  3,1,0,0,0,1,0,94.5,171.2,65.5,52.4,2823,1,1,152,0,2.68,3.47,9.0,154,5000,19,26,16500.0,2,0
5  4,2,0,0,1,2,1,0,99.8,176.6,66.2,54.3,2337,2,0,100,0,3.19,3.4,10.0,102,5500,24,30,13950.0,3,1
6  5,2,0,0,1,2,2,0,99.8,176.6,66.2,54.3,2824,2,2,136,0,3.19,3.4,8.5,110,5500,19,25,17710.0,4,1
7  6,2,0,0,0,2,1,0,99.8,177.3,66.3,53.1,2507,2,2,136,0,3.19,3.4,8.5,110,5500,19,25,15250.0,5,1
8  7,1,0,0,1,2,1,0,105.8,192.7,71.4,55.7,2844,2,2,136,0,3.19,3.4,8.5,110,5500,19,25,17710.0,4,1
9  8,1,0,0,1,3,1,0,105.8,192.7,71.4,55.7,2954,2,2,136,0,3.19,3.4,8.5,110,5500,19,25,18920.0,6,1
10 9,1,0,0,1,1,2,0,105.8,192.7,71.4,55.9,3086,2,2,131,0,3.13,3.4,8.3,140,5500,17,20,23875.0,7,1
11 10,0,0,1,0,1,2,0,99.5,178.2,67.9,52.0,3053,2,2,131,0,3.13,3.4,7.0,160,5500,16,22,17850.167,8,1
12 11,2,0,0,0,2,0,0,101.2,176.8,64.8,54.3,2395,2,0,108,0,3.5,2.0,8.0,101,5800,23,29,16430.0,9,2
13 12,0,0,0,1,2,0,0,101.2,176.8,64.8,54.3,2395,2,0,108,0,3.5,2.0,8.0,101,5800,23,29,16025.0,9,2
14 13,0,0,0,0,2,0,0,101.2,176.8,64.8,54.3,2710,2,1,164,0,3.31,3.19,9.0,121,4250,21,28,20970.0,10,2
15 14,0,0,0,1,2,0,0,101.2,176.8,64.8,54.3,2765,2,1,164,0,3.31,3.19,9.0,121,4250,21,28,21105.0,11,2
16 15,1,0,0,1,2,0,0,103.5,189.0,66.9,55.7,3055,2,1,164,0,3.31,3.19,9.0,121,4250,20,25,24565.0,12,2
17 16,0,0,0,1,2,0,0,103.5,189.0,66.9,55.7,3230,2,1,209,0,3.62,3.39,8.0,182,5400,16,22,30760.0,13,2
18 17,0,0,0,0,2,0,0,103.5,193.8,67.9,53.7,3380,2,1,209,0,3.62,3.39,8.0,182,5400,16,22,41315.0,14,2
19 18,0,0,0,1,2,0,0,110.0,197.0,70.9,56.3,3505,2,1,209,0,3.62,3.39,8.0,182,5400,15,20,36880.0,11,2
20 19,2,0,0,0,1,1,0,88.4,141.1,60.3,53.2,1488,3,3,61,1,2.91,3.03,9.5,48,5100,47,53,5151.0,15,3
21 20,1,0,0,0,1,1,0,94.5,155.9,63.6,52.0,1874,2,0,90,1,3.03,3.11,9.6,70,5400,38,43,6295.0,16,3
22 21,0,0,0,1,2,1,0,94.5,188.0,63.6,52.0,1909,2,0,90,1,3.03,3.11,9.6,70,5400,38,43,6575.0,17,3
23 22,1,0,0,0,1,1,0,93.7,157.3,63.8,50.8,1876,2,0,90,1,2.97,3.23,9.41,68,5500,37,41,5572.0,18,4
24 23,1,0,0,0,1,1,0,93.7,157.3,63.8,50.8,1876,2,0,90,1,2.97,3.23,9.4,68,5500,31,38,6377.0,19,4
25 24,1,0,0,1,1,0,93.7,157.3,63.8,50.8,2128,2,0,98,0,3.03,3.39,7.6,102,5500,24,30,7957.0,20,4
26 25,1,0,0,1,1,0,93.7,157.3,63.8,50.6,1967,2,0,90,1,2.97,3.23,9.4,68,5500,31,38,6229.0,21,4
27 26,1,0,0,1,2,1,0,93.7,157.3,63.8,50.6,1989,2,0,90,1,2.97,3.23,9.4,68,5500,31,38,6692.0,22,4
28 27,1,0,0,1,2,1,0,93.7,157.3,63.8,50.6,1989,2,0,90,1,2.97,3.23,9.4,68,5500,31,38,7609.0,21,4
29 28,1,0,0,1,2,1,0,93.7,157.3,63.8,50.6,2191,2,0,98,0,3.03,3.39,7.6,102,5500,24,30,8558.0,23,4
30 29,1,0,0,1,3,1,0,103.3,174.6,64.6,59.8,2535,2,0,122,1,3.34,3.46,8.5,88,5000,24,30,8921.0,23,4
