

Osnove statističkog programiranja

Ak. god. 2023./2024.

Neki naslov

Dokumentacija

Julijana Kolarec, Lucija Topolko

siječanj 2024., Zagreb

Nastavnik: *prof.dr.sc. Damir Pintar*

Sadržaj

1	Prediktivni modeli primjenom strojnog učenja	2
1.1	Priprema podataka	2
	Indeks slika i dijagrama	9

1. Prediktivni modeli primjenom strojnog učenja

Kako bismo bolje razumjeli što film čini uspješnim ili neuspješnim, provele smo analizu dobivenog skupa podataka primjenom strojnog učenja. Cilj nam je bio razviti model koji može čim točnije predviđati uspjeh filma na temelju njegovih karakteristika.

1.1 Priprema podataka

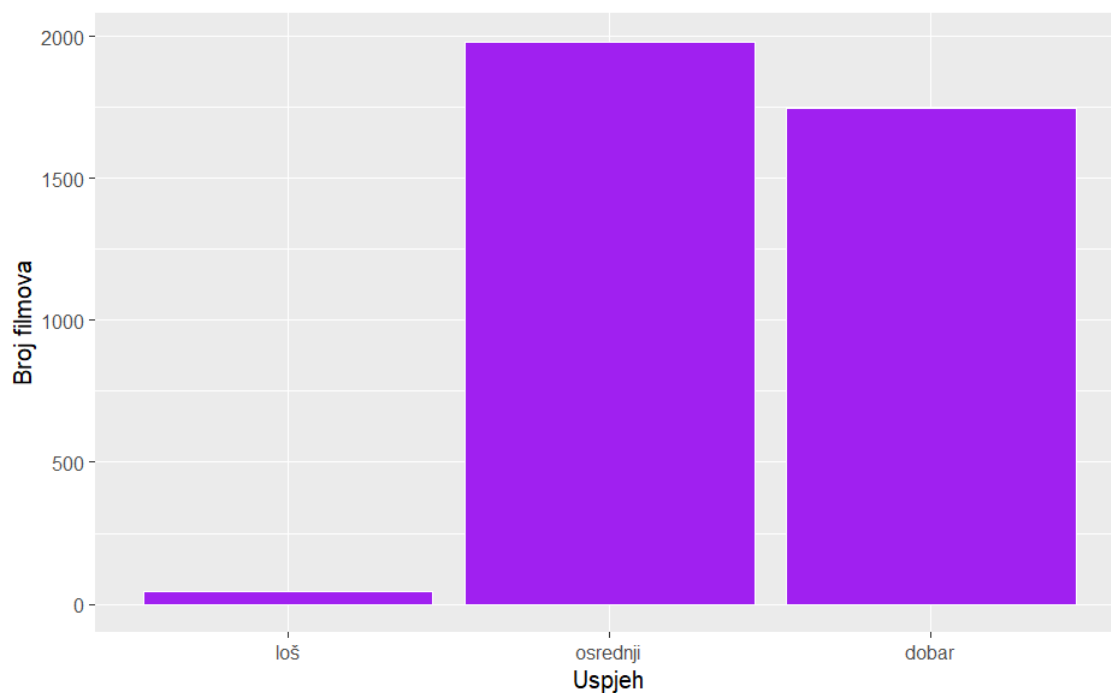
Iz dobivenih podataka izbacile smo retke kojima su nedostajali neki podaci. Takvih je redaka bilo 1261. Također, uklonile smo stupce koji su sadržavali jedinstvene ili skoro jedinstvene vrijednosti (*movie_title*, *movie_imdb_link*, *plot_keywords*, *genres*). Još smo izbacile tekstualne stupce koji su bili prekorelirani s nekim numeričkim stupcem. Na primjer, *actor_1_name* je prekoreliran s *actor_1_facebook_likes*.

```
1 columns <- c('duration', 'director_facebook_likes', 'actor_1_facebook_likes', 'actor_2_facebook_likes', 'actor_3_facebook_likes', 'num_user_for_reviews', 'num_critic_for_reviews', 'num_voted_users', 'cast_total_facebook_likes', 'movie_facebook_likes', 'facenumber_in_poster', 'color', 'title_year', 'language', 'country', 'content_rating', 'aspect_ratio', 'gross', 'budget', 'imdb_score')
2
3 mldata <- data[,columns]
4
5 mldata <- na.omit(mldata)
```

