

Polovni telefoni

Ime Prezime, index, e-mail

UVOD

Kupovina i prodaja korišćenih telefona i tableta nekada su se obavljale na samo nekoliko online tržišnih sajtova. Međutim, tržište polovnih i obnovljenih uređaja značajno je poraslo tokom poslednje decenije. Ovaj rast se može pripisati porastu potražnje za polovnim telefonima i tabletima koji nude značajne uštede u poređenju sa novim modelima. Obnovljeni i korišćeni uređaji i dalje pružaju ekonomične alternative kako potrošačima tako i poslovnim korisnicima koji žele da uštede novac prilikom kupovine. Postoje mnoge druge prednosti koje su povezane sa tržištem korišćenih uređaja. Polovni i obnovljeni uređaji mogu se prodavati sa garancijama, a mogu se osigurati i dokazom o kupovini. Treće strane prodavci/platforme, poput A1, Telenora, MTS-a, itd., nude atraktivne ponude za obnovljene uređaje. Povećanje trajnosti uređaja kroz trgovinu polovnih uređaja takođe smanjuje njihov ekološki uticaj i pomaže u recikliranju i smanjenju otpada.

I. BAZA PODATAKA

Baza podataka sadrži 15 atributa(kolona) i 3454 podataka(redova)(baza podataka je dimenzija 3454x15). Ti atributi su device_name (ime brenda), os (operativni sistem), screen_size (velicina ekrana u cm), 4g, 5g, front_camera_mp (rezolucija prednje kamere u mega pikselima), rear_camera_mp (rezolucija zadnje kamere u mega pikselima), internal_memory (interna memorija u gigabajtima), ram (količina operativne memorije u gigabajtima), battery (energetski kapacitet baterije u mAh), weight (težina u gramima), release_year (godina proizvodnje), days_used (koliko je uređaj dana korišćen), normalized_new_price (normalizovana cena novog uređaja), normalized_used_price (cena polovnog uređaja).

Cilj nam je da odredimo cenu polovnog uređaja (zavisna promenljiva) koristeći neki od regresionih modela.

II. ANALIZA PODATAKA

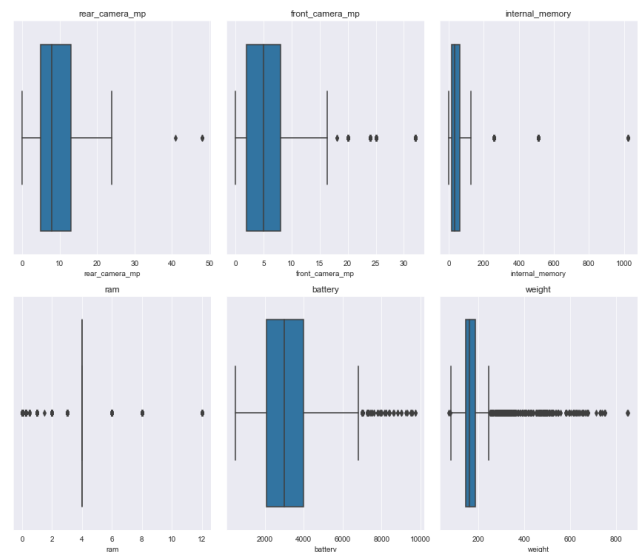
A. Tip podataka

Atributi kategorickog tipa koji su prisutni u nasoj bazi su 4g, 5g, ime brenda i operativni sistem. Ostali atributi su numerckog tipa. Atributi 4g i 5g imaju 2 moguće vrednosti (da ili ne), pa su njihove vrednosti zamenjene sa 0 i 1 respektivno. Za attribute ime brenda i operativni sistem smo primenili 'one hot encoding' zato što ne postoji ordinalni odnos izmedju njihovih vrednosti, pa 'integer encoding' nije dovoljan. Ime brenda ima 34 unikatne vrednosti, stoga imaćemo 34 nove kolone za svaku marku telefona koja je prisutna u nasoj bazi. Operativni

sistem ima 4 unikatne vrednosti (android, ios, windows, other), stoga imaćemo 4 nove kolone koje će predstavljati odgovarajući operativni sistem telefona. Posle prevodjenja svih kategorickih obeležja u numericka, dimenzija baze je 3454x48.

