


```

In [67]: #PurchaseDate, VehicleAge, Make, Model, SubModel, Color, Transmission, WheelType, VehOdo, Size, MMRAcquisitionAveragePrice

RefId      PurchaseDate    VehicleAge    Make     Model   SubModel    Color  Transmission  WheelType  VehOdo    Size    MMRAcquisitionAveragePrice
-----
1    1/7/2009          0.12         DODGE    STRATUS V6 4D SEDAN SXT FV MAROON AUTO Covers 73807 MEDIUM
2    1/7/2009          0.12         DODGE    NEON       4D SEDAN SILVER AUTO Alloy 65617 COMPACT
3    1/7/2009          0.12         FORD     FOCUS      2D COUPE ZK3 SILVER MANUAL Covers 69367 COMPACT
4    1/7/2009          0.12         5 MITSUBISHI GALANT 4C 4D SEDAN WHITE AUTO Covers 81054 MEDIUM
5    1/21/2009        0.12         DODGE    RAM PICKUP 2WD QUAD CAB 4.7L SLT WHITE AUTO Alloy 76173 LARGE TRUCK

In [68]: # Analiza ostalih kategorickih atributa

In [69]: data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6797 entries, 1 to 6798
Data columns (total 22 columns):
# Column Non-Null Count Dtype
---
0 PurchDate 6797 non-null object
1 VehicleAge 6797 non-null int64
2 Make 6797 non-null object
3 Model 6797 non-null object
4 SubModel 6797 non-null object
5 Color 6797 non-null object
6 Transmission 6797 non-null object
7 WheelType 6797 non-null object
8 VehOdo 6797 non-null int64
9 Size 6797 non-null object
10 MMRAcquisitionAuctionAveragePrice 6797 non-null float64
11 MMRAcquisitionAuctionCleanPrice 6797 non-null float64
12 MMRAcquisitionRetailAveragePrice 6797 non-null float64
13 MMRAcquisitionRetailCleanPrice 6797 non-null float64
14 MMRCurrentAuctionAveragePrice 6797 non-null float64
15 MMRCurrentAuctionCleanPrice 6797 non-null float64
16 MMRCurrentRetailAveragePrice 6797 non-null float64
17 MMRCurrentRetailCleanPrice 6797 non-null float64
18 VehCost 6797 non-null int64
19 IsOnlineSale 6797 non-null int64
20 WarrantyCost 6797 non-null int64
21 IsBadBuy 6797 non-null int64
memory usage: 1.2+ MB

In [70]: # PurchaseDate

In [71]: data['PurchaseDate'].value_counts()

1/30/2010 0:09 42
1/25/2009 0:02 41
1/20/2010 0:10 40
1/4/2009 0:03 36
1/8/2010 0:12 36

1/16/2009 0:01 1
1/26/2010 0:03 1
1/19/2010 0:04 1
1/27/2010 0:09 1
1/25/2009 0:05 1
Name: PurchaseDate, Length: 503, dtype: int64

In [72]: # Za nasu analizu znacajna nam je godina aukcije, tako da cemo napraviti novu varijablu gde cemo iz ove izvuci

In [73]: data['PurchaseYear'] = pd.DatetimeIndex(data['PurchaseDate']).year

In [74]: data = data.drop(['PurchaseDate',axis=1])

In [75]: data['PurchaseYear'].value_counts()

2010 3504
2009 3293
Name: PurchaseYear, dtype: int64

In [76]: # Vidimo da imamo samo dve godine i da je broj opservacija priblizno isti za obe, pa cemo sada proveriti kakva
# o ovog atributa sa nekim drugim morda dati bolje rezultate, tako da cemo je za sada ostaviti u datasetu

In [77]: sns.countplot(x='PurchaseYear', hue='IsBadBuy', data=data)

Out[77]: <AxesSubplot:xlabel='PurchaseYear', ylabel='count'>



In [78]: # Zakljuujemo da ne postoji znacajna razlika u kvalitetu kupovine u odnosu na ove dve godine, ali ce nam kombi
# ovog atributa sa nekim drugim morda dati bolje rezultate, tako da cemo je za sada ostaviti u datasetu

In [79]: #VehicleAge

In [80]: data['VehicleAge'].value_counts()

4 1598
3 1453
# 1214
2 792
6 753
# 429
1 295
8 207
# 66
Name: VehicleAge, dtype: int64

In [81]: sns.countplot(x='VehicleAge', hue='IsBadBuy', data=data)

