

Used Car Price Prediction

FINAL SEMESTER PROJECT - SUMMER 2021

Syed Muhammad Sahib

Project Statement

For the US Cars Dataset, the analysis pertains to predict the price of the car based on the appropriate features.

Dataset Information

The dataset is from the Kaggle Repository and was scraped from auctionexport.com. The dataset includes information about 28 brands of clean and used vehicles located in the US. The dataset includes 12 features: price, years, brand, model, color, state, mileage, VIN, title status, lot, and the condition.

Feature Description

- Price: The sale price of the vehicle
- Year: Vehicle registration year
- Brand: Brand/Make of the vehicle
- Model: Model of the vehicle
- Color: Color of the vehicle
- State: The location where the vehicle is for sale
- Mileage: Miles traveled by the vehicle
- VIN: Vehicle identification number
- Title status: Clean or salvaged title
- lot: lot number from the manufacturer
- condition: time

Library Import

```
In [92]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from numpy import mean
from numpy import std
from pandas import read_csv

from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.preprocessing import Normalizer

import xgboost as xgb

from sklearn.model_selection import KFold, GridSearchCV
from sklearn.metrics import accuracy_score, make_scorer
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import StratifiedKFold
from sklearn.metrics import roc_curve, precision_recall_curve
from sklearn.metrics import f1_score, precision_score, recall_score
from sklearn.metrics import classification_report, confusion_matrix, auc

from sklearn.linear_model import LinearRegression
from sklearn.neighbors import KNeighborsRegressor
from sklearn.svm import LinearSVR, SVR
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from xgboost import XGBRegressor
```

Loading the Dataset

```
In [131]: df = pd.read_csv("C:/Users/smaab/Desktop/UCF/STA - 6704/dataset/USA_cars_datasets.csv")
df

Out[131]:
```

Unnamed: 0	price	brand	model	year	title_status	mileage	color	vin	lot	state	country	condition	
0	0	6300	toyota	cruiser	2008	clean vehicle	274117.0	black	1jctz1168K007763	159348797	new jersey	usa	10 days
1	1	2899	ford	se	2011	clean vehicle	190552.0	silver	2fmdk3gc0b0602217	166851262	tennessee	usa	6 days
2	2	5350	dodge	mpv	2018	clean vehicle	39590.0	silver	3c4pdgg9j9346413	167655728	georgia	usa	2 days
3	3	25000	ford	door	2014	clean vehicle	64146.0	blue	1ftw1e1ef4e623745	167733655	virginia	usa	22 hours
4	4	27700	chevrolet	1500	2018	clean vehicle	6654.0	red	3gcpncr2g473991	167763266	florida	usa	22 hours
...
2484	2484	7800	nissan	versa	2019	clean vehicle	23609.0	red	3n1cn7ap9k880319	167722215	california	usa	1 day
2495	2495	9200	nissan	versa	2018	clean vehicle	34553.0	silver	3n1cn7ap9j884088	167722225	florida	usa	21 hours
2496	2496	9200	nissan	versa	2018	clean vehicle	31594.0	silver	3n1cn7ap9j884191	167762226	florida	usa	21 hours
2497	2497	9200	nissan	versa	2018	clean vehicle	32557.0	black	3n1cn7ap3j883263	167762227	florida	usa	2 days
2498	2498	9200	nissan	versa	2018	clean vehicle	31371.0	silver	3n1cn7ap4j884311	167762228	florida	usa	21 hours

2499 rows x 13 columns

Exploratory Data Analysis

Based on the knowledge of cars and in my experience, the price of a used car usually depends on the brand, model, year, title status, mileage, color, and what state the car is located in. For the purpose of this analysis, lets drop the VIN, lot number, country, and condition column.