B. Nedostajuci podaci

Podaci koji nedostaju su prednja kamera (5.18% podataka), zadnja kamera (0.057% podataka), interna memorija (0.11% podataka), ram memorija (0.11% podataka), baterija (0.17% podataka). Izbacivanje ovih podataka ne bi bilo ispravno rešenje posto bismo imali manje podataka, i potencijalno bismo izgubili značajne informacije. Da bismo odlučili da li ćemo podatke zameniti medianom ili srednjom vrednošću moramo se prvo upoznati sa distribucijom podataka i utvrditi da li postoje outlieri. Na osnovu slike 1 možemo utvrditi da outlieri postoje. Takođe, distribucija ovih podataka je asimetrična. Najbolja praksa je da kada su outlieri prisutni nedostajuće podatke zamenimo sa medianom, jer je mediana robustnija i manje osetljiva na outlieri nego srednja vrednost. Nedostajuće vrednosti svakog atributa smo zamenili, tako što smo svaki atribut grupisali po njemu respektivnom brendu i odredili medianu te grupe. Atribut zadnja kamera je i dalje imala nedostajuće vrednosti, jer za brend 'Celkon' ne postoje podaci o zadnjoj kameri. Preostale nedostajuće podatke za atribut zadnja kamera smo zamenili uopštenom medianom.



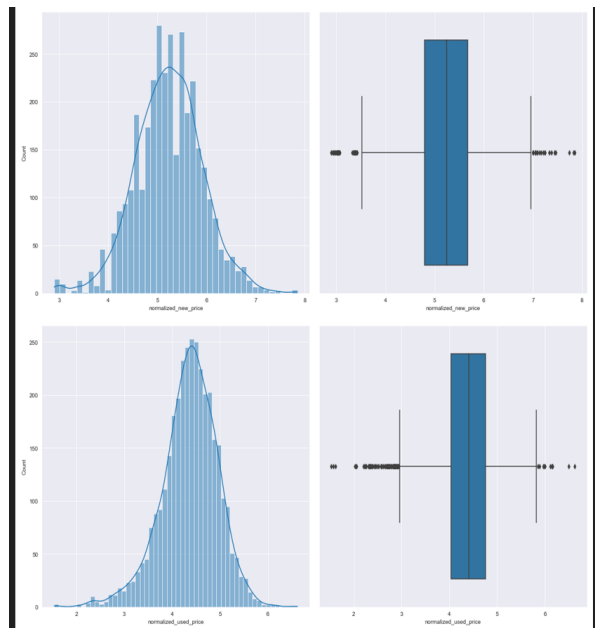
Slika 1. Boxplot sa outlierima za attribute zadnja kamera, prednja kamera, interna memorija, ram memorija, baterija i težina.

C. Outlieri i distribucija podataka

Atribut smo predstavili funkcijom boxplot, gde smo uvideli njihove interkvartalne opsege, medianu i outlieri. Takođe, smo za svaki atribut odredili njegovu distribuciju. Na slici 1 su već prikazani interkvartalni opsezi, mediana i outlieri za nedostajuće vrednosti.

Na slici 2 vidimo da normalizovana cena korišćenih uređaja i normalizovana cena novih uređaja imaju normalnu (simetričnu) distribuciju.

Ostali atributi imaju asimetričnu distribuciju. Na osnovu izracunatih interkvartalnih opsega outlieri su izbaceni kako bi povećali linernost u nasim podacima i izvršili bolju predikciju.



Slika 2. Hisplot sa distribucijom i boxplot sa outlierima za attribute normalizovana cena novih uređaja i normalizovana cena koriscenih uređaja.

