

Analiza podataka

- ➔ 18 atributa ima NA vrednosti u celom DataSetu-u.
- ➔ **IsBadBuy je IZLAZ** : da li je kupovina auta pametan izbor za dalju preprodaju.

Izdvajanje atributa:

- **PurchDate** (godina kada je auto stavljen na aukciju) ima godinu u sebi, izdvojili smo godine u poseban atribut jer su veoma bitne za odluku. Pravimo atribut **PurchYear** . Postoje 2 godine 2009 i 2010 god aukcija, i relativno jednak broj automobile u obe godine.
- **Auction** se odnosi na dobavljača automobila i ima nedostajuće vrednosti i to 3815. Dakle preko pola ima NA i još oko 1598 auta ima dobavljača OTHER, **sto ovaj atribut cini neprakticnim za dalju analizu.**
- **VehYear** predstavlja godinu proizvodnje automobila a **VehicleAge** predstavlja starost automobile u godinama. **Ova dva atributa su visoko korelisana, te izbacujemo jednog od njih.**
- **Make** predstavlja marku automobila, i imamo 32 marke automobila. Pored modela automobila imamo i **model** automobila i to 632 modela. To je previše, te želimo da ovo svedemo na manji broj automobila da ne dodje do pretreniranosti. Vidimo da u **okviru naziva modela automobila postoji i oznaka za vrstu pogona i broj cilindara na motoru.** Postoje 3 vrste oznaka – AWD, 4WD i 2WD za vrstu pogona kao i V8, V6, 6C, i6 za broj cilindara. Shodno tome, **kreirali smo još 2 atributa DriveType i NumCyls .**
- **Model** koji sadrži u sebi već ova 2 podatka (pogon i broj cilindara), filtriramo i ostavljamo samo naziv modela. Smanjujemo broj vrednosti za model sa prethodno 632 na samo 185.
- **Trim** predstavlja nivo kvaliteta opreme automobila. Kako je rec o kategorickom atributu popunjavamo nedostajuće vrednosti **onom vrednoscu koja se NAJCESCE pojavljuje i za to se koristi ova funkcija mode()**.
- **SubModel** atribut ima 465 jedinstvenih vrednosti i zbog tog velikog broja atribut se izostavlja iz dalje analize.
- **Color i Transmission** su samo plotovani i iskazan je neki odnos izmedju njih.
- **WheelTypeID** se odnosi na to da li točak ima **poklopac ili felnu** i ima vrednosti 1 i 2 a **WheelType** se odnosi na samu vrstu točka. Kako se prvi atribut odnosi isključivo na ID izostavlja se iz dalje analize (**WheelTypeID**). *2 je Covers a 1 je Alloy* bilo Relativno sličan odnos ove dve vrednosti 3376 i 3083. **Kako alloy ima veću frekvenciju NA se popunjavaju tom vrednoscu.** Dakle, 339 Na vrednosti se popunjava alloy (3376 > 3083).
- **VehOdo** se odnosi na predjenu kilometražu a **Nationality** se odnosi na zemlju porekla automobila, i ima JEDNU NA vrednost koju popunjavamo najfrekventnijom nacionalnoscu, a to je AMERICAN.
- **Crosstab-ovajem MAKE (MARKA) i NATIONALITY (ZEMLJA POREKLA) vidimo da su ova dva atributa VISOKO KORELISANA i jednog mozemo da izbacimo iz dalje analize.**
- **Size i TopThreeAmerican** atributa se dalje analiziraju. Size ima 1 NA vrednost i pripada marki JEEP, pa gledamo koja je najfrekventnija vrednost ovog atributa za JEEP marku i vidimo da je to **Small SUV** i na mesto NA stavljamo upravo to.

- **TopThreeAmerican** atribut ima 1 NA vrednost marke JEEP I model PATRIOT. Opet, crosstabovanjem vidimo da je najcesci za TopThreeAmerican za JEEP upravo CHRYSLER.
- **Medjutim MARKA I TOPTHREEAMERICAN su visoko korelisana te jedan izbacujemo iz dalje analize.**
- Dalje analiziraju kretanje cena po razlicitim osnovama.
- Sve NA vrednosti popunjavamo **medijanom tog atributa (kolone).**
- **PRIMEUNIT I AUCGUART** imaju veliki broj NA vrednosti te se ova dva atributa izostavljaju iz dalje analize.
- **BYRNO I VNZIP1** se odnose na sifre I kodove pa se takodje izostavljaju.
- **VNST** se odnosi na zemlju porekla kupca.
- **Plotovanjem IsBadBuy I VehichleAge vidimo da sto je auto stariji to je kupovina bila gora odluka (logicno).** *Na ovom grafiku se uočava da starost automobila utiče na to da bude klasifikovan kao nekvalitetan. Što je automobil mlađi manja je verovatnoća da će se klasifikovati kao nekvalitetan.*
- **Dakle, izostavljaju se atributi**

```
data = data.drop(columns=['RefId', 'PurchDate', 'Auction', 'VehYear', 'WheelTypeID', 'Nationality',
'TopThreeAmericanName',
'PRIMEUNIT', 'AUCGUART', 'BYRNO', 'VNZIP1', 'VNST', 'SubModel', 'Trim'])
```

Ubacujemo dummie vrednosti za Kategoricke attribute (binarne I nebinarne).

```
df_dummies = pd.get_dummies(data, columns = ['Make', 'Model', 'Color', 'Size', 'DriveType'])
pd.set_option('display.max_columns', 200)
df_dummies
```

```
df_dummies = df_dummies.replace({'Transmission' : {'MANUAL':0, 'AUTO':1}})
df_dummies = df_dummies.replace({'WheelType' : {'Covers':0, 'Alloy':1}})
```

Promena tipa izlaza – IsBadBuy

Posto je int izlaz menjamo ga u bool atribut.

```
df_dummies['IsBadBuy'] = df_dummies['IsBadBuy'].astype('bool')
```

Dakle sada nije 1 I 0 vec True I False.