Out[81]: <AxesSubplot:xlabel='VehicleAge', ylabel='count'>



In [82]: # Vidimo da razlicita starost vozila ima drugaciji uticaj na kupovinu, pa mozemo da zakljucimo da nam je ovaj
# znacajan za dalju analizu

In [83]: # Make

In [84]: data['Make'].value_counts()

CHEVROLET 1574
DODGE 1187
FORD 1059
CHRYSLER 827
NISSAN 585
KIA 232
SATURN 211
ACURA 184
HYUNDAI 181
JEEP 150
SUBARU 132
MAZDA 112
TOYOTA 109
MITSUBISHI 91
MERCURY 83
BUICK 61
HONDA 52
OLDSMOBILE 22
VOLKSWAGEN 14
ISUZU 8
LINCOLN 7
PONTIAC 6
MINI 5
SCION 5
VOLVO 4
ACURA 4
SUBARU 1
CADILLAC 1
HUMMER 1
PLYMOUTH 1
LEXUS 1
Name: Make, dtype: int64

In [85]: # Model

In [86]: data['Model'].value_counts()

PT CRUISER 216
IMPALA 194
BENTLEY 128
CARAVAN GRAND FWD V6 118
CALIBER 117

CMRV V6 3.0L / 3.3L ...
MARINER 2WD 4C 1
STARLINE T73 1
VERONA 1
PATRIOT 2WD 4C 2.0L 1
Name: Model, Length: 632, dtype: int64

In [87]: pd.crosstab([data['Model'], data['SubModel']],sum())

Out[87]: SubModel
2D CONVERTIBLE 30
2D CONVERTIBLE CLS 3
2D CONVERTIBLE GTC 2
2D CONVERTIBLE TOURING 7
2D COUPE 102

WAGON 3.5I SXT 2
WAGON DX 1
WAGON LX 1
WAGON SXT 1
Length: 465, dtype: int64

In [88]: # Vidimo da za obe varijable imamo veliki broj razlicitih vrednosti, pa cemo obe varijable iskljuciti iz dalje

In [89]: data.drop(['Model','SubModel'], axis=1, inplace=True)

In [90]: data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6797 entries, 1 to 6798
Data columns (total 20 columns):
# Column Non-Null Count Dtype
---
0 VehicleAge 6797 non-null int64
1 Make 6797 non-null object
2 Color 6797 non-null object
3 Transmission 6797 non-null object
4 WheelType 6797 non-null object
5 VehOdo 6797 non-null int64
6 Size 6797 non-null object
7 MMRAcquisitionAuctionAveragePrice 6797 non-null float64
8 MMRAcquisitionAuctionCleanPrice 6797 non-null float64
9 MMRAcquisitionRetailAveragePrice 6797 non-null float64
10 MMRAcquisitionRetailCleanPrice 6797 non-null float64
11 MMRCurrentAuctionAveragePrice 6797 non-null float64
12 MMRCurrentAuctionCleanPrice 6797 non-null float64
13 MMRCurrentRetailAveragePrice 6797 non-null float64
14 MMRCurrentRetailCleanPrice 6797 non-null float64
15 VehCost 6797 non-null int64
16 IsOnlineSale 6797 non-null int64
17 WarrantyCost 6797 non-null int64
18 IsBadBuy 6797 non-null int64
19 PurchaseYear 6797 non-null int64
dtypes: float64(9), int64(6), object(5)
memory usage: 1.1+ MB

In [91]: # Color

In [92]: data['Color'].value_counts()

SILVER 1379
WHITE 1118
BLACK 930
BLUE 735
GREY 710
RED 570
GOLD 504
GREEN 297
KAWASAKI 210
BEIJING 170
BROWN 52
CORLE 32
CRANGE 29
YELLOW 28
OTHER 28
NOT AVAIL 5
Name: Color, dtype: int64

In [93]: sns.countplot(x='Color', hue='IsBadBuy', data=data)

Out[93]: <AxesSubplot:xlabel='Color', ylabel='count'>



In [94]: data.loc[data['Color'] == 'NOT AVAIL','Color']

Out[94]: RefId      NOT AVAIL
4072 NOT AVAIL
4160 NOT AVAIL
4910 NOT AVAIL
4912 NOT AVAIL
Name: Color, dtype: object

In [95]: # Vidimo da za 5 opservacija imamo vrednost not avail, sto moramo zameniti
# Medjutim, necemo menjati najstapljenijom bojom, vec cemo ih dodeliti vrednost other

In [96]: data.loc[data['Color'] == 'NOT AVAIL','Color'] = 'OTHER'

In [97]: #Transmission

In [98]: data['Transmission'].value_counts()

AUTO 6546
MANUAL 251
Name: Transmission, dtype: int64

In [99]: sns.countplot(x='Transmission', hue='IsBadBuy', data=data)

Out[99]: <AxesSubplot:xlabel='Transmission', ylabel='count'>


