

# Projektni zadatak: Iznajmljivanje bicikala

Lazar Mitrovic, IN20/2018, [lazar.mitrovic@gmail.com](mailto:lazar.mitrovic@gmail.com)

Aleksa Skrbic, IN29/2018, [aleksa.shrbic@gmail.com](mailto:aleksa.shrbic@gmail.com)

## I Uvod

Izveštaj se bavi analizom meteoroloških podataka i njihovom uticaju na iznajmljivanje bicikala u Njujorku tokom 2011-te i 2012-te godine. Analizom nam dostupnih podataka moguće je kreirati model za predikciju što može doprineti optimizaciji celokupnog procesa iznajmljivanja bicikala i samim tim maksimizacijom profita.

## II Baza podataka

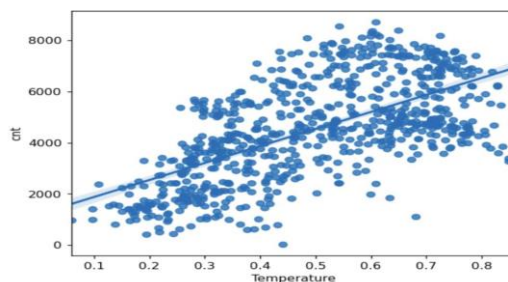
Na raspolaganju imamo imamo 2 dataseta s nazivom **day** i **hour**. **Hour** dataset sadrži podatke o broju iznajmljenih biciklova za svaki sat zajedno sa meteorološkim podacima za dati sat. Sadrži **17378** uzoraka i **17** obeležja od kojih je **5** numeričko i **12** kategoričkih. Kategorička obeležja su redni broj uzorka (**instant**), datum (**dteday**), godišnje doba (**season**), godina (**yr**), mesec (**mnth**), sat (**hr**), informacija o tome da li je bio praznik ili ne (**holiday**), dan u nedelji (**weekday**), informacija o tome da li je bio radni dan ili ne (**workingday**), ocena kvaliteta vremena (**weathersit**). Numerička obeležja su ukupan broj iznajmljenih bicikala za dati sat (**cnt**), broj iznajmljenih bicikala neregistrovanih korisnika (**casual**), broj iznajmljenih bicikala registrovanih korisnika (**registered**), temperatura (**temp**), vlažnost vazduha (**humid**), brzina vetra (**windspeed**) i subjektivni osećaj temperature (**atemp**). Dataset **day** sadrži agregirane podatke dataseta **hour** odnosno sabrani su brojevi iznajmljivanja bicikala za svaki sat i spojeni u jedan dan. Obeležja ova 2 dataseta su skoro identična s tim da **day** ne sadrži obeležje **hr**.

## III Analiza podataka

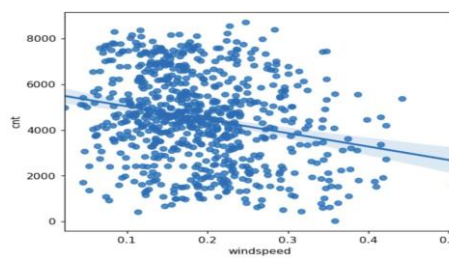
Nijedan od datasetova ne sadrži nedostajuće vrednosti. Obeležja **casual** i **registered** su uklonjena jer ćemo se fokusirati na predikciju i analizu samo ukupnog broja bicikala za sve korisnike. Takođe, obeležje **instant** je uklonjeno jer predstavlja samo redni broj uzorka i ne doprinosi analizi, obeležje datum je uklonjeno jer već imamo podatke o mesecu i godini. Na slikama 1, 2...6 biće analizirana obeležja dataseta **day**.

### A. Uticaj vremenskih prilika na broj iznajmljenih bicikala

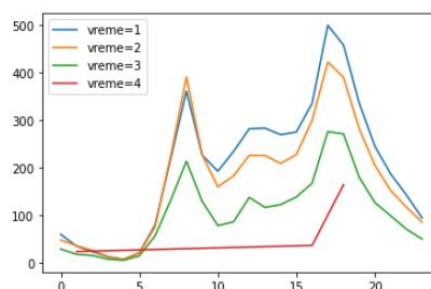
Posmatrajuci sliku 1 možemo uočiti pozitivnu korelaciju između temperature i broja iznajmljenih bicikala odnosno porastom temperature raste i broj iznajmljenih bicikala što je i logično jer ako je vreme lepše, više ljudi će zeleći da vozi. Na slici 2 vidimo suprotan slučaj, porastom brzine vetra opada broj iznajmljenih bicikala što takođe ima smisla jer jači vetar ometa i smanjuje kvalitet vožnje bicikla.