Preostale nenumeričke stupce pretvorile smo u tip integer.

```
1 label_encode <- function(column) {
2   as.integer(factor(column, levels = unique(column)))
3 }
```

Značajka koju predviđamo je *imdb_score*. To broj zaokružen na jednu decimalu, pa smo za bolje rezultate uspjeh filma podijelile u tri skupine: loš, osrednji i dobar,



Slika 1.1: Podjela filmova po uspjehu

a stupac `imdb_score` smo zbog prekorreliranosti uklonile.

```
1 mldata$score <- ifelse(mldata$imdb_score < 3.33, "los", ifelse(mldata$  
  imdb_score < 6.66, "osrednji", "dobar"))
```

Graf 1.1 prikazuje omjer broja filmova po uspjehu. Filmova koji su ocijenjeni kao loši znatno je manje od ostalih. Točnije, loših je filmova 43, osrednjih 1981, a dobrih 1746.

```
1 Model <- train(score ~ ., data = training_set,  
2               method = "svmPoly",  
3               na.action = na.omit,  
4               preProcess=c("scale","center"),  
5               trControl= trainControl(method="none"),  
6               tuneGrid = data.frame(degree=1, scale=1, C=1)  
7 )
```

Confusion Matrix and Statistics

		Reference		
Prediction		loš	osrednji	dobar
loš		0	0	0
osrednji		8	339	141
dobar		0	54	205

Overall Statistics

Accuracy : 0.7282
95% CI : (0.6948, 0.7599)
No Information Rate : 0.5261
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.4518

Mcnemar's Test P-Value : NA

Statistics by Class:

	Class: loš	Class: osrednji	Class: dobar
Sensitivity	0.00000	0.8626	0.5925
Specificity	1.00000	0.5791	0.8653
Pos Pred Value	NaN	0.6947	0.7915
Neg Pred Value	0.98929	0.7915	0.7111
Prevalence	0.01071	0.5261	0.4632
Detection Rate	0.00000	0.4538	0.2744
Detection Prevalence	0.00000	0.6533	0.3467
Balanced Accuracy	0.50000	0.7208	0.7289

```
1 Model_rf <- randomForest(score ~ ., data = training_set, ntree = 500,  
  importance = TRUE)
```

Confusion Matrix and Statistics

		Reference		
Prediction		loš	osrednji	dobar
loš		0	0	0
osrednji		7	325	86
dobar		1	68	260

Overall Statistics

Accuracy : 0.7831
95% CI : (0.7518, 0.8122)
No Information Rate : 0.5261
P-Value [Acc > NIR] : <2e-16