```



```
In [284]: data_original.loc[:, 'IsBadBuy'] = data_original.loc[:, 'IsBadBuy'] * (-1)

In [285]: x_o_n = data_original.loc[:, data_original.columns != 'IsBadBuy']
x_o_n = data_original.loc[:, 'IsBadBuy']

In [286]: # Moramo i drugi dataset sa svim KMR podacima da podelimo na trening i test

In [287]: from sklearn.model_selection import train_test_split
x_train_o, x_test_o, y_train_o, y_test_o = train_test_split(x_o_n, y_o_n, test_size = 0.3, random_state=10)

In [288]: lista1 = list()
for model in (model_nb, model_knn, model_dt, model_log, model_voting, model_bagging, model_rf):
    scores = cross_val_score(model, x_train_o, y_train_o, cv=10, scoring='roc_auc')
    lista1.append(round(np.mean(scores)*100, 2))
data_alg_rez['Potpuni skup podataka'] = lista1

In [289]: data_alg_rez.sort_values

Out[289]: <bound method DataFrame.sort_values of
Naivni Bajes      50.34      59.91      14.36
KNN              55.23      55.11      85.68
Stablo odlucivanja  51.51      52.55      78.05
Logisticka regresija  67.14      67.20      87.16
Voting           61.01      61.73      87.18
Bagging          66.96      67.03      87.16
Random forest    62.77      63.36

In [210]: # Dalje cemo raditi sa potpunim skupom podataka, jer vidimo da nema znacajne razlike u rezultatima

In [276]: lista2 = list()
for model in (model_nb, model_knn, model_dt, model_log, model_voting, model_bagging, model_rf):
    scores = cross_val_score(model, x_train_o, y_train_o, cv=10, scoring='accuracy')
    lista2.append(round(np.mean(scores)*100, 2))
data_alg_rez['Potpuni skup podataka sa ACC'] = lista2

In [277]: data_alg_rez

Out[277]:
```