31 30,3,0,1,0,1,1,0,95.0,172.2,66.2,50.2,2811,2,0,150,2,2.6,3.9,7.0,145,5000,19,24,12664.0,21,4
32 31,2,0,0,0,1,1,0,86.6,144.6,63.0,50.8,1713,2,0,92,3,2.91,3.41,9.6,58,4800,49,54,6470.0,24,5
33 32,2,0,0,0,1,1,0,86.6,144.6,63.0,50.8,1819,2,0,92,3,2.91,3.41,9.2,76,6000,31,38,6855.0,25,5
34 33,1,0,0,0,1,1,0,93.7,150.0,64.0,52.6,1837,2,0,79,3,2.91,3.07,10.1,60,5500,38,42,5390.0,24,5
35 34,1,0,0,0,1,1,0,93.7,150.0,64.0,52.6,1940,2,0,92,3,2.91,3.41,9.2,76,6000,30,34,6529.0,25,5
36 35,1,0,0,0,1,1,0,93.7,150.0,64.0,52.6,1956,2,0,92,3,2.91,3.41,9.2,76,6000,30,34,7129.0,25,5
37 36,0,0,0,1,2,1,0,96.5,163.4,64.0,54.5,2010,2,0,92,3,2.91,3.41,9.2,76,6000,30,34,7295.0,26,5
38 37,0,0,0,1,3,1,0,96.5,157.1,63.9,58.3,2024,2,0,92,3,2.92,3.41,9.2,76,6000,30,34,7295.0,27,5
39 38,0,0,0,0,1,1,0,96.5,167.5,65.2,53.3,2236,2,0,110,3,3.15,3.58,9.0,86,5800,27,33,7895.0,28,5
40 39,0,0,0,0,1,1,0,96.5,167.5,65.2,53.3,2289,2,0,110,3,3.15,3.58,9.0,86,5800,27,33,9095.0,29,5
41 40,0,0,0,0,2,1,0,96.5,175.1,65.2,54.3,2304,2,0,110,3,3.15,3.58,9.0,86,5800,27,33,9045.0,30,5
```

Figure 1: Baza de date pentru test, cu valorile string transformate în numeric.

Contextul bazei de date și al proiectului

Cerinte

Proiectul prezent este un program de Machine Learning care, în baza unui set complex de date, învață să prezică prețul unei mașini. Baza de date cuprinde informații relevante de listă a 205 mașini. Informațiile relevante sunt ID, numele brand-ului, numele modelului, symboling, tipul de combustibil folosit, dacă sunt aspirate sau nu, numărul de uși, tipul de caroserie, tracțiunea, localizarea motorului, wheelbase, lungimea mașinii, lățimea mașinii, înălțimea mașinii, curbweight, tipul de motor, numărul de cilindrii, mărimea motorului, sistemul de combustibil, boreratio, stroke, compressionratio, cai putere, peakrpm, consum în oraș, consum pe autostrada și preț. Cu toate că par multe date, acestea sunt suficiente pentru o predicție de bază a prețului unei mașini. Aplicația nu va putea prezice prețuri pentru mașinile second hand și nu va putea diferenția două mașini identice, dar cu specificații diferite ce țin de confort. Cerințele proiectului sunt să analizeze setul de date, ca mai apoi să extragă din acesta atributele cele mai relevante în contextul prețului unei mașini. În cele din urmă, programul va fi antrenat pentru a putea să fie capabil, ca folosindu-se de acele date care sunt cele mai relevante, să prezică prețul real al unor mașini despre care știe doar specificațiile tehnice, fără preț.