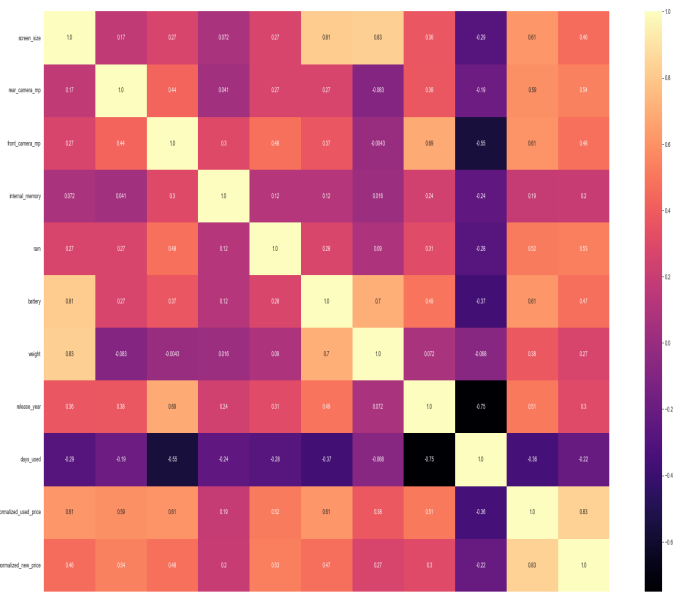
D. Korelacija atributa

Stepen korelacije se meri pomoću koeficijenta korelacije, koji može da varira od -1 do 1. Koeficijent korelacije od -1 označava jaku negativnu korelaciju, koeficijent od 0 označava da nema korelacije, a koeficijent od 1 označava jaku pozitivnu korelaciju. Na slici 3. je prikazana heatmapa pomocu koje mozemo utvrditi koeficijente korelacije izmedju nezavisnih atributa i zavisne promenljive (normalizovana cena koriscenih uređaja). Najveci koeficijent korelacije sa zavisnom promenljivom ima normalizovana cena novih uređaja (0.83). Posle njega najveće pozitivne koeficijente korelacije imaju velicina ekrana (0.61), prednja kamera (0.61), kapacitet baterije (0.61), zadnja kamera (0.59), godina proizvodnje (0.51). Jedini atribut koji ima negativnu korelaciju je koliko je dana telefon koriscen (-0.36). Na slici 4. je prikazana korelacija za attribute normalizovana cena novih uređaja i normalizovana cena koriscenih uređaja. Vidimo vrlo dobru linearnu vezu izmedju ova dva atributa. Kada se cena novog telefona povećava, cena polovnog telefona se linerno povećava. Kod drugih atributa ne mozemo uociti dobru linernu povezanost sa normalizovanom cenom koriscenih uređaja (primer slika 5.).

E. Generalni zaključci o podacima

Samsung je bio najčešće ponovo korišćen telefon, odmah posle brendova koji se kategorizuju kao 'Ostali'. Najviše ljudi koristi Android operativni sistem, a najmanje njih Windows (slika 7.).

Veličina ekrana je uglavnom između 10 i 15 cm-a. Najveću cenu imaju telefoni marke 'One Plus' (pa slede telefoni marke 'Apple') dok najmanju cenu imaju telefoni marke 'Ceikon' (slika 6.). Korisnici 4G i 5G mreže su češći. Telefoni bez 4g imaju nize cene (slika 8). Raspon megapiksela za zadnju kameru uglavnom je između 5-15 mega piksela. Raspon megapiksela za prednju kameru uglavnom je između 0-10 mega piksela. Interni memorija uglavnom iznosi između 0-100 GB. RAM memorija telefona uglavnom iznosi između 3-5 gb. Baterija traje uglavnom oko 3000 mAh. Težina telefona je oko 150-210 grama. Većina telefona je objavljena između 2013-2015 godine. Ljudi su koristili telefone u proseku između 600-800 dana, što je otprilike 2-2,5 godine.