Slika 1. Uticaj temperature na broj iznajmljenih bicikala



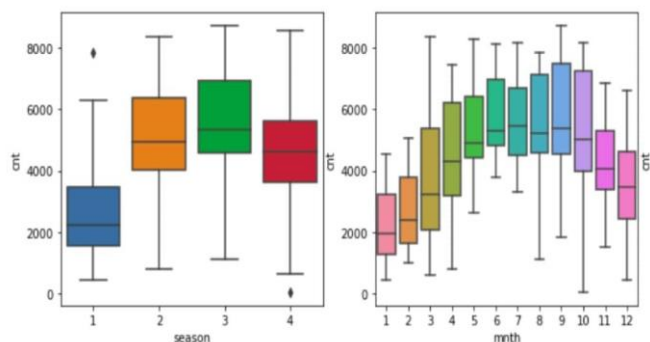
Slika 2. Uticaj brzine vetra na broj iznajmljenih bicikala



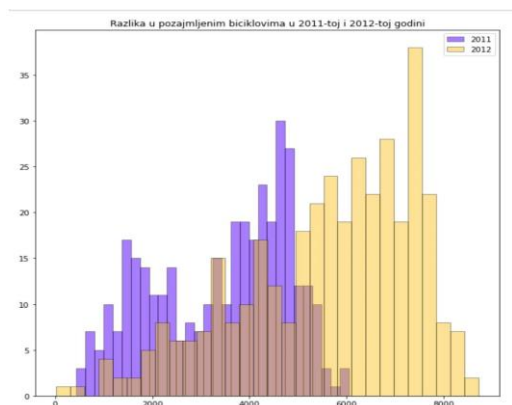
Slika 3. Prosecan broj iznajmljenih bicikala u periodu od 24 casa za razlicite vremenske uslove

Na slici 3 vidimo kako se menja prosek iznajmljenih bicikala s ocenom kvaliteta vremena (1 je najbolja ocena, 4 najgora). Plavom bojom oznacena je kriva prosecnog broja iznajmljenih bicikala kada je vreme ocenjeno najbolje i tada je prosečno najviše biciklova iznajmljeno dok kada je vreme ocenjeno sa 4, prosek je ubedljivo najmanji.

Na slici 4 vidimo da je potražnja bicikala najveća u jesen (mesec septembar), tada su i vremenski uslovi najbolji i imaju ocenu 1, takođe potražnja je visoka i na leto što se pripisuje lepom i suncanom vremenu. Potraživanje drastično opada tokom zimskih meseci.

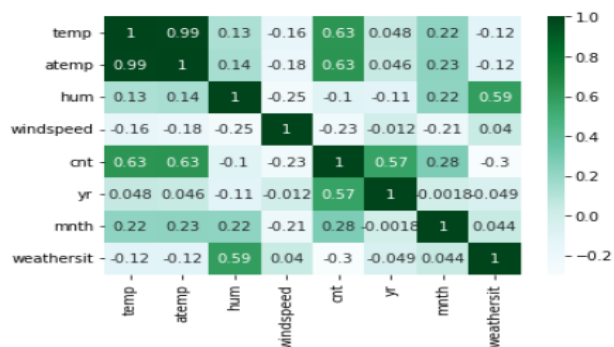


Slika 4 – Levo se nalazi broj iznajmljenih bicikala po godišnjim dobima dok je desno broj iznajmljenih bicikala po mjesecima



Slika 5. Promena u broju iznajmljenih bicikala između 2011-te i 2012-te godine

Sa *slike 5* uočljiva je promena u broju iznajmljenih bicikala u 2012-toj u odnosu na 2011-tu godinu. Jasno se može videti da je dosta veći broj iznajmljenih bicikala u 2012-toj godini. Jedan od validnih zaključaka može biti da su građani svesniji kako upotreba fosilnih goriva uništava ekosistem i zagađuje vazduh te su se opredelili za zdraviju i jeftiniju varijantu od voznje.