Obiective

Scopul acestui proiect este de a obține în mod automat un preț cât mai apropiat, dacă nu chiar exact, de cel real. Aplicația va funcționa doar pentru mașinile noi. Programul va trebui să învețe să interpreteze corelațiile dintre atribute și prețul mașinilor, iar pe baza acestora, va trebui să pună la dispoziție un preț pe care acesta îl consideră ca fiind cel corect pentru vehiculul respectiv. Mai mult decât atât, aplicația va prezice, cu o acuratețe destul de ridicată, prețuri pentru mașini, altele decât cele cu care a fost antrenată. Și pentru acestea, va trebui să estimeze un preț cu o precizie cât mai ridicată, cu toate că este prima dată când programul întâlnește configurația respectivă de atribute. Această aplicație va putea fi folosită atât de persoane fizice, cât și de persoane juridice care vor să achiziționeze o mașină și nu vor să plătească pentru aceasta mai mult decât aceasta merită.

Aspecte teoretice relevante

Inteligența artificială începe să acapareze toate industriile. Modelele de Machine Learning și Inteligență Artificială sunt intens folosite în analiza comportamentală a utilizatorilor. Pe baza acestei analize, companiile de analitică, au tras anumite concluzii și au venit cu soluții sau îmbunătățiri pentru anumite ramuri ale vânzărilor auto. Câteva ramuri pe care eu le consider relevante și în care este folosită inteligența artificială sunt:

- În analiza predictivă [5] a prețurilor. De la bun început, companiile au vânat vânzările rapide, cu prețuri cât mai bune. Pe măsură ce tehnica a avansat, oamenii au început să culeagă informații despre vânzările pe care le-au făcut. Mai târziu, acele informații au stat la baza unor îmbunătățiri care au scăpat unele companii din faliment. În prezent, modelele de Machine Learning pot analiza datele istorice de vânzări, comportamentul clienților și tendințele de pe piață și informațiile pot fi folosite pentru a prezice comportamentul clienților sau pentru a optimiza strategia de vânzări.
- În personalizarea experienței clienților[2]. Orice persoană tinde să cumpere mai mult atunci când se simte bine. Această senzație de bine poate fi dată prin multe moduri, de la o ambianță plăcută în showroom, până la un angajat care îți vorbește frumos și are răbdare cu tine, în alegerea unei mașini pe site. Prin utilizarea procesării limbajului natural, pot fi dezvoltati chatboți care să răspundă la întrebările clienților și să îi ajute în procesul de achiziționare a unui vehicul, precum ar face-o un angajat real, cu experiență în vânzări, psihologie și relații cu publicul.
- În detectarea fraudelor[4]. În paralel cu dorința de a crește vânzările, implicit banii care intră în conturile companiei, s-a dezvoltat și dorința de minimizare a pierderilor. O bună parte din aceste pierderi este dată de persoane rău intenționate care, profitând de anumite vulnerabilități în procesul de achiziție, reușesc să prejudicieze companiile. Machine Learning-ul poate identifica tipare de comportament fraudulos în procesul de finanțare și de achiziționare a unui vehicul prin analizarea datelor [1] tranzacțiilor și a istoricului clienților.
- În previzionarea cererii[3]. Factorul cel mai important în optimizarea vânzărilor, este cererea și oferta. Companiile s-au prins destul de rapid că dacă vor vinde ceea ce nu dorește nimeni să cumpere, nu vor avea

câștiguri prea mari. Evoluția tehnicii a dat o soluție în acest sens. Învățarea automată poate prevedea cererea pentru anumite mașini.

- În campanii de publicitate[2]. Campaniile de publicitate ale companiilor vânzătoare de mașini nu sunt mereu de succes. Unul dintre motive este și lipsa publicului țintă. Inteligența Artificială a venit cu o soluție și în acest sens. Algoritmii de Machine Learning pot analiza comportamentul online al clienților și pot să le afișeze reclame relevante.

Referințele științifice aferente informațiilor de mai sus:

- Dua, S., Du, X. (2019). Data Mining and Machine Learning in Cyber-security. CRC Press.