	inicijalni skup podataka	Potpuni skup podataka	Potpuni skup podataka sa ACC
Algoritam			
Naivni Bajes	50.34	59.91	14.36
KNN	55.23	55.11	85.68
Stablo odlucivanja	51.51	52.55	78.05
Logisticka regresija	67.14	67.20	87.16
Voting	61.01	61.73	87.18
Bagging	66.96	67.03	87.16
Random forest	62.77	63.36	87.11

```
In [213]: model_voting2 = VotingClassifier(voting='soft', estimators=[('knn',model_knn), ('tree', model_dt), ('log', model_log), ('voting', model_voting)])

In [214]: cross_val_score(model_voting2, x_o_n, y_o_n, cv=10, scoring='accuracy').mean()
0.8540561379190853

Out[214]:
```

In [215]: # Vidimo da nam Voting radi bolje ukoliko isbacimo Naivnog Bajesa jer on pretpostavlja da su promenljive medjus se a znano da imamo 3 varijabli koje su bas jako korelirane

```
In [218]: from sklearn.classifier import StackingClassifier

In [219]: model_stacking = StackingClassifier(classifiers=[model_knn, model_dt, model_log], meta_classifier=LogisticRegression())

In [220]: cross_val_score(model_stacking, x_o_n, y_o_n, cv=10, scoring='roc_auc').mean()
0.64124322262424308

Out[220]:
```

```
In [221]: cross_val_score(model_stacking, x_o_n, y_o_n, cv=10, scoring='accuracy').mean()
0.7785792255046349

Out[221]:
```

In [222]: # Vidimo da nam najbolje rezultate daje logisticka regresija i KNN, kao i kombinacija ova dva algoritma i drve # koristeci Voting i Stacking, ali cemo Voting isbaciti is dalje analize i ostalim modelima uraditi predviđanj # kako bismo poredili rezultate i odlucile se za najbolji model

Predviđanja

```
In [224]: model_log.fit(X_train_o, Y_train_o)
predictions_log = model_log.predict(X_train_o)
pd.DataFrame({'stvarno':Y_train_o, 'predvidjeno':predictions_log}).head(13)

Out[224]:
```

stvarno	predvidjeno
554	1.0
5591	1.0
2516	1.0
3588	-0.0
3760	1.0
53	1.0
1775	1.0
2950	1.0
4738	1.0
3012	1.0
2137	1.0
3226	1.0
3253	1.0

```
In [315]: model_knn.fit(X_train_o, Y_train_o)
predictions_knn = model_knn.predict(X_train_o)
pd.DataFrame({'stvarno':Y_train_o, 'predvidjeno':predictions_knn}).head(10)

Out[315]:
```

stvarno	predvidjeno
554	1.0
5591	1.0
2516	1.0
3588	-0.0
3760	1.0
53	1.0
1775	1.0
2950	1.0
4738	1.0
3012	1.0
2137	1.0
3226	1.0
3253	1.0

```
In [226]: model_stacking.fit(X_train_o, Y_train_o)
predictions_stacking = model_stacking.predict(X_train_o)
pd.DataFrame({'stvarno':Y_train_o, 'predvidjeno':predictions_stacking}).head(13)

Out[226]:
```