Slika 6. Matrica korelacije za pojedina obeležja

Na *slici 6* vidimo matricu korelacije za pojedina obeležja. Radi preglednosti, neka od obeležja čija je korelacija s obeležjem **cnt** a i ostalim obeležjima jako mala nisu prikazana. Kao što smo i mogli zaključiti na osnovu slike 1 obeležja **temp** i **atemp** imaju veću pozitivnu korelaciju s obeležjem **cnt** od svih ostalih obeležja. Takođe, obeležje **yr** je pozitivno korelisano s **cnt**. Obeležja **temp** i **atemp** su međusobno jako visoko korelisana što je i logično jer predstavljaju skoro jednu te istu stvar i retko mnogo odstupaju jedna od druge.

U narednom delu izveštaja bavićemo se predikcijom obeležja **cnt** i to koristeći tehnike linearne regresije, stabla odluke i neuralnih mreža. Navedene tehnike biće primenjene na oba dataseta.

## IV Linearna Regresija

### A. Linearna regresija na day datasetu

Podaci su podeljeni na skup za trening koji sadrži 85% od ukupnog broja uzoraka i test skup koji sadrži 15% od ukupnog broja uzoraka. Za svaki od algoritama za obučavanje modela urađena je i kros validacija na trening skupu s parametrom **n\_splits=5** te izračunata srednja vrednost za svaku od ocena performansi i prikazana na narednim slikama. Test skup od 15% je ostavljen za krajnje testiranje modela. Za obučavanje prvog modela korišćena je klasa **LinearRegression**, obučen je na trening skupu i testiran na test skupu. Takođe, izvršena je i standardizacija obeležja koja je doprinela malim poboljšanjima performansi. Njegove performanse možemo videti na *slici 7*.

#	result
<b>mse</b>	669960.963
<b>mae</b>	618.832
<b>rmse</b>	817.515
<b>r2 score</b>	0.803
<b>r2 adj</b>	0.798

Slika 7. LinearRegression

#	result
<b>mse</b>	753980.332
<b>mae</b>	546.490
<b>rmse</b>	841.447
<b>r2 score</b>	0.786
<b>r2 adj</b>	0.781

Slika 8. Hipoteza samo interakcija

#	result
<b>mse</b>	618821.932
<b>mae</b>	517.097
<b>rmse</b>	769.324
<b>r2 score</b>	0.824
<b>r2 adj</b>	0.820

Slika 9. Hipoteza interakcije i kvadrata

#	result
<b>mse</b>	621365.605
<b>mae</b>	505.393
<b>rmse</b>	768.253
<b>r2 score</b>	0.824
<b>r2 adj</b>	0.820

Slika 10. Lasso regresija

Na *slikama 7,8,9 i 10* vidimo performanse 4 različita modela linearne regresije. Za sada najbolje performanse ima Lasso regresija za koju je koriscen parameter  $\alpha=0.01$ . Ima najmanje odstupanje predviđenih od stvarnih vrednosti i najveći r2 score odnosno pokriva najviše udela ukupne varijanse. Primećuje se da svi modeli osim prvog (*slika 7*) daju slične performanse, s malim razlikama. Međutim, najbolje performanse je dala Ridge regresija čije su performanse predstavljene na *slici 11*.

#	result
<b>mse</b>	470363.497
<b>mae</b>	496.728
<b>rmse</b>	684.757
<b>r2 score</b>	0.862
<b>r2 adj</b>	0.859

Slika 11. Ridge regresija

#	result
<b>mse</b>	503400.507
<b>mae</b>	520.300
<b>rmse</b>	709.507
<b>r2 score</b>	0.882
<b>r2 adj</b>	0.865

Slika 12. Ridge regresija na krajnjem test skupu

Ridge regresija po svim parametrima pobeđuje ostale modele pa je ona testirana i na krajnjem test skupu. Treba uzeti u obzir da ceo dataset **day** ima samo 730 uzoraka te sam test skup nije veliki i performanse modela dobijene nad njime treba uzeti s rezervom. Performanse finalnog modela mozemo videti na *slici 12*. Za Ridge regresiju koriscen je parametar **alfa=13** koji je dobijen metodom **RidgeCV**. Takodje, pokusana je redukcija dimenzionalnosti metodom **PCA**, medjutim doprinela je samo smanjenju performansi pa o njoj nije diskutovano.