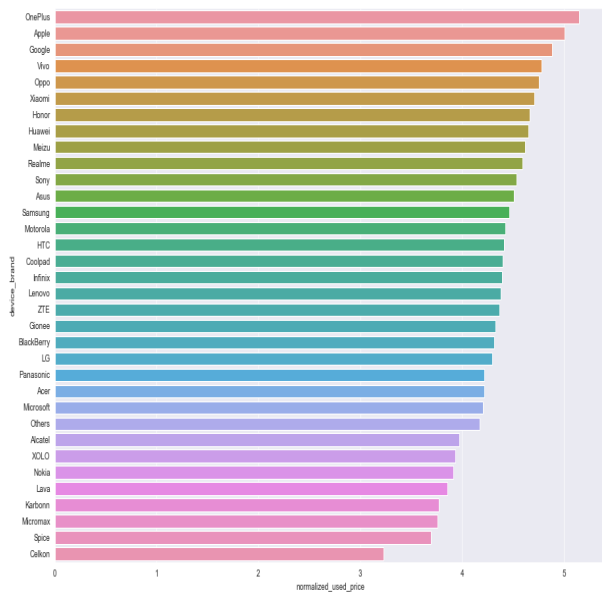
Slika 3. Heatmapa sa koeficientima korelacije izmedju atributa



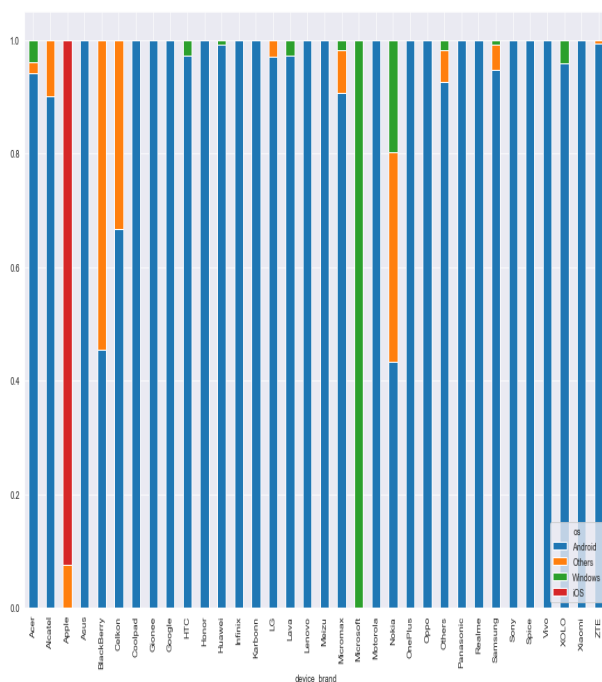
Slika 4. Korelacija izmedju cene novoh telefona i cenr koriscenog telefona



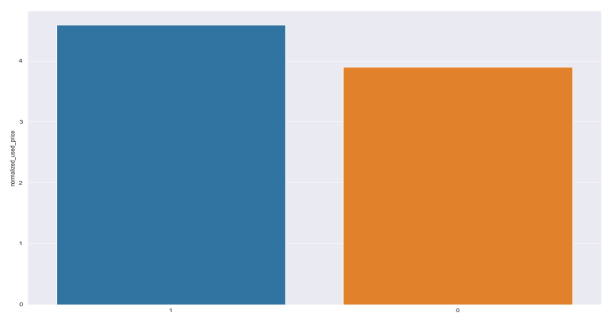
Slika 5. Korelacija izmedju tezindr telefona i cene koriscenog telefona



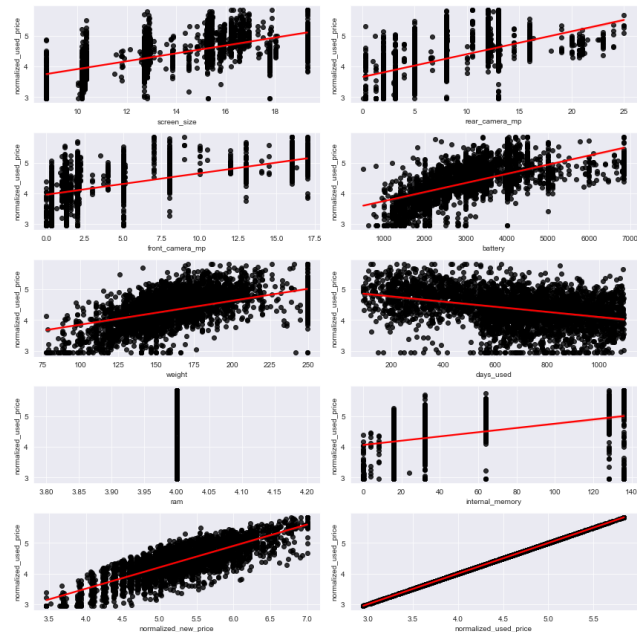
Slika 6. Cena koriscenih uredjaja i naziv brenda