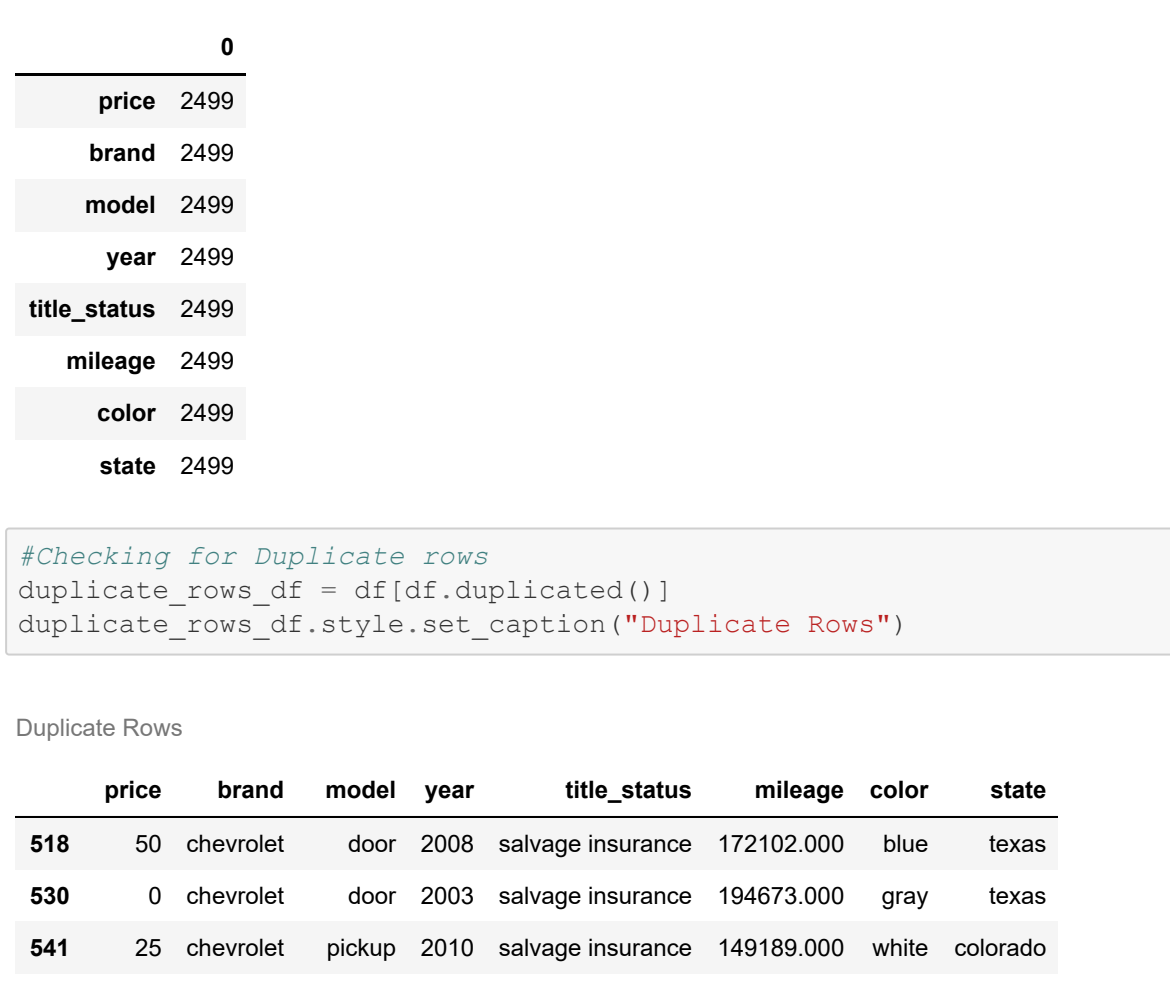
```
In [132]: #dropping the Unrelevant columns
df = df.drop(["Unnamed: 0", "vin", "lot", "country", "condition"], axis=1)
```

```
In [133]: df.head(10).style.set_properties(**{'background-color': 'black',
                                           'color': 'lightgreen',
                                           'border-color': 'white'})
```

```
Out[133]:
```