## B. Linearna regresija na hour datasetu

**Hour** dataset sadrzi puno vise uzoraka nego **day** pa u ovom slucaju test skup sadrzi 10% od ukupnog broja uzoraka dok ostalih 90% pripada trening skupu. Sve sto je radjeno na **day** datasetu ponovo je radjeno na **hour** datasetu, medjutim modeli dobijeni na hour datasetu daju slabije performanse nego na **day** pa se o njima u ovom izvestaju nece posebno diskutovati.

## V Stabla odluke

### A. Stabla odluke na day datasetu

U ovom poglavlju ce biti diskutovano o modelima obucenim na bazi stabala odluke, konkretno **DecisionTreeRegressor** i **RandomForestRegressor**. Zbog potreba ovih algoritama kreirane su dummy varijable za obelezja **season**, **mnth**, **weathersit** i **weekday**. Podela je kao i do sada izvršena u odnosu 15% test skup i 85% trening skup, takodje radjena je kros validacija s 5 particija.

#	result
<b>mse</b>	885320.779
<b>mae</b>	732.156
<b>rmse</b>	938.049
<b>r2 score</b>	0.741
<b>r2 adj</b>	0.725

Slika 13. DecisionTreeRegressor

#	result
<b>mse</b>	771151.957
<b>mae</b>	687.400
<b>rmse</b>	876.557
<b>r2 score</b>	0.775
<b>r2 adj</b>	0.761

Slika 14. RandomForestRegressor

Na *slikama 13. i 14.* vidimo performanse modela dobijenih metodama **DecisionTreeRegressor** i **RandomForestRegressor**. Nesto bolje performanse daje **RandomForestRegressor** sto je u neku ruku logično jer slučajna suma koristi više stabala odluke kako bi dosla do boljih rezultata. Optimalni parametri za ove 2 metode dobijeni su koriscenjem metode **GridSearchCV** kojoj su prosledjene razlicite kombinacije parametara i uz pomoc malo rucnih korekcija. Optimalni parametri za **DecisionTree** su:  
**max\_depth=3, criterion='friedman\_mse', max\_features='auto', max\_leaf\_nodes=None, min\_samples\_leaf=1, min\_weight\_fraction\_leaf=0.1, splitter='best'**, dok je za **RandomForest** dobijeno:  
**n\_estimators=100, max\_depth=5, max\_features=None, max\_leaf\_nodes=50, min\_samples\_leaf=1, min\_weight\_fraction\_leaf=0.1, criterion='friedman\_mse'**.

Takodje, pokusana je redukcija dimenzionalnosti **PCA** metodom ali nije dovela do poboljsanja performansi.

### B. Stabla odluke na hour datasetu

Podela je izvršena tako da test skup sadrzi 10% a trening 90% uzoraka. Kros validacija je radjena na 3 particije zbog velicine dataseta. Ostala predobrada podataka i korisceni algoritmi su isti kao i na stablima odluke nad **day** datasetom. Rezultati dobijeni stablima odluke nad **hour** datasetom vidno su bolji od onih nad **day** datasetom jer po svojoj prirodi stable odluke i slucajne sume bolje rade na vise uzoraka pa se to i u ovom slucaju ispostavilo kao tacno. Performanse modela mozemo videti na *slikama 15. i 16.*

#	result
<b>mse</b>	3211.717
<b>mae</b>	33.951
<b>rmse</b>	56.667
<b>r2 score</b>	0.901
<b>r2 adj</b>	0.901

Slika 15. DecisionTreeRegressor(hour)

#	result
<b>mse</b>	1883.863
<b>mae</b>	26.447
<b>rmse</b>	43.382
<b>r2 score</b>	0.942
<b>r2 adj</b>	0.942

Slika 16. RandomForestRegressor(hour)