stvarno	predvidjeno
554	1.0
5591	1.0
2516	1.0
3588	-0.0
3760	1.0
53	1.0
1775	1.0
2950	1.0
4738	1.0
3012	1.0
2137	1.0
3226	1.0
3253	1.0

```
In [227]: from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, classification_report, confusion_matrix

In [228]: print('Recall logisticka regresija: ', recall_score(Y_test_o, model_log.predict(X_test_o)))
print('Recall KNN: ', recall_score(Y_test_o, model_knn.predict(X_test_o)))
print('Recall Stacking: ', recall_score(Y_test_o, model_stacking.predict(X_test_o)))

print('Precision logisticka regresija: ', precision_score(Y_test_o, model_log.predict(X_test_o)))
print('Precision KNN: ', precision_score(Y_test_o, model_knn.predict(X_test_o)))
print('Precision Stacking: ', precision_score(Y_test_o, model_stacking.predict(X_test_o)))

print('F1 logisticka regresija: ', f1_score(Y_test_o, model_log.predict(X_test_o)))
print('F1 KNN: ', f1_score(Y_test_o, model_knn.predict(X_test_o)))
print('F1 Stacking: ', f1_score(Y_test_o, model_stacking.predict(X_test_o)))

Recall logisticka regresija: 0.48878923766816146
Recall KNN: 0.9854260896860999
Recall Stacking: 0.8806053801459192
Precision logisticka regresija: 0.8745098039215686
Precision KNN: 0.9854260896860999
Precision Stacking: 0.877245089820359
F1 logisticka regresija: 0.68365650965291
F1 KNN: 0.928194297782471
F1 Stacking: 0.6810392708917554
```

```
In [229]: # Menjamo granicu odlucivanja jer nam dijablena klasa izlaze promjenljive znatno utice na predviđanje

In [230]: prob_log = model_log.predict_proba(X_test_o[:, 1]) >= 0.9
prob_knn = model_knn.predict_proba(X_test_o[:, 1]) >= 0.9
prob_stacking = model_stacking.predict_proba(X_test_o[:, 1]) >= 0.9

In [231]: print('Recall logisticka regresija: ', recall_score(Y_test_o, prob_log))
print('Recall KNN: ', recall_score(Y_test_o, prob_knn))
print('Recall Stacking: ', recall_score(Y_test_o, prob_stacking))

print('Precision logisticka regresija: ', precision_score(Y_test_o, prob_log))
print('Precision KNN: ', precision_score(Y_test_o, prob_knn))
print('Precision Stacking: ', precision_score(Y_test_o, prob_stacking))

print('F1 logisticka regresija: ', f1_score(Y_test_o, prob_log))
print('F1 KNN: ', f1_score(Y_test_o, prob_knn))
print('F1 Stacking: ', f1_score(Y_test_o, prob_stacking))

Recall logisticka regresija: 0.48878923766816146
Recall KNN: 0.9854260896860999
Recall Stacking: 0.8806053801459192
Precision logisticka regresija: 0.8745098039215686
Precision KNN: 0.9854260896860999
Precision Stacking: 0.877245089820359
F1 logisticka regresija: 0.68365650965291
F1 KNN: 0.928194297782471
F1 Stacking: 0.6810392708917554
```

```
In [232]: confusion_matrix(Y_test_o, model_log.predict(X_test_o))
array([[ 0, 256],
       [ 0, 1784]], dtype=int64)

In [233]: confusion_matrix(Y_test_o, model_knn.predict(X_test_o))
array([[ 10, 246],
       [ 26, 1758]], dtype=int64)

Out[233]:
```

```
In [234]: confusion_matrix(Y_test_o, model_stacking.predict(X_test_o))
array([[ 45, 211],
       [ 213, 1571]], dtype=int64)

Out[234]:
```

In [235]: # Uocavamo da smo dobili zadovoljavajuce rezultate za precision sto nam je svakako bitnije i mozemo videti da u regresiji definitivno radi najbolje

In [236]: # Posto su logisticka regresija i Stacking dali dobre rezultate, sad cemo izvrsti optimizaciju parametara sa stablo odlucivanja i KNN kako bismo i njih pokušale da poboljšamo

```
In [284]: data_alg_rez

Out[284]:
```

	inicijalni skup podataka	Potpuni skup podataka	Potpuni skup podataka sa ACC
Algoritam			
Naivni Bajes	50.34	59.91	14.36
KNN	55.23	55.11	85.68
Stablo odlucivanja	51.51	52.55	78.05
Logisticka regresija	67.14	67.20	87.16
Voting	61.01	61.73	87.18
Bagging	66.96	67.03	87.16
Random forest	62.77	63.36	87.11

```
In [285]: selection = VarianceThreshold(0.05)
selection.fit(x_o_n)
x_o_filter = x_o_n.loc[:, selection.get_support()]

In [286]: x_o_filter

Out[286]:
```

