



УНИВЕРЗИТЕТ У БЕОГРАДУ  
ФАКУЛТЕТ ОРГАНИЗАЦИОНИХ НАУКА

# Предвиђање цена половних аутомобила применом алгоритама машинског учења

Ментор: Др Сандро  
Радовановић

Алекса Радојичић,  
2019/0165

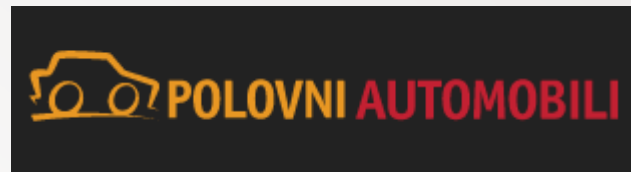
# Садржај

Увод	1	4	Резултати
Теоријске основе	2	5	Расправа
Опис истраживања	3	6	Закључак

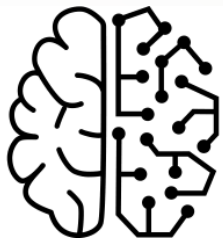
# 1

## Увод

- 77.973 путничких возила на платформи *polovniautomobili*
- Прецењеност и потцењеност кола и 90%-ни интервали предвиђаја



## 2 Теоријске основе



### Машинско учење [МУ]

Направити математички модел који предвиђа излазне вредности на основу улазних вредности. Користити познате податке за прескаживање сродних и непознатих података.

Целокупни подаци

Скуп за обуку

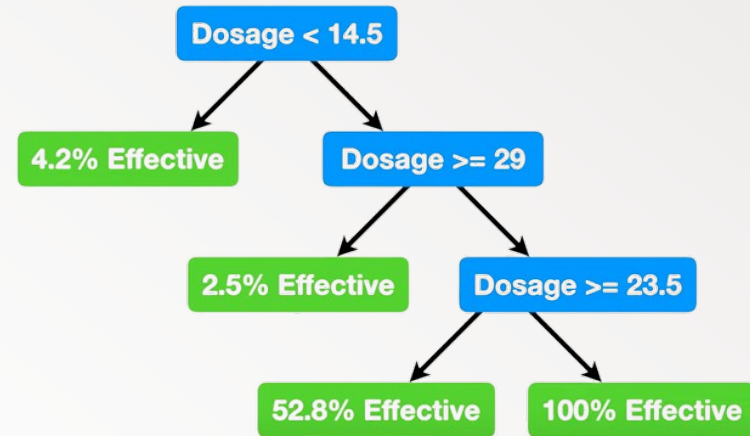
Скуп за тестирање

**Потприлагодба спрам преприлагоде**

## 2.1

# Дрво одлучивања [DT]

- ❑ Састоји се од чворова и грана. Чвор који се не грана је лист.
- ❑ Врши се гранање тако да се највише повећа информациона добит.
- ❑ Нелинеаран алгоритам.
- ❑ Тумачљив и временски и меморијски исплатив.

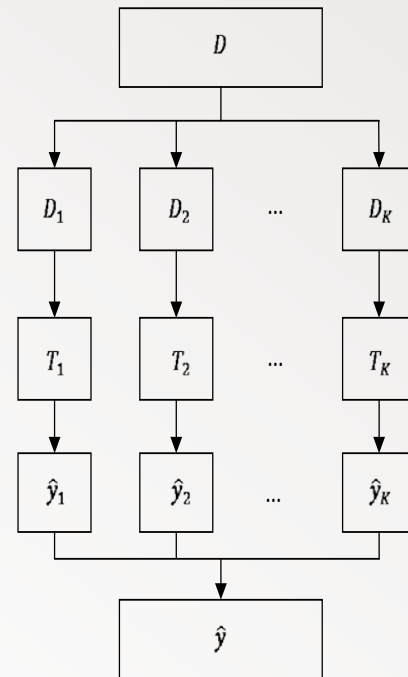


Преузета слика са [StatQuest](#) снимка

## 2.2

## Случајна шума [RF]

- ❑ Начинити више подскупа са различитим редовима и / или ступцима (**багинг**) и над сваким обучити дрва одлучивања.
- ❑ Боље перформансе и честа употреба у пракси.
- ❑ Мање тумачљив алгоритам и меморијски скуп.
- ❑ Подржава паралелизацију.