Vidimo da je r2 score kod obe metode dosta veci nego na **day** datasetu odnosno da pokrivaju veci udeo od ukupne varijanse u odnosu na prethodne modele. Takodje se vidi da je odstupanje od stvarnih vrednosti dosta malo. Za odredjivanje optimalnih parametara takodje je koriscena metoda **GridSearchCV** uz malo rucnog podesavanja. Za **DecisionTreeRegressor** dobijeni su parametri:

**max\_depth=12, max\_features='auto', max\_leaf\_nodes=None, min\_samples\_leaf=1, splitter='best', criterion='friedman\_mse'** a za **RandomForestRegressor** **n\_estimators=350, max\_depth=None, criterion='friedman\_mse', max\_features='auto', max\_leaf\_nodes=None, min\_samples\_leaf=1**.

**PCA** redukcija dimenzionalnosti ni ovoga puta nije dovela do poboljsanja performansi, stavise znatno ih je smanjila dajuci **r2 score** 0.38. Posto je **RandomForestRegressor** na **hour** datasetu dao najbolje performanse on jer testiran i na test skupu koji smo ostavili za krajnje testiranje. Performanse ovog modela mozemo videti na *slici 17.*

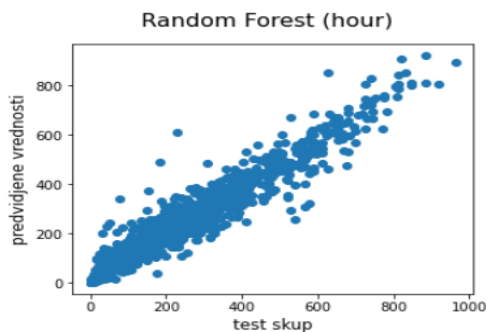
#	result
<b>mse</b>	1916.810
<b>mae</b>	26.339
<b>rmse</b>	43.781
<b>r2 score</b>	0.943
<b>r2 adj</b>	0.943

Slika 17. RandomForestRegressor

Primecujemo da su za **RandomForestRegressor** performanse slične na kros validaciji i na krajnjem test skupu sto znaci da nas model dobro generalizuje i da je dobro obucen.

Na *slici 18.* mozemo videti odnos pravih i predvidjenih vrednosti na krajnjem test skupu. Vidimo da su vrednosti jako slične sto se i moze zakljuciti iz performansi finalnog modela

sa *slike 17*. I vidimo da bi se oni mogli lepo aproksimirati pravom.



Slika 18. Odnos pravih I predviđenih vrednosti u modelu dobijenom metodom slučajne sume

## VI Neuralna mreža

U ovom poglavlju vice diskutovano u neuralnoj mrezi implementiranoj pomocu metode **MLPRegressor**. Test skup sadrzi 10% a trening 90% od ukupnog broja uzoraka. Takodje, I u ovom slucaju je radjena kros validacija sa 3 particije. Performanse neuralne mreze na datasetu **day** su se pokazale jako lose, sto je posledica manjka uzorka tog dataseta da bi se neuralna mreza adekvatno obucila. Stoga, u ovom poglavlju ce biti predstavljeni rezultati modela samo na datasetu **hour**.

Optimalni parametra za neuralnu mrezu dobijeni su opet metodom **GridSeachCV** I uz odredjenu dozu rucnih korekcija. Primeceno je da porastom broja skrivenih slojeva performanse modela se poboljsavaju dok ne dodje do 4 skrivena sloja gde performanse dostizu svoj peak. Takodje, porastom broja neurona u slojevima performanse se poboljsavaju.

Parametri **MLPRegressor-a** koji su korisnici su: **hidden\_layer\_sizes=(300,300,300,300)**, **activation='relu'**, **solver='adam'**, **batch\_size=50**, **learning\_rate='adaptive'**, **learning\_rate\_init=0.001**, **max\_iter=100**, **shuffle=True**, **random\_state=42**, **early\_stopping=True**, **n\_iter\_no\_change=10**, **validation\_fraction=0.1**, **verbose=False**.