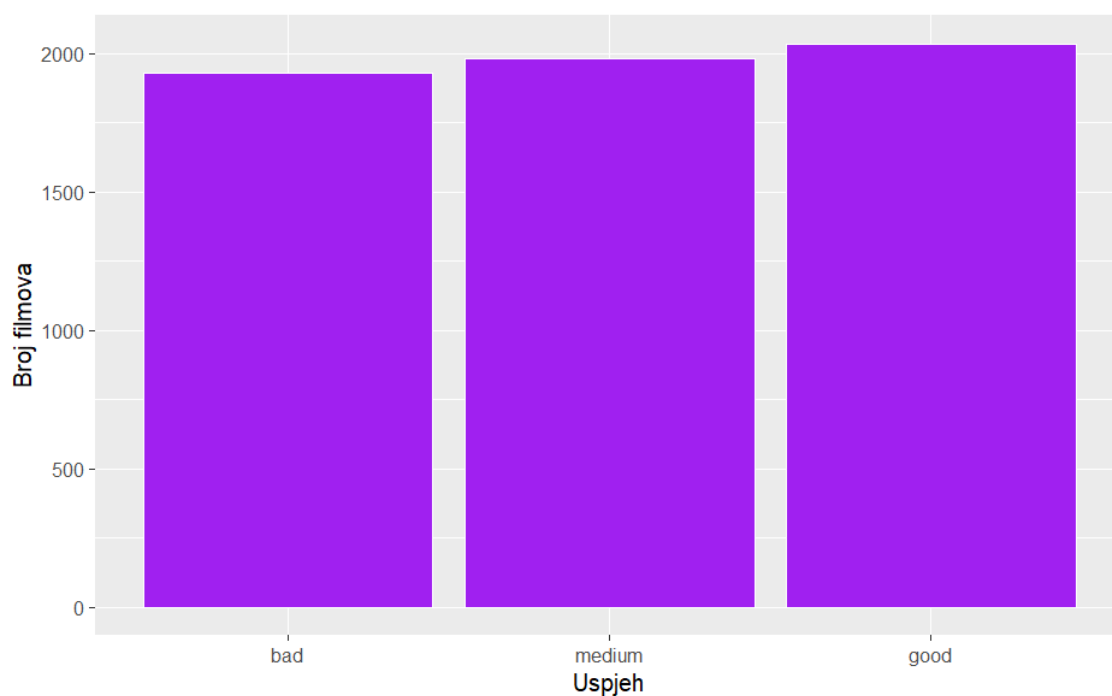
Kappa : 0.5677

McNemar's Test P-Value : 0.0177

Statistics by Class:

	Class: loš	Class: osrednji	Class: dobar
Sensitivity	0.00000	0.8270	0.7514
Specificity	1.00000	0.7373	0.8279
Pos Pred Value	NaN	0.7775	0.7903
Neg Pred Value	0.98929	0.7933	0.7943
Prevalence	0.01071	0.5261	0.4632
Detection Rate	0.00000	0.4351	0.3481
Detection Prevalence	0.00000	0.5596	0.4404
Balanced Accuracy	0.50000	0.7821	0.7897

```
oversample <- ovun.sample(score~., data = over, method = "both", N =  
3932)$data
```



Confusion Matrix and Statistics

Reference			
Prediction	loš	osrednji	dobar
loš	374	69	34
osrednji	7	256	109
dobar	0	68	261

Overall Statistics

Accuracy : 0.7564
95% CI : (0.7308, 0.7806)
No Information Rate : 0.343
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.6352

McNemar's Test P-Value : < 2.2e-16

Statistics by Class:

	Class: loš	Class: osrednji	Class: dobar
Sensitivity	0.9816	0.6514	0.6460
Specificity	0.8708	0.8522	0.9121
Pos Pred Value	0.7841	0.6882	0.7933
Neg Pred Value	0.9900	0.8300	0.8316
Prevalence	0.3234	0.3336	0.3430
Detection Rate	0.3175	0.2173	0.2216
Detection Prevalence	0.4049	0.3158	0.2793
Balanced Accuracy	0.9262	0.7518	0.7791

Confusion Matrix and Statistics

Prediction	Reference		
	loš	osrednji	dobar
loš	381	0	0
osrednji	0	330	44
dobar	0	63	360

Overall Statistics

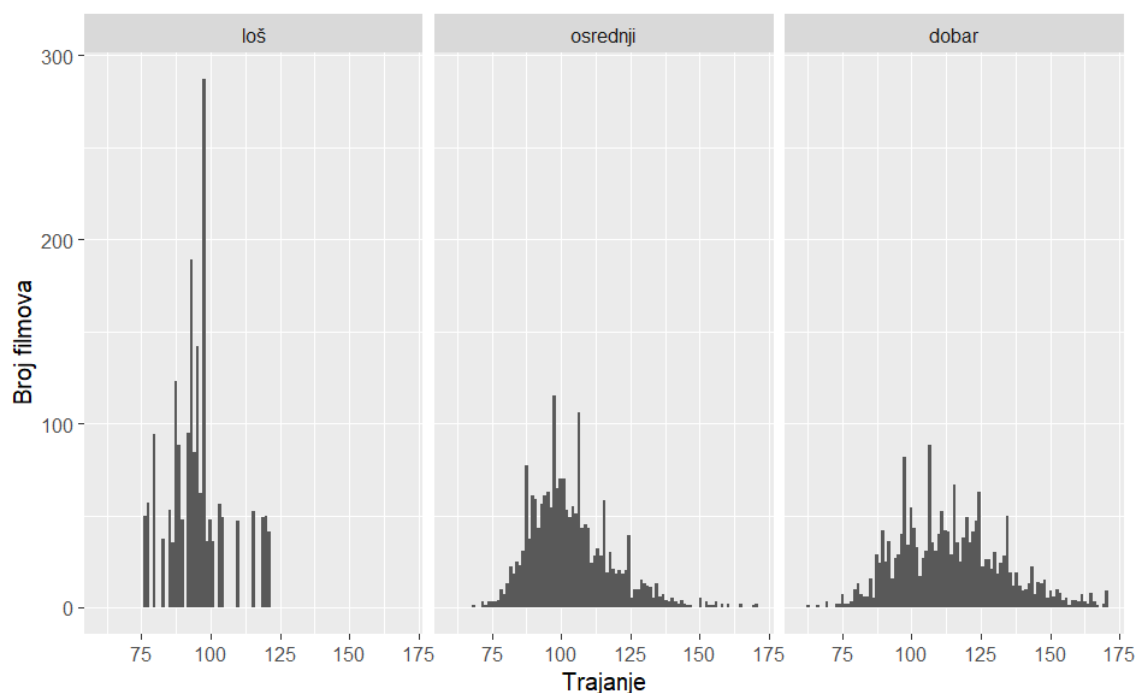
Accuracy : 0.9092
95% CI : (0.8913, 0.925)
No Information Rate : 0.343
P-Value [Acc > NIR] : < 2.2e-16