## 2.3

# Модел градијентног појачавања [GBM]

- ❑ Редно (итеративно) се творе дрва одлучивања тако да свако идуће буде боље од претходног (**појачавање**).
- ❑ Максимални саоднос (корелација) са негативним градијентом функције губитка.
- ❑ Омогућавају творење интервала предвиђаја.
- ❑ Не подржава паралелизацију.

# Опис истраживања

- 1) Прикуп података
- 2) Увид у податке
- 3) Предобрада података
- 4) Обука модела
- 5) Оптимизација модела





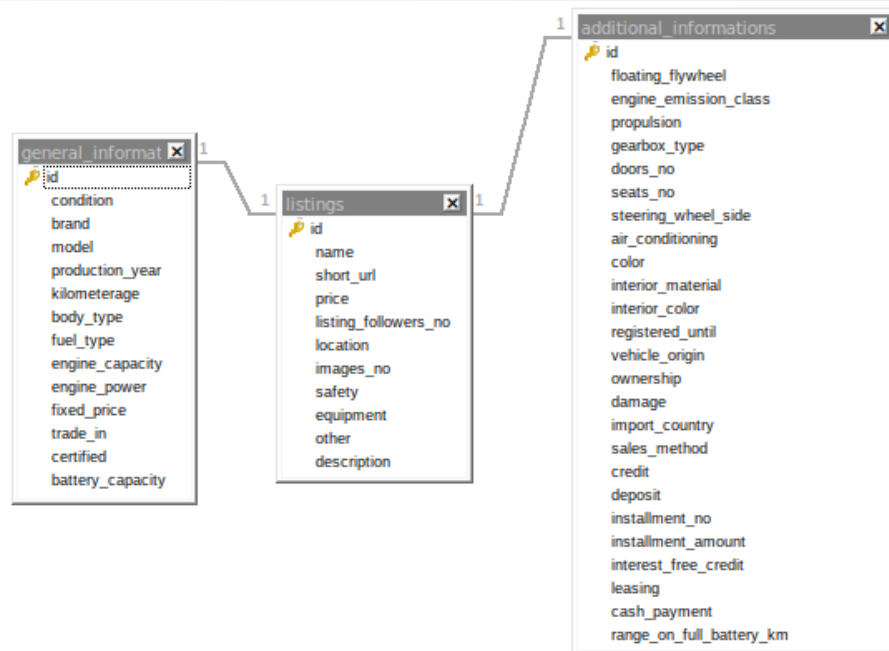
## 3.1 Прикуп података



BeautifulSoup



- Прикупљено 30.788 редова и 50 стубаца (колона).