Slika 7. Operativni sistem i brend telefona



Slika 8. 4g(plavi plot telefon podrzava 4g, narandzasti telefon ne podrzava 4g) i cena koriscenih telefona



Slika 9. Regresione linije za numericka obelezja

III. Regresija

Na osnovu slike 9. mozemo zakljuciti da dani koriscenja uredjaja imaju opadajucu(negativnu) regresionu liniju, dok ostala obelezja imaju rastucu(pozitivnu) regresionu liniju.

Kako bi se izvršilo predviđanje normalizovane cene polovnih telefona obuceno je 5 modela linerne regresije. Definisane su mere uspesnosti modela: srednja kvadratna greska(mse), koren srednje kvadratne greske(rmse), srednja apsolutna greska(mae) i R^2 . Izvršena je podela podataka na 70% podataka za treniranje modela, 15% podataka za validaciju, dok je preostalih 15% korisceno za test skup i korišćeno tek za estimaciju performansi konačnog modela, o čemu će biti reči u nastavku. Za svaki od isterniranih modela je izvršena evaluacija nad validacionim skupom, i izabran je model koji daje najbolje rezultate. Pre procesiranja podataka izvršena je standardizacija obelezja, da bi se skalirale i trasformisale karekteristike na zajednicki opseg(jedna karakteristika sa vecom skalom moze dominirati nad ostalim karakteristikama) i da bi distribucija karakteristika bila vise normalizovana(Gaussian distribusion).

Prvi model regresije koji je isteran je osnovni model linearne regresije bez selekcije obelezja. Ovaj najosnovniji model je dao odlicne rezultate(slika 10. nulti model). Nakon primene linearnog regresijskog modela, izvršena je regularizacija modela korišćenjem Lasso i Ridge regresije bez selekcije obelezja. Unakrsna validacija je korišćena kako bi se izabrao optimalni parametar alpha, koji kontroliše jačinu regularizacije. Kroz analizu rezultata, primećeno je da je Lasso regresija dala bolje rezultate u odnosu na Ridge regresiju(slika 10. prvi i drugi model). Unakrsna validacija je omogućila odabir optimalne vrednosti parametra alpha, što je rezultiralo poboljšanjem performanse modela u odnosu na linearnu regresiju. Ovo ukazuje na efikasnost Lasso i Ridge regresije u smanjivanju overfittinga i poboljšanju generalizacije modela u slučajevima kada postoji više ulaznih karakteristika ili visoka multikolinearnost.

Unakrsna validacija je pokazala da je najbolja vrednost alpha parametra za Lasso regresiju 0.001, a za Ridge regresiju 10. Nakon primene regularizacije, isprobani su i drugi modeli, kao što su Decision Tree i Random Forest bez selekcije obelezja. Međutim, Decision Tree model je pokazao loše rezultate u poređenju sa ostalim modelima, što ukazuje na njegovu slabiju sposobnost za generalizaciju na novim podacima (slika 10. model tri). Sa druge strane, Random Forest model je pokazao slično dobre rezultate kao Lasso regresija (slika 10. model četiri). Dodatno, izvedena je analiza "feature importance" u Random Forest modelu, koja je omogućila identifikaciju značajnih karakteristika (atributa) koje su doprinele predviđanju ciljne promenljive (slika 11.). Ovaj pristup je koristan u razumevanju uticaja pojedinačnih atributa na modeliranje. Zaključeno je da karakteristike koje dominiraju su u opadajućem redosledu: normalizovana cena novog uređaja (najveći "feature importance"), veličina ekrana, interna memorija, zadnja kamera, baterija, težina, dani koriscenja uređaja, prednja kamera. Random Forest model se sastoji od više stabala odlucivanja koji se treniraju na razlicitim podskupovima podataka. Konacna predikcija se dobija agregacijom predikcija svakog stabla, sto cini ovaj model stabilnim i otpornim na overfitting. Hiperparametri koji su odabrani za Random forest regresiju dobijeni su unakrsnom validacijom i vrednosti tih parametara variraju, ali najbolje su se pokazale vrednosti: "n_estimators"=150, "max_depth"=10, "min_samples_split"=2.