- Siegel, E. (2016). Predictive Analytics: The Power to Predict Who Will Click, Buy, Lie, or Die. Wiley.
- Chen, T., Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.
- Xu, R., Wunsch, D. (2005). Clustering. IEEE Transactions on Neural Networks.
- Phua, C., Lee, V., Smith, K., Gayler, R. (2010). A comprehensive survey of data mining-based fraud detection research. Artificial Intelligence Review.
- Kouki, Y., Dalle, O., Zouinkhi, A. (2014). Optimization of Inventory Management with Machine Learning. International Journal of Engineering Business Management.
- Dale, R. (2016). The return of the chatbots. Natural Language Engineering.
- Shmueli, G., Patel, N. R., Bruce, P. C. (2010). Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner. Wiley.
- Jardine, A. K., Lin, D., Banjevic, D. (2006). A review on machinery diagnostics and prognostics implementing condition-based maintenance. Mechanical Systems and Signal Processing.
- Fader, P. S., Hardie, B. G. (2009). Probability Models for Customer-Base Analysis. Journal of Interactive Marketing.

Implementarea aspectelor teoretice în cadrul proiectului

Am acordat atenție sporită proiectării și planificării fiecărui task din cadrul proiectului. Acesta a fost împărțit în X pași. Primul pas l-a reprezentat prelucrarea setului de date. După câteva zile în care am analizat datele și structura bazei de date, am ajuns la concluzia că acestea vor trebui prelucrate pentru a fi utilizabile într-un model de Machine Learning. Prima provocare a acestui proiect a fost transformarea valorilor de tip string în valori numerice.

```
17 def map_non_numeric_values(column):
18     unique_values = column.unique()
19     mapping = {value: index for index, value in enumerate(unique_values)}
20     return column.map(mapping)
21
22 # înlocuirea valorilor non-numerice cu valori numerice pentru fiecare coloană non-numerică
23 for feature in non_numeric_features:
24     df[feature] = map_non_numeric_values(df[feature])
```

Figure 2: Funcția care transformă atributele string în atribute numerice

Pasul următor l-a reprezentat calcularea Gini Index pentru setul meu de date. Gini Index este utilizat în Machine Learning ca o măsură de impuritate pentru a calcula calitatea unui split. În contextul proiectului, este o metrică foarte importantă, pentru alegerea atributelor în nodurile de decizie. După ce am calculat indicele pentru toate campurile din baza mea de date, am ales primele trei cele mai bune caracteristici din punct de vedere al Gini Index. Aceste trei atribute sunt: curbweight, carmodel și carlength.

```
✓ Primii 3 cei mai buni atribuți în funcție de Gini Index:
1. curbweight: 0.0976
2. carmodel: 0.4146
3. carlength: 0.5966
```

Figure 3: Primele 3 atribute din Gini Index

Ulterior, pentru a avea mai multe atribute de calitate pe baza cărora să se efectueze antrenarea modelului, am calculat și Information Quantity. Information Quantity este o măsură a impurității sau dezordinii dintr-un set de date. Această metrică clasifică atributele din setul de date, după relevanța pe care o au în contextul modificării unui atribut vizat. Acest atribut, în proiectul de față este prețul, iar cele trei cele mai bune atribute din punct de vedere al Information Quantity, sunt: enginelocation, aspiration și cylindernumber.

```
Cele mai bune 3 atribute în funcție de cantitatea de informație:  
engineLocation: 0.0000  
aspiration: -0.0098  
cylindernumber: -0.0098
```

Figure 4: Primele 3 atribute din Information Quantity

Pasul următor în proiectul de față, a fost reprezentat de alegerea a 3 modele de Machine Learning pe care ulterior setul de date să fie testat. Cele trei modele alese au fost Random Forest, Decision Tree și Linear Regression. Am efectuat pentru fiecare din cele trei modele o testare cu un set de date care conținea 205 mașini. Pașii precedenți au influențat în mod direct aceste teste. În toate cele trei modele, au fost folosite cele 3 atribute care sunt cele mai bune din punct de vedere al indicelui Gini și cele trei atribute care sunt cele mai bune din punct de vedere al Information Quantity.

Primul model testat a fost Random Forest. Acesta este un algoritm utilizat pentru clasificare și regresie care se bazează pe construirea mai multor arbori decizionali, cu scopul de a obține o precizie cât mai mare. Este avantajos pentru eliminarea problemei overfitting-ului.