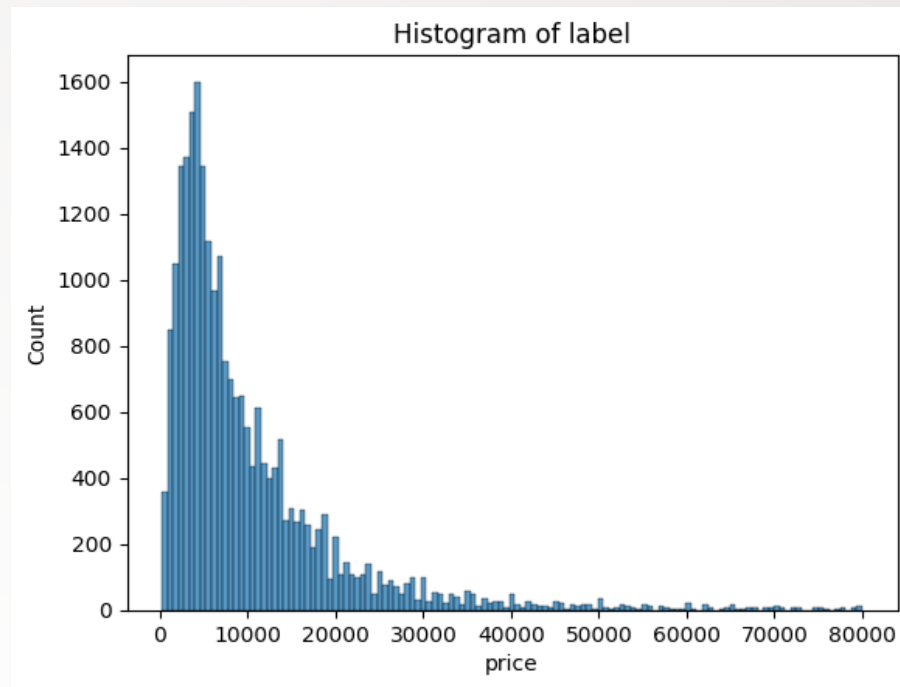


## 3.2 Увид у податке

	фрекв.	удео [%]
<b>„марка“</b>		
Volkswagen	3734	15,716150
Audi	2529	10,644387
BMW	2151	9,053411
Opel	1862	7,837030
Peugeot	1582	6,658529

	фрекв.	удео [%]
<b>„локација“</b>		
Београд	3828	16,111789
Нови Сад	1709	7,193064
Ниш	954	4,015321
Крагујевац	883	3,716486
Чачак	730	3,072520

	фрекв.	mean	std	min	25%	50%	75%	max
<b>„година производње“</b>	23759	2010,59472	6,25557	1953	2007	2011	2015	2024

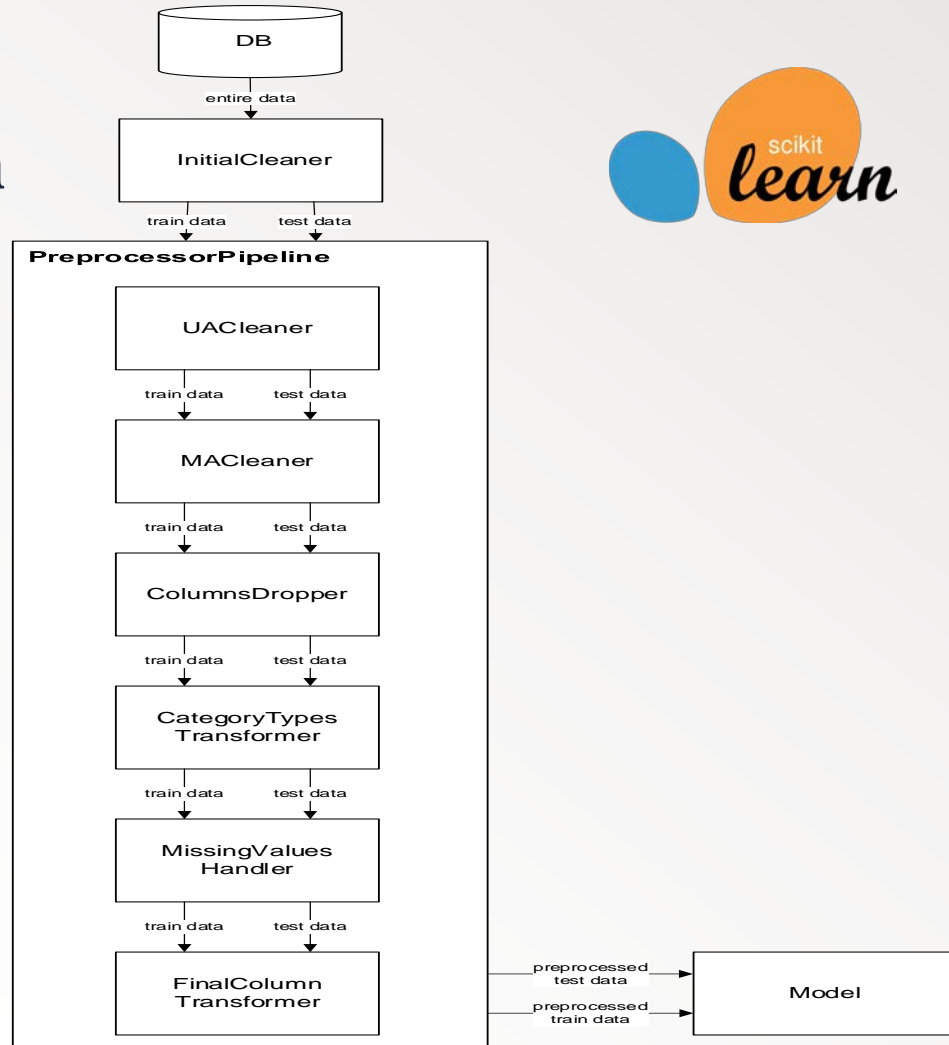


## 3.3

# Предобрада података



- ❑ Димензије података после уводног чишћења: 29.699, 164.
- ❑ Цевовод за предобраду:  
уклања слабо употребљиве колоне,  
обрађује типове података, начин  
попуне недостајућих вредности...
- ❑ Обрађује засебно скупове за обуку и  
тестирање.