U nastavku je eksperimentisano sa redukcijom dimenzionalnosti primenom PCA algoritma. Na slici 15. je predstavljena scree plot metoda objasnjenja varijanse i broja komponenti, koji je sličan elbow metodi, i analizira grafikon padanja objasnjenjene varijanse sa svakom novom komponentom. Cilj je naci lakat, tj. broj komponenti koji se nalazi pre 'ravnog' dela grafa. U našem slučaju optimalan broj komponenti je 8. Ako ne specificiramo broj komponenti, primeni ce se granicni uslov, tj. da broj komponenti bude onaj koji zadržava 95% varijabilnosti uzoraka. U našem slučaju se boje pokazao defaultni metod nego prosledjivanje broja komponenti kao parametar pca algoritma. Performanse svakog od regresionih modela gde je koriscen pca algoritam se nalaze na slici 10 (modeli 5-9).

Nakon isprobavanja različitih metoda za redukciju dimenzionalnosti, analizirali smo feature importance i mutual information kako bi identifikovali attribute koji imaju najveću korelaciju sa zavisnom promenljivom. Feature importance se odnosi na ocenjivanje značaja svakog atributa u modelu, dok mutual information meri međuzavisnost između atributa i ciljne promenljive (slika 12.). Korišćenjem ovih metoda, izdvojeni su atributi koji su najviše korelisani sa zavisnom promenljivom, koje smo koristili u daljoj analizi podataka. Identifikacija ključnih atributa može biti od velike važnosti u razumevanju suštine problema i optimizaciji modela, kao i u smanjenju dimenzionalnosti podataka na suštinske informacije koje su relevantne za predikciju ciljne promenljive. Eksperimentisano je sa vrednostima k prilikom selekcije broja obelezja i najbolji rezultat je bio pri

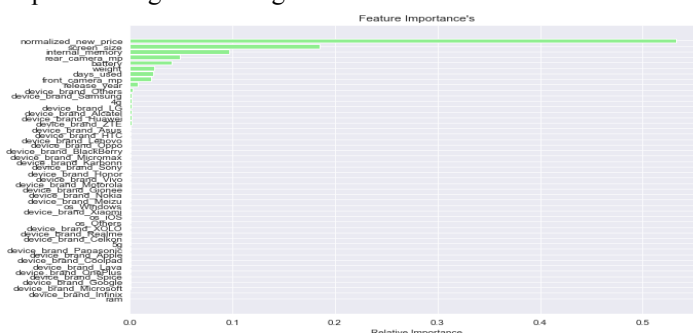
selekciji 8 najbitnijih obelezja. Nakon izdvajanja 8 najbitnijih atributa na osnovu feature importance i mutual information, ponovili smo obuku već ranije koriscenih regresionih modela. Prilikom smanjenja broja atributa, redukcija dimenzionalnosti podataka je rezultirala poboljšanjem performansi Ridge, Lasso i Linear regresionih modela u poređenju sa prethodno korišćenim modelima koji su koristili sve attribute (slika 10. modeli 5-9).

Ovo ukazuje na to da su izdvojeni atributi bili dovoljni za efikasniju predikciju ciljne promenljive u poređenju sa upotrebom svih dostupnih atributa. Na slici 10. su predstavljeni regresioni modeli sa izvojenih top 8 featura (modeli 10-14).