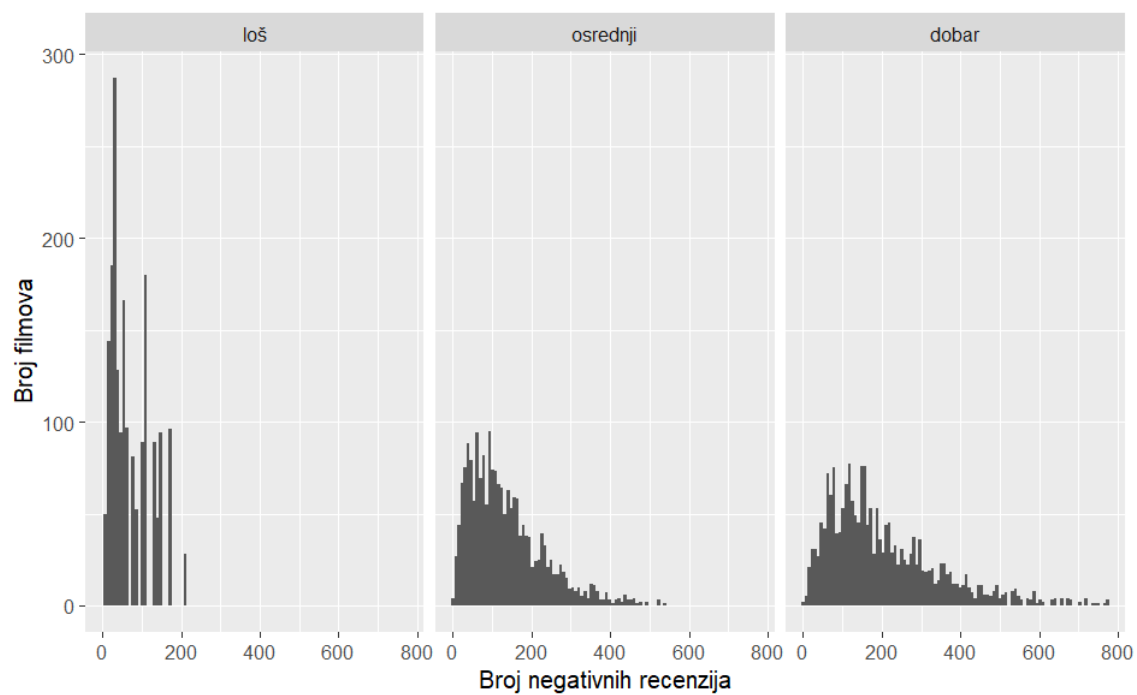
Kappa : 0.8637

Mcnemar's Test P-Value : NA

Statistics by Class:

	Class: loš	Class: osrednji	Class: dobar
Sensitivity	1.0000	0.8397	0.8911
Specificity	1.0000	0.9439	0.9186
Pos Pred Value	1.0000	0.8824	0.8511
Neg Pred Value	1.0000	0.9216	0.9417
Prevalence	0.3234	0.3336	0.3430
Detection Rate	0.3234	0.2801	0.3056
Detection Prevalence	0.3234	0.3175	0.3591
Balanced Accuracy	1.0000	0.8918	0.9048





Indeks slika i dijagrama

1.1	Podjela filmova po uspjehu	3
-----	--------------------------------------	---