	price	brand	model	year	title_status	mileage	color	state
0	6300	toyota	cruiser	2008	clean vehicle	274117.000	black	new jersey
1	2899	ford	se	2011	clean vehicle	190552.000	silver	tennessee
2	5350	dodge	mpv	2018	clean vehicle	39590.000	silver	georgia
3	25000	ford	door	2014	clean vehicle	64146.000	blue	virginia
4	27700	chevrolet	1500	2018	clean vehicle	6654.000	red	florida
5	81900	dodge	mpv	2018	clean vehicle	45681.000	white	txas
6	7300	chevrolet	pk	2010	clean vehicle	149050.000	black	california
7	13550	gmc	door	2017	clean vehicle	232525.000	gray	georgia
8	14600	chevrolet	malibu	2018	clean vehicle	53711.000	silver	florida
9	5250	ford	mpv	2017	clean vehicle	63416.000	black	txas

```
In [134]: #lets look at the spread of the initial data
boxplot_data = df.plot.box(title="Box and Whisker Plot")
plt.show()
```



At this point, the data is raw, and not been pre-processed. The box and whisker plot will not show promising results.

Investigating Numeric and Categorical Variables

```
In [135]: #Type of the variables
df.dtypes
```

```
Out[135]: price          int64
brand            object
model            object
year            int64
title_status     object
mileage         float64
color            object
state           object
dtype: object
```

```
In [136]: #Descriptive Statistics
pd.set_option('precision',3)
df.describe(include='all')
```

```
Out[136]:
```

	price	brand	model	year	title_status	mileage	color	state
count	2499.000	2499	2499	2499.000	2499	2.499e+03	2499	2499
unique	NaN	28	127	NaN	2	NaN	40	44
top	NaN	ford	door	NaN	clean vehicle	NaN	white	pennsylvania
freq	NaN	1235	651	NaN	2336	NaN	707	299
mean	18767.671	NaN	NaN	2016.714	NaN	5.230e+04	NaN	NaN
std	12116.095	NaN	NaN	3.443	NaN	5.971e+04	NaN	NaN
min	0.000	NaN	NaN	1973.000	NaN	0.000e+00	NaN	NaN
25%	10200.000	NaN	NaN	2016.000	NaN	2.147e+04	NaN	NaN
50%	16900.000	NaN	NaN	2018.000	NaN	3.536e+04	NaN	NaN
75%	25555.000	NaN	NaN	2019.000	NaN	6.347e+04	NaN	NaN
max	84900.000	NaN	NaN	2020.000	NaN	1.018e+06	NaN	NaN

```
In [137]: # Checking for unique columns
print("\n\nUnique column values in dataset\n", df.nunique())
```

Unique column values in dataset
price 790
brand 28
model 127
year 30
title_status 2
state 44
mileage 2439
color 49
state 44
dtype: int64

Categorical Variables are brand, model, title, color, and state.

Checking for duplicate rows and deleting them

At this point, lets delete any duplicate rows.

```
In [138]: #Count of the variables before deleting duplicate rows
count_before_removing_dup = pd.DataFrame(df.count())
count_after_removing_dup.style.set_caption("Count of variables before removing duplicates")
```

```
Out[138]:
```

	price	brand	model	year	title_status	mileage	color	state
0	2499	2499	2499	2499.000	2499	2.499e+03	2499	2499
brand	2499	2499	2499	2499.000	2499	2.499e+03	2499	2499
model	2499	2499	2499	2499.000	2499	2.499e+03	2499	2499
year	2499	2499	2499	2499.000	2499	2.499e+03	2499	2499
title_status	2499	2499	2499	2499.000	2499	2.499e+03	2499	2499
mileage	2499	2499	2499	2499.000	2499	2.499e+03	2499	2499
color	2499	2499	2499	2499.000	2499	2.499e+03	2499	2499
state	2499	2499	2499	2499.000	2499	2.499e+03	2499	2499

```
In [139]: #Checking for Duplicate rows
duplicate_rows_df = df[df.duplicated()]
duplicate_rows_df.style.set_caption("Duplicate Rows")
```

```
Out[139]:
```

	price	brand	model	year	title_status	mileage	color	state
518	50	chevrolet	door	2008	salvage insurance	172102.000	blue	txas
530	0	chevrolet	door	2010	salvage insurance	194673.000	gray	txas
541	25	chevrolet	pickup	2003	salvage insurance	149189.000	white	colorado
661	26900	chevrolet	traverse	2018	clean vehicle	42941.000	black	missouri

```
In [140]: print("number of duplicate rows: ", duplicate_rows_df.shape)
```

number of duplicate rows: (4, 8)