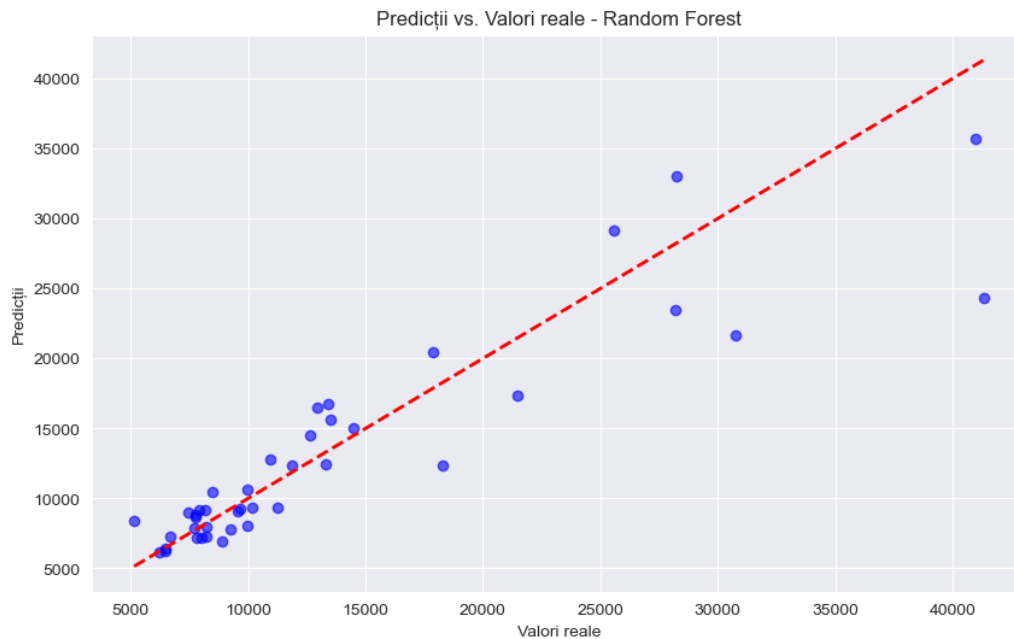


Figure 5: Graficul pentru Random Forest

Cel de al doilea model testat a fost Decision Tree. Asemenea primului model testat, Decision Tree este folosit la clasificare și regresie. Marea diferență între cele două este dată de modul de a funcționa a acestora. De-

cision Tree, spre deosebire de Random Forest, se folosește de o structură de tip arbore în luarea deciziilor.

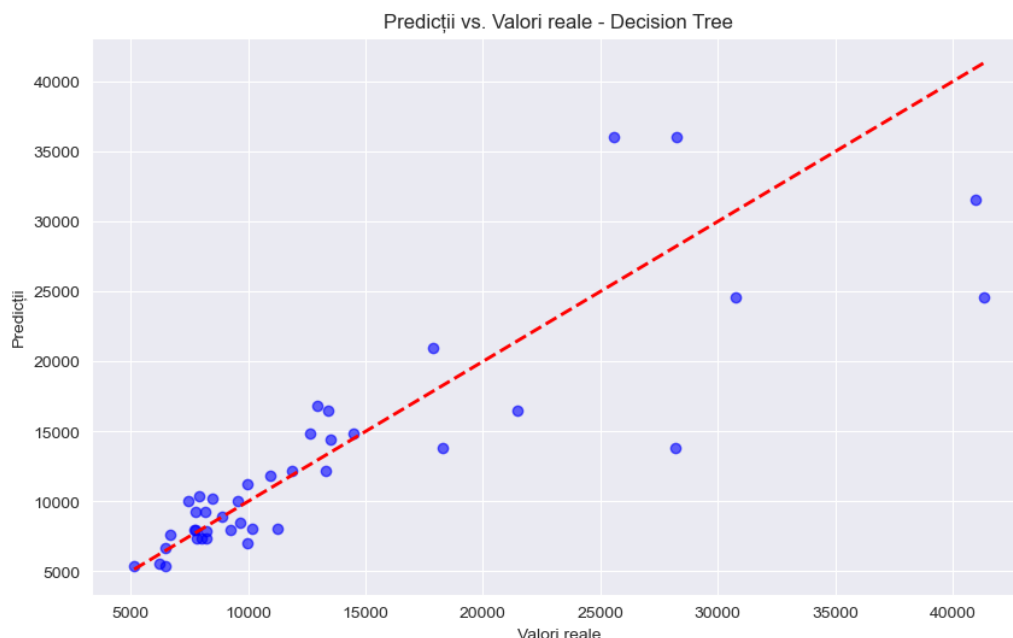


Figure 6: Graficul pentru Decision Tree

Cel din urmă model testat a fost Linear Regression. Acesta este un algoritm utilizat în modelarea relației dintre un atribut cheie și alte câteva atribute.