#	result
<b>mse</b>	1607.403
<b>mae</b>	24.888
<b>rmse</b>	40.061
<b>r2 score</b>	0.951
<b>r2 adj</b>	0.951

Slika 19. MLPRegressor

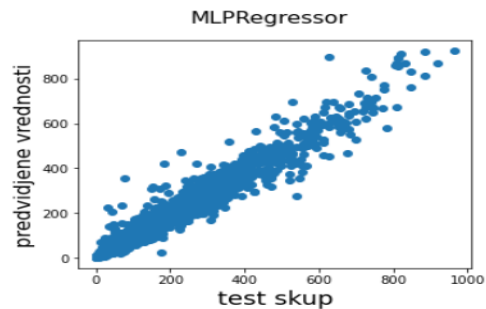
#	result
<b>mse</b>	1747.919
<b>mae</b>	26.413
<b>rmse</b>	41.808
<b>r2 score</b>	0.948
<b>r2 adj</b>	0.948

Slika 20. MLPRegressor(krajnji test skup)

Pokusana je I redukcija dimenzionalnosti metodom **PCA** ali ni ovoga puta nije dala poboljsanja performansi, iako je ovaj put smanjenje performansi neznatno, dajuci **r2 score** 0.913. Rezultati dobijeni na kros validaciji I na konacnom test skupu se ne razlikuju puno sto znaci da model dobro generalizuje I da je dobro obucen.

Primecujemo da neuralna mreza daje najbolje performanse, bolje od svih ostalih algoritama. To moze biti posledica

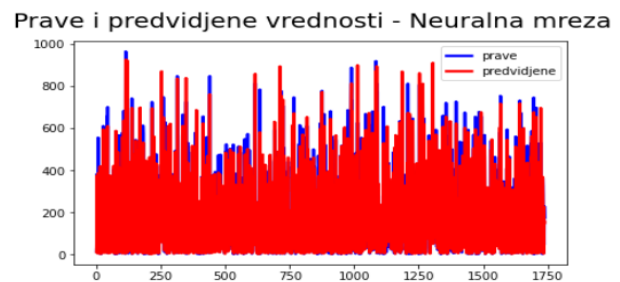
velikog broja adekvatnih parametara koriscenih za obucavanje iste uz dovoljan broj uzoraka da bi obuka bila uspesna. Prosledivsi joj dovoljan broj skrivenih slojeva I neurona mreza je sama uspela da nauči kompleksne veze izmedju obelezja cega su produkt najbolje perofmanse od svih modela.



Slika 21. Odnos pravih I predviđenih vrednosti u modelu dobijenom neuralnom mrežom

Na *slici 21*. Primecujemo istu stvar kao na *slici 18*. a to je da se prave I predviđene vrednosti ne razlikuju puno I da se mogu aproksimirati pravom s tim da je rasipanje u ovom slucaju malo manje jer su performanse modela ipak malo bolje.

Na *slici 22*. uocavamo veliko preklapanje izmedju predviđenih I pravih vrednosti sto smo mogli I da zakljucimo u dosadasnjoj analizi.



Slika 22. Preklapanje pravih I predviđenih vrednosti

## VII Zakljucak

Na osnovu analize podataka iz treceg poglavlja mozemo zakljuciti da bi kompanija koja se bavi iznajmljivanjem bicikala trebalo da se fokusira na prosirivanje biznisa tokom letnjih I jesenjih meseci jer je potraznja tada najveća. Tokom zimskih meseci I losev vremena bi trebalo vise da se fokusiraju na servisiranje I odrzavanje bicikala kako bi bili spremni onda kada ljudima najvise trebaju. Takodje, prognozan je rast potraznje za iznajmljivanjem bicikala tokom godina sto moze biti u korelaciji s potencijalnim porastom ekoloske svesti kod ljudi.

Vidimo da metode **slučajne sume** I **neuralne mreze** daju vidno bolje performanse od linearne regresije ali isto tako zahtevaju veliki broj parametara cega je posledica puno vremena za njihovu obuku. Takodje, za neuralnu mrezu, konkretno **MLPRegressor** je potreban veci broj uzoraka da bi se obucila pa se nije mogao adekvatno iskoristiti na **day** datasetu. Stoga, mozemo zakljuciti da linearna regresija ce u ovom slucaju linearna regresija dati relativno solidan rezultat na datasetu manjeg broja uzoraka dok je za neuralnu mrezu I slučajne sume potreban veci broj uzoraka I vise vremena za obuku al ce dati znatno bolje performanse.