## 3.4

## Обука модела

## – почетни модели –

- ❑ 80% кола (23.759) за обуку, 20% (5.940) за тестирање.
- ❑ Коришћено свих 153 стубаца.
- ❑ Случајна шума се најбоље показала над скупом за тестирање.

	Метрике над скупом за тестирање		
Модел	RMSE	MAE	R2
DT	4198,384	2050,968	0,84021
RF	<b>2786,725</b>	<b>1364,816</b>	<b>0,92960</b>
GBM $Q_{0,50}$	2921,088	1656,298	0,92265
GBM $Q_{0,05}$	6795,749	3460,504	0,58134
GBM $Q_{0,95}$	5720,118	3635,153	0,70339

	Метрике над скупом за обуку		
Модел	RMSE	MAE	R2
RF	<b>1045,089</b>	<b>510,826</b>	<b>0,99013</b>
GBM $Q_{0,50}$	2727,873	1583,053	0,93277
GBM $Q_{0,05}$	5593,659	3502,329	0,57025
GBM $Q_{0,95}$	6896,814	3615,054	0,71731

## 3.5

## Оптимизација модела

- ❑ Употребљено 28 стубаца.
- ❑ Цена логаритмована.
- ❑ Решеткаста претрага за добијање оптималних хиперпараметара (80% обука, 20% вредновање).

Оптимални хиперпараметри за RF

Хиперпараметри	Вредности
<i>n_estimators</i>	<b>245</b>
<i>min_samples_split</i>	<b>2</b>

Оптимални хиперпараметри за GBM  $Q_{0,50}$

Хиперпараметри	Вредности
<i>learning_rate</i>	<b>0,05</b>
<i>n_estimators</i>	<b>500</b>
<i>min_samples_split</i>	<b>15</b>
<i>max_depth</i>	<b>12</b>
<i>criterion</i>	<b>„friedman_mse“</b>

## 3.5

## Оптимизација модела – интервали предвиђаја –

Оптимални хиперпараметри за GBM  $Q_{0,05}$   
(долњи интервал предвиђаја)

Хиперпараметри	Вредности
<i>learning_rate</i>	<b>0,1</b>
<i>n_estimators</i>	<b>2000</b>
<i>max_depth</i>	<b>28</b>
<i>criterion</i>	<b>„friedman_mse“</b>

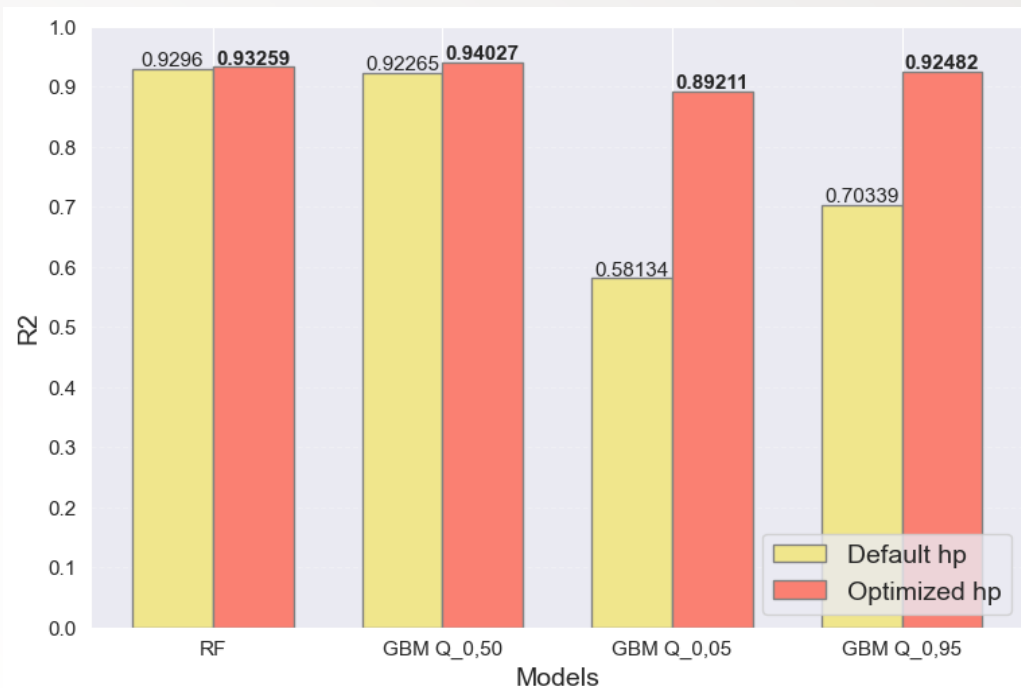
Оптимални хиперпараметри за GBM  $Q_{0,95}$   
(горњи интервал предвиђаја)