	Year	VehicleAge	WheelType	Make_CHEVROLET	Make_CHRYSLER	Make_DODGE	Make_FORD	Make_PONTIAC	Color_BLACK	Color_WHITE
0	0.0	0.428571	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0
1	0.0	0.571429	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0
2	0.0	0.428571	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
3	0.0	0.571429	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.0	0.571429	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0
...
6792	0.0	0.285714	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
6793	0.0	0.428571	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
6794	0.0	0.857143	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
6795	0.0	0.285714	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
6796	0.0	0.857143	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0
6797 rows x 20 columns										

```
In [287]: x_train_f, x_test_f, y_train_f, y_test_f = train_test_split(x_o_filter, y_o_n, test_size = 0.3, random_state=10)

In [288]: lista5 = list()
for model in (model_nb, model_knn, model_dt, model_log, model_voting, model_bagging, model_rf):
    scores = cross_val_score(model, x_train_f, y_train_f, cv=10, scoring='roc_auc')
    lista5.append(round(np.mean(scores)*100, 2))
data_alg_rez['Potpuni skup podataka filter'] = lista5

In [290]: lista6 = list()
for model in (model_nb, model_knn, model_dt, model_log, model_voting, model_bagging, model_rf):
    scores = cross_val_score(model, x_train_f, y_train_f, cv=10, scoring='accuracy')
    lista6.append(round(np.mean(scores)*100, 2))
data_alg_rez['Potpuni skup podataka filter ACC'] = lista6

In [291]: data_alg_rez

Out[291]:
```

	inicijalni skup podataka	Potpuni skup podataka	Potpuni skup podataka sa ACC	Potpuni skup podataka filter	Potpuni skup podataka filter ACC
Algoritam					
Naivni Bajes	50.34	59.91	14.36	61.44	81.46
KNN	55.23	55.11	85.68	56.20	85.31
Stablo odlucivanja	51.51	52.55	78.05	53.12	77.63
Logisticka regresija	67.14	67.20	87.16	65.53	87.18
Voting	61.01	61.73	87.18	63.25	87.18
Bagging	66.96	67.03	87.16	65.45	87.18
Random forest	62.77	63.36	87.11	59.78	83.18

Optimizacija parametara

```
In [311]: knn_params = {'n_neighbors': [1,2,3,4,5,6,7,8,9,10,12,20,30,50,100]}
grid = GridSearchCV(model_knn, knn_params, cv=10, scoring='roc_auc')
grid.fit(X_train_f, Y_train_f)

print('Best param: ', grid.best_params_)
Best param: {'n_neighbors': 100}

In [301]: max_depths = {'max_depth': np.linspace(1, 40, 40, endpoint=True)}

In [302]: grid = GridSearchCV(model_dt, max_depths, cv=10, scoring='roc_auc')
grid.fit(X_train_f, Y_train_f)

print('Best param: ', grid.best_params_)
Best param: {'max_depth': 2.0}

In [305]: rez_optimizovano = pd.DataFrame()
rez_optimizovano['rez_optimizacije'] = data_alg_rez.iloc[:,3]
rez_optimizovano

Out[305]:
```

	Pre optimizacije
Algoritam	
Naivni Bajes	61.44
KNN	56.20
Stablo odlucivanja	53.12
Logisticka regresija	65.53
Voting	63.25
Bagging	65.45
Random forest	59.78

```
In [312]: model_knn_new = KNeighborsClassifier(n_neighbors=100)
model_dt_new = DecisionTreeClassifier(max_depth=2)
model_voting2 = VotingClassifier(voting='soft', estimators=[('knn',model_knn_new), ('tree', model_dt_new), ('log', model_log)])
lista_new = list()
scores = cross_val_score(model, x_train_f, y_train_f, cv=10, scoring='roc_auc')
lista_new.append(round(np.mean(scores)*100, 2))
rez_optimizovano['Nakon optimizacije'] = lista_new

In [313]: rez_optimizovano

Out[313]:
```

	Pre optimizacije	Nakon optimizacije
Algoritam		
Naivni Bajes	61.44	61.44
KNN	56.20	61.08
Stablo odlucivanja	53.12	63.19
Logisticka regresija	65.53	65.53
Voting	63.25	65.12
Bagging	65.45	65.38
Random forest	59.78	59.01

In [316]: # Analizom svih rezultata dolazimo do zakljucka da nam Logisticka regresija daje najbolje rezultate i da cemo R kao model koji cemo koristiti za predviđanje