Dintre toate cele trei modele testate, precizia cea mai bună pe care programul a obținut-o, a fost pe modelul Random Forest, prin urmare, acesta a fost modelul final, folosit în acest proiect.

Testare și validare

Metode de testare

În ceea ce privește testarea modelului acesta de Machine Learning, am urmărit obținerea unei precizii cât mai apropiate de cea pe care am obținut-o pe setul de date de antrenament. Am folosit două variante de testare. În prima, am împărțit setul inițial de date în 70% date de antrenament și 30% date de testare. A doua metoda de testare folosită a fost cea manuală. Am introdus manual într-un nou fișier CSV, specificațiile a 15 mașini, cu diferențe față de

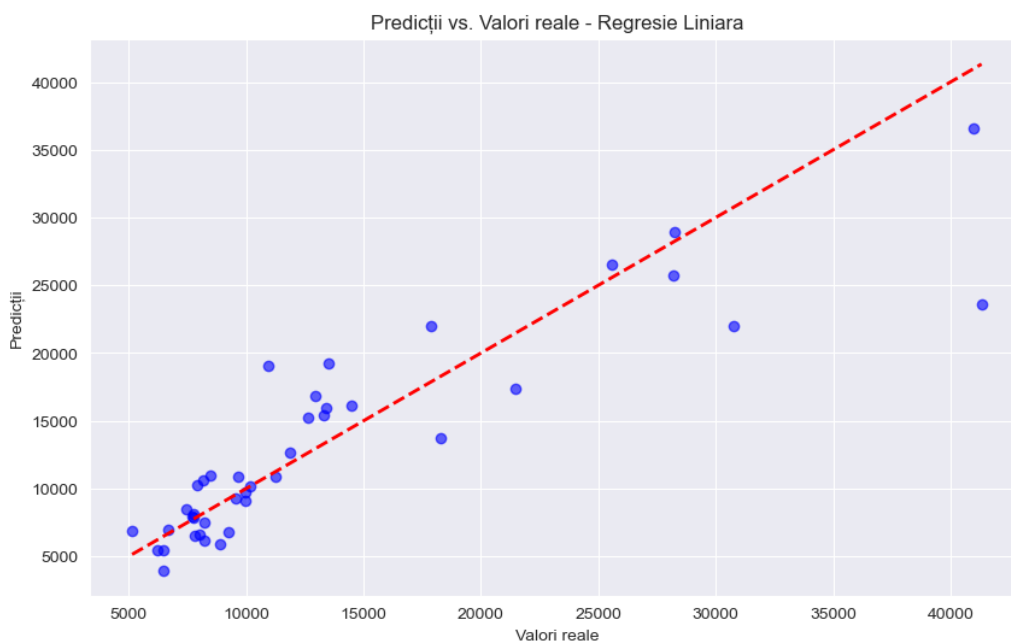


Figure 7: Graficul pentru regresie liniară

cele introduse în setul de date de antrenament. Din nou, atributele au fost transformate din string în numeric, unde a fost cazul și s-a început testarea.

Rezultate ale testării

Rezultatele testării au fost conforme cu așteptările. Atât pentru prima modalitate de testare, cât și pentru a doua, au fost obținute rezultate similare, o precizie în jur de 80-90%. Precizia mai bună a avut-o cea de a doua variantă testată, adică 89.22%. Am observat, de asemenea, că, pe măsură ce prețul vehiculului crește, predicția acestuia este supusă la o eroare mai mare din punct de vedere al diferenței de bani. Acest lucru vine din marja de eroare pe care o mașină de \$1000 o poate avea, în contrast cu una de \$100000. Cu toate acestea, diferențele sunt din nou asemănătoare, din punct de vedere procentual, raportat la prețul mașinii.