Deleting 4 Duplicate rows

```
In [141]: #Deleting the Duplicate rows
df = df.drop_duplicates()
duplicate_rows_df = df[df.duplicated()]
print("number of duplicate rows: ", duplicate_rows_df.shape)
```

number of duplicate rows: (0, 8)

```
In [143]: #Count of the data after deleting duplicate rows
count_after_removing_dup = pd.DataFrame(df.count())
count_after_removing_dup.style.set_caption("Count of variables after removing duplicates")
```

```
Out[143]:
```

	price	brand	model	year	title_status	mileage	color	state
0	2495	2495	2495	2495.000	2495	2.495e+03	2495	2495
brand	2495	2495	2495	2495.000	2495	2.495e+03	2495	2495
model	2495	2495	2495	2495.000	2495	2.495e+03	2495	2495
year	2495	2495	2495	2495.000	2495	2.495e+03	2495	2495
title_status	2495	2495	2495	2495.000	2495	2.495e+03	2495	2495
mileage	2495	2495	2495	2495.000	2495	2.495e+03	2495	2495
color	2495	2495	2495	2495.000	2495	2.495e+03	2495	2495
state	2495	2495	2495	2495.000	2495	2.495e+03	2495	2495

```
In [12]: #New shape of the dataset
print(" dataset: ( rows, columns) = ", df.shape)
```

dataset: (rows, columns) = (2495, 8)

Checking for any missing values

```
In [13]: #Finding the null values
print(df.isnull().sum())
```

```
price          0
brand          0
model          0
year          0
title_status   0
mileage        0
color          0
state          0
dtype: int64
```

There are no missing values

Modifying the Year Variable

```
In [14]: #Fixing the year variable to where it shows how old the car is
def mod_year(x):
    return 2021 - x
df['year'] = mod_year(df.year)
df

<ipython-input-14-67f72da5d522>4: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
g.mod_year(df.year) = mod_year(df.year)
```

```
Out[14]:
```

	price	brand	model	year	title_status	mileage	color	state
0	6300	toyota	cruiser	13	clean vehicle	172102.000	black	new jersey
1	2899	ford	se	10	clean vehicle	190552.0	silver	tennessee
2	5350	dodge	mpv	3	clean vehicle	39590.0	blue	georgia
3	25000	ford	door	7	clean vehicle	64146.0	blue	virginia
4	27700	chevrolet	1500	3	clean vehicle	6654.0	red	florida
...
2484	7800	nissan	versa	2	clean vehicle	23609.0	red	california
2495	9200	nissan	versa	3	clean vehicle	34553.0	silver	florida
2496	9200	nissan	versa	3	clean vehicle	31594.0	silver	florida
2497	9200	nissan	versa	3	clean vehicle	32557.0	black	florida
2498	9200	nissan	versa	3	clean vehicle	31371.0	silver	florida

2495 rows x 8 columns

Visualize the top brands of the vehicle by count

```
In [15]: brand_count = df['brand'].value_counts().to_frame().reset_index()
brand_count.rename(columns={'index':'brand', 'brand':'count'}, inplace=True)
fig = plt.figure(figsize=(10, 5))
axes = fig.add_axes([0, 0, 1, 1])
sns.barplot(y = brand_count['brand'], x = brand_count['count'], ax=axes)
```

```
Out[15]: <matplotlib.axes._axes.Axes at 0x1ee46c1640>
```