Linearna, Ridge i Lasso regresija sa izvojenih top 8 featura imaju isti skor i greske. Sva tri modela su jednako efektivni u predviđanju zavisne varijable. Međutim, razlike u njihovim pristupima regularizaciji i selekciji funkcija mogu dovesti do razlika u njihovoj interpretabilnosti i generalizabilnosti. U nastavku cemo izvršiti evaluaciju nad test skupom ova tri modela i odlucimo koji je najbolji model.

# models	R2	RMSE	MSE	MAE
0 Linear regression all features	0.8329	0.2226	0.0496	0.1829
1 Ridge regression all features	0.8329	0.2226	0.0496	0.1828
2 Lasso regression all features	0.8334	0.2223	0.0494	0.1829
3 Decision Tree regression all features	0.6864	0.3050	0.0930	0.2370
4 Random Forest regression all features	0.8338	0.2220	0.0493	0.1801
5 PCA Linear regression	0.8329	0.2226	0.0496	0.1829
6 PCA Ridge regression	0.8337	0.2221	0.0493	0.1824
7 PCA Lasso regression	0.7146	0.2909	0.0846	0.2306
8 PCA Decision Tree regression	0.3451	0.4407	0.1943	0.3336
9 PCA Random Forest regression	0.6313	0.3307	0.1094	0.2530
10 Linear regression top 8 features	0.8346	0.2215	0.0491	0.1835
11 Ridge regression top 8 features	0.8346	0.2215	0.0491	0.1835
12 Lasso regression top 8 features	0.8346	0.2215	0.0491	0.1835
13 Random Forest regression top 8 features	0.8320	0.2232	0.0498	0.1786
14 Decision Tree regression top 8 features	0.6836	0.3063	0.0938	0.2376

Slika 10. evaluacija performansi nad validacionim skupom svakog koriscenog modela



Slika 11. feature importance

#	MI Scores
normalized_new_price	0.603288
screen_size	0.475478
battery	0.418283
front_camera_mp	0.409601
rear_camera_mp	0.363841
internal_memory	0.352839
weight	0.317154
release_year	0.205303
4g	0.181775

Slika 12. mutual information

IV. ZAKLJUČAK

Nakon testiranja Ridge regresionog modela sa top 8 featura na test podacima, dobili smo sledeće rezultate: R^2 skor za test podatke je 0.8417, RMSE (srednja kvadratna greška) je 0.2274, MSE (srednja kvadratna greška) je 0.0517, a MAE (srednja apsolutna greška) je 0.1777.

Nakon testiranja Linearne regresije sa top 8 featura na test podacima, dobili smo sledeće rezultate: R^2 skor za test podatke je 0.8419, RMSE (srednja kvadratna greška) je 0.2273, MSE (srednja kvadratna greška) je 0.0516, a MAE (srednja apsolutna greška) je 0.1776.

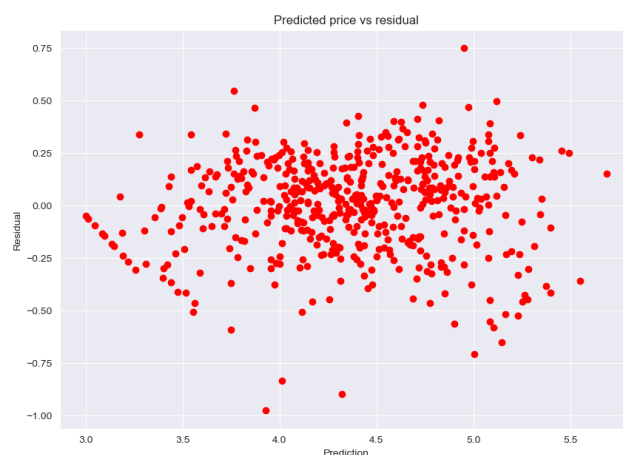
Nakon testiranja Lasso regresionog modela sa top 8 featura na test podacima, dobili smo sledeće rezultate: R^2 skor za testiranje je 0.8419, RMSE (srednja kvadratna greška) je 0.2273, MSE (srednja kvadratna greška) je 0.0516, a MAE (srednja apsolutna greška) je 0.1776. Lasso regresija se pokazala bolje na test skupu u poređenju sa validacionim skupom. Hiperparametri Lasso regresije koji su se najbolje pokazali su: alpha parametar 0.001, i selection criterion 'aic' (Akaike Information Criterio). Naši rezultati pokazuju da Lasso regresija je dala odlične performanse na test podacima, sa visokim R^2 skorom i relativno niskim vrednostima RMSE, MSE i MAE.