Хиперпараметри	Вредности
<i>n_estimators</i>	<b>1000</b>
<i>min_samples_split</i>	<b>5</b>
<i>max_depth</i>	<b>28</b>
<i>criterion</i>	<b>„friedman_mse“</b>

## 4

## Резултати

	Метрике над скупом за тестирање		
Модел	RMSE	MAE	R2
RF	2726,88 4	1311,138	0,93259
GBM $Q_{0,50}$	<b>2566,82</b> <b>2</b>	<b>1230,770</b>	<b>0,94027</b>
GBM $Q_{0,05}$	3449,85 7	1715,0888	0,89211
GBM $Q_{0,95}$	2879,73 6	1624,346	0,92482

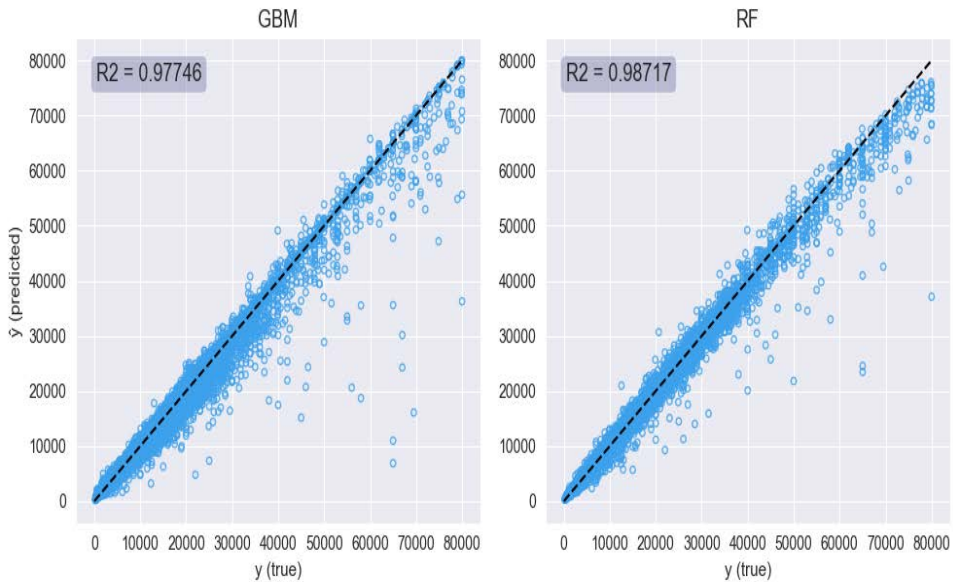