Rezultate

Rezultatele finale ale proiectului sunt conforme cu așteptările pe care le-am impus la începerea acestuia. În ceea ce privește prelucrarea datelor și

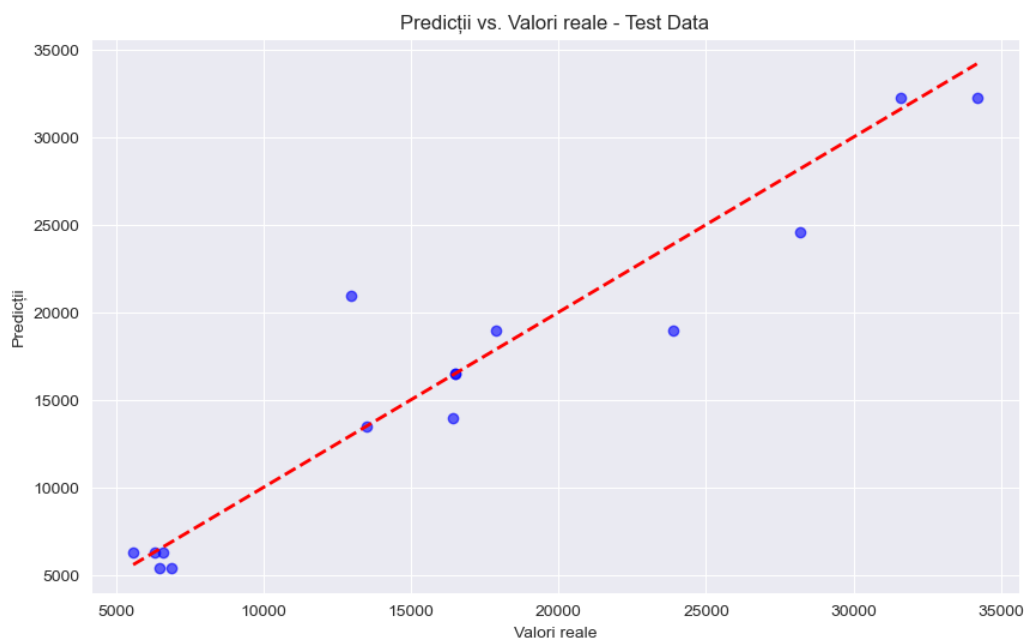


Figure 8: Graficul datelor de testare

antrenarea modelului, metricele rezultate din prelucrarea setului de date au fost de calitate și relevante pentru rezultatul final al proiectului. Mai departe, cele trei modele testate, au fost cele mai bune variante posibile pentru problema dată. Din aceste trei variante, cea mai bună s-a dovedit a fi modelul Random Forest. Chiar și în această situație mai sunt loc de îmbunătățiri. Acestea ar putea fi făcute în partea de prelucrare a setului de date. În loc de transformarea datelor din numeric în string, s-ar putea crea câmpuri noi pentru toate acele date, iar conținutul acestora să fie de tip 0 sau 1, în funcție de valoarea de adevăr a acestora.

Concluzii

Concluziile în urma acestui proiect sunt clare. În ceea ce privește calitățile și pregătirea mea tehnică, am învățat noțiunile introductive din domeniul Machine Learning și Data Analysis. Aceste noțiuni pot fi baza unei cariere în adevăratul sens al cuvântului, pe mai departe, ținând cont de amploarea pe care aceste două domenii o iau. Pe de alta parte, în ceea ce proiectul în sine, rezultatele obținute au fost bune. Ținând cont de acest aspect, putem cataloga proiectul ca fiind de succes. Bineînțeles, programul are loc de mai bine. Acesta ar putea fi îmbunătățit pe viitor. Un exemplu de funcționalitate

pe care programul ar putea-o avea, ar fi aceea de a prezice prețurile și pentru mașinile second-hand, ținând cont de alte atribute, cum ar fi kilometrajul. O alta îmbunătățire ar fi ca programul să țină cont și de atribute ce țin de aspectul sau confortul mașinii, precum opțiunile interioare.