Sva tri regresiona modela su se bolje na test skupu u poređenju sa validacionim skupom. Linearna regresija i Lasso regresija imaju identične rezultate, ali ćemo izabrati Lasso regresioni model kao konacni model zbog moći regularizacije i visoke interpretabilnosti.

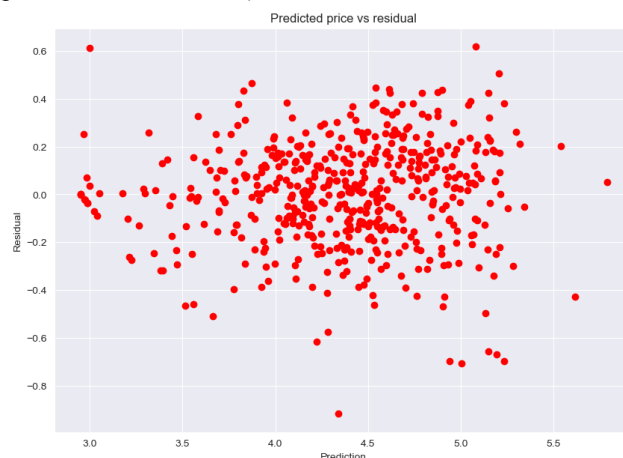
Random forest se sa svim atributima se takodje dobro pokazao i treba izvršiti njegovu evaluaciju nad test skupom. Nakon testiranja Random Forest regresionog modela na test podacima, dobili smo sledeće rezultate: R^2 skor za testiranje je 0.8544, RMSE (srednja kvadratna greška) je 0.2182, MSE (srednja kvadratna greška) je 0.0477, a MAE (srednja apsolutna greška) je 0.1791.

Na slikama 13. i 14. možemo primetiti da Random Forest model ima više scattera skoncentrisanih blizu nuli i da se bolje fitovao nad podacima nego Lasso model.

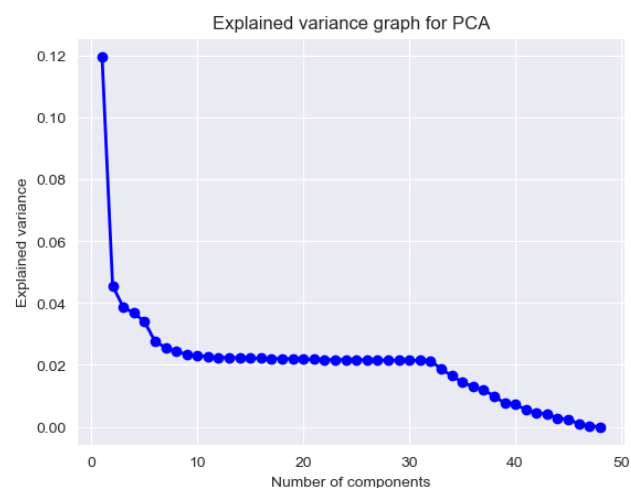
Takodje, vidimo da se Random Forest model nad test podacima ima viši R^2 skor i niže vrednosti RMSE, MSE i MAE od Lasso modela i možemo zaključiti da je najbolji regresioni model Random Forest.



Slika 13. Rezidual (razlika između posmatrane i prediktovane vrednosti) lasso modela



Slika 14. Rezidual (razlika između posmatrane i prediktovane vrednosti) Random forest modela



Slika 15. Scree plot grafik zavisnosti % objasnjene varianse i broj pca komponentni