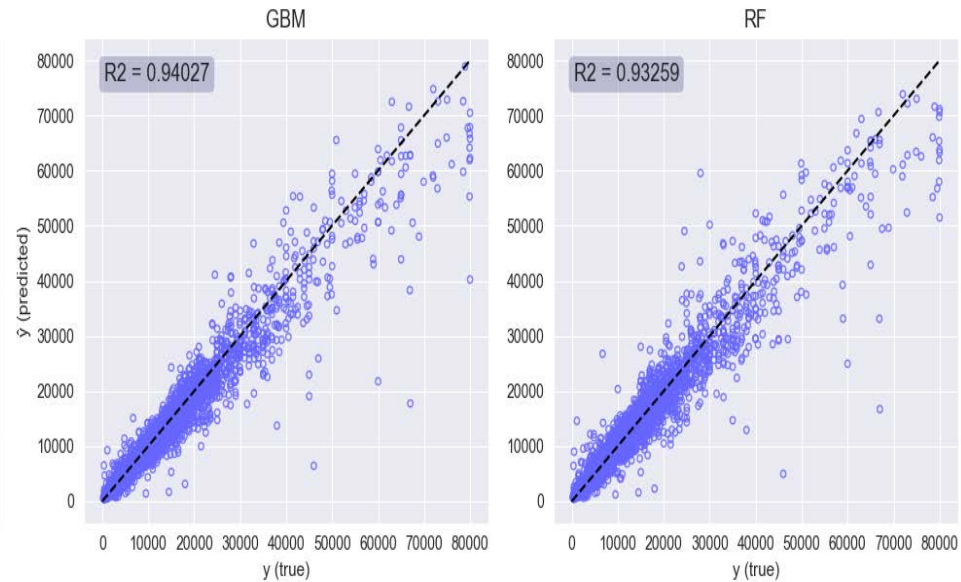
## 4

## Результати

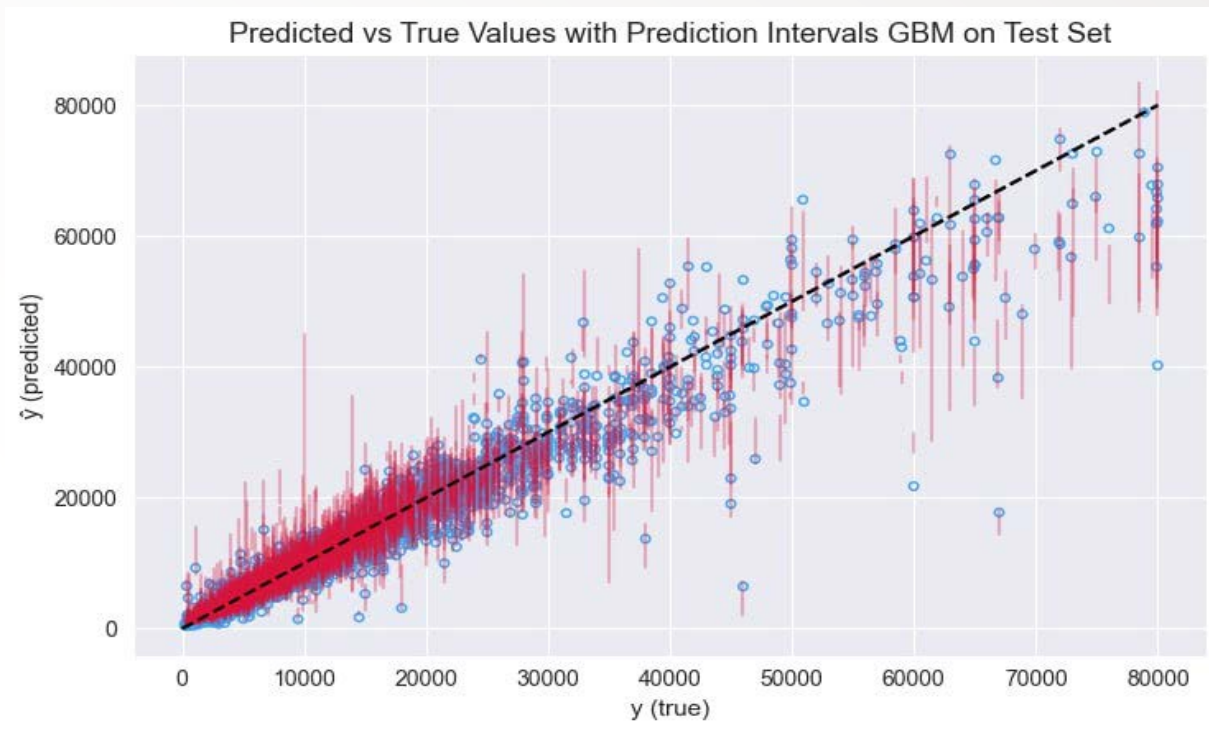
Predicted vs True Values for Train Set

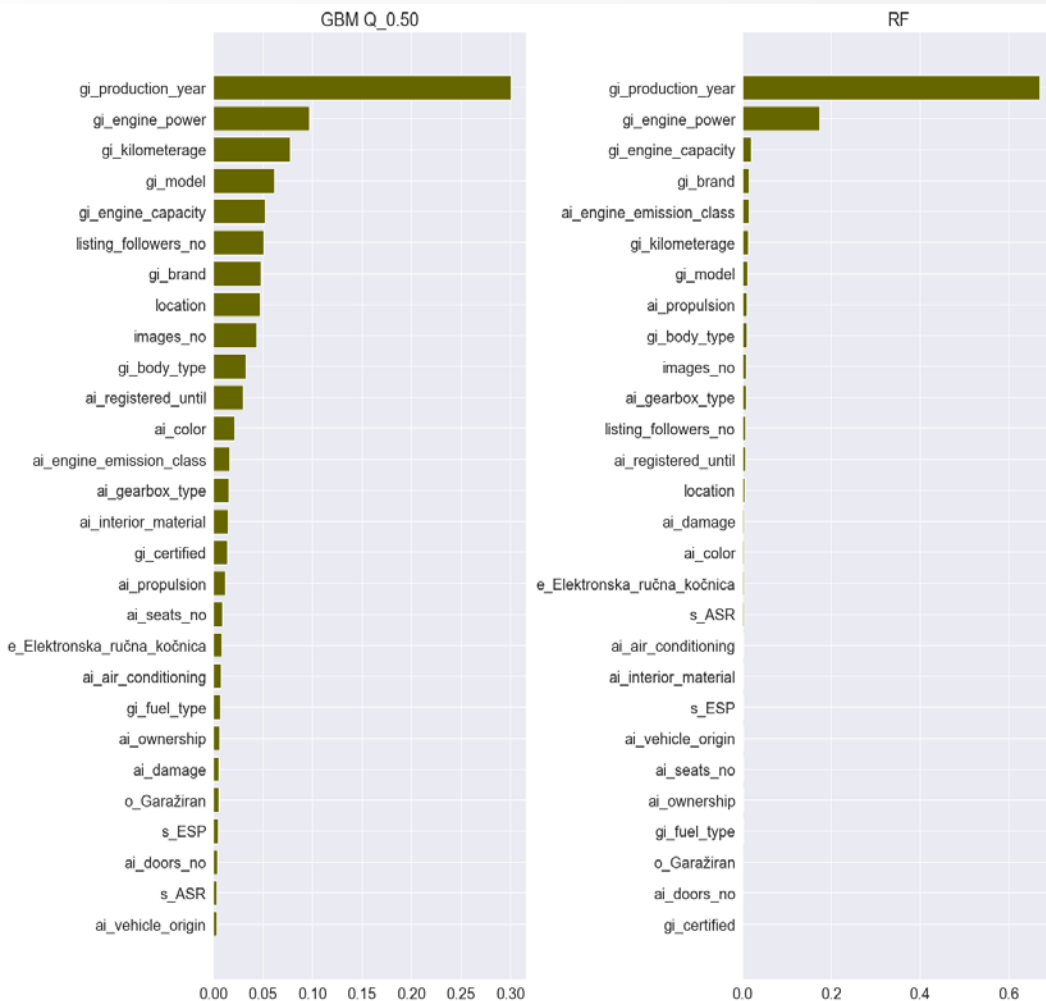


Predicted vs True Values for Test Set



## Резултати – приказ интервала предвиђаја





4

## Резултати Значајности стубаца

- ❑ Обележје „година производње“ јесте најзначајније, а затим „снага мотора“ (KS).
- ❑ GBM знатно боље користи информације из свих стубаца за предвиђај него RF.

## 5 Расправа

- ❑ Мали узорак за кола скупља од 60.000€.
- ❑ Да ли су некоји продавци намерно дигли цену својих кола?
- ❑ Недаћа: непостојање података о тачном износу цене по којој је возило продато.



- ❑ Интервали предвиђаја слабо употребљиви (свега 43,38384% покривеност од очекиваних ~90%).
- ❑ Користити текст и слике приликом предвиђања (вештачке неуронске мреже), искористити више стубаца приликом обуке, сажети обележја у мањи број...
- ❑ Велики потенцијал метода за машинско учење уз доста података.