Since we have a larger number of vehicles that are Ford, Dodge, Nissan, Chevrolet, and GMC, lets only take the datapoints that incorporate the top 5 brands for the dataset

```
In [16]: # Set the index of the DataFrame to the brand name
df1 = df.set_index("brand")
#Drop the brands other than top 5
df1 = df1.drop(["jeep", "chrysler", "dodge", "hyundai", "buick", "kia", "honda", "infiniti", "cadillac",
               "mercedes-benz", "nissan", "audi", "land", "peterbilt", "acura", "lincoln", "mazda",
               "lexus", "jaguar", "ram", "maserati", "harley-davidson", "toyota"])

In [17]: #Reset the index to have brand as a column/feature
df1.reset_index(level=0, inplace=True)
```

```
In [18]: #Visualize the shape of the dataset
df1
```

```
Out[18]:
```

	brand	price	model	year	title_status	mileage	color	state
0	ford	2899	se	10	clean vehicle	190520.0	silver	tennessee
1	dodge	5350	mpv	3	clean vehicle	39590.0	silver	georgia
2	ford	25000	door	7	clean vehicle	64146.0	blue	virginia
3	chevrolet	27700	1500	3	clean vehicle	6654.0	red	florida
4	dodge	5700	mpv	3	clean vehicle	45681.0	white	txas
...
2309	nissan	7800	versa	2	clean vehicle	23609.0	red	california
2310	nissan	9200	versa	3	clean vehicle	34553.0	silver	florida
2311	nissan	9200	versa	3	clean vehicle	31594.0	silver	florida
2312	nissan	9200	versa	3	clean vehicle	32557.0	black	florida
2313	nissan	9200	versa	3	clean vehicle	31371.0	silver	florida

2314 rows x 8 columns

```
In [19]: brand_count = df1['brand'].value_counts().to_frame().reset_index()
brand_count.rename(columns={'index':'brand', 'brand':'count'}, inplace=True)
fig = plt.figure(figsize=(10, 5))
axes = fig.add_axes([0, 0, 1, 1])
sns.barplot(y = brand_count['brand'], x = brand_count['count'], ax=axes)
```

```
Out[19]: <matplotlib.axes._axes.Axes at 0x1ee46c1640>
```

```
In [20]: #Number of records removed
no_of_records_removed = df.count() - df1.count()
no_of_records_removed
```

```
Out[20]: brand          181
color          181
mileage        181
model          181
price          181
state          181
title_status   181
year          181
dtype: int64
```

181 records removed

Visualize the numeric variables

Univariate graphs

```
In [21]: #Box Plots for Numeric Variables
fig, axes = plt.subplots(1, 3, figsize=(20, 5))
fig.suptitle('Numeric Variables', fontsize=20)
```

```
# Price
sns.boxplot(ax=axes[0], x=df1['price'], palette='rocket')
axes[0].set_title("Price")

# Year
sns.boxplot(ax=axes[1], x=df1['year'], palette='magma')
axes[1].set_title("Year")

# Mileage
sns.boxplot(ax=axes[2], x=df1['mileage'], palette='viridis')
axes[2].set_title("Mileage")

Text(0.5, 1.0, 'Mileage')
```



```
In [144]: #Histogram for Numeric Variables
fig, axes = plt.subplots(1, 3, figsize=(20, 10))
fig.suptitle('Numeric Variables', fontsize=20)
```

```
# Price
sns.histplot(ax=axes[0], x=df1['price'], color='purple')
axes[0].set_title("Price")

# Year
sns.histplot(ax=axes[1], x=df1['year'], color='orange')
axes[1].set_title("Year")

# Mileage
sns.histplot(ax=axes[2], x=df1['mileage'], color='green')
axes[2].set_title("Mileage")

Text(0.5, 1.0, 'Mileage')
```



```
In [23]: #Density Plot for Numeric Variables
fig, axes = plt.subplots(1, 3, figsize=(20, 10))
fig.suptitle('Numeric Variables', fontsize=20)
```

```
# Price
sns.distplot(ax=axes[0], x=df1['price'], color='purple')
axes[0].set_title("Price")

# Year
sns.distplot(ax=axes[1], x=df1['year'], color='orange')
axes[1].set_title("Year")

# Mileage
sns.distplot(ax=axes[2], x=df1['mileage'], color='green')
axes[2].set_title("Mile
```


In [27]: fig, axes = plt.subplots(1, 3, figsize=(20, 5))
fig.suptitle('Visualizing Brand', fontsize=20)

```
#Brand and Price
sns.boxplot(ax=axes[0], x="brand", y="price", data=df1)
axes[0].set_title("Brand and Price")

#Brand and Year
sns.boxplot(ax=axes[1], x="brand", y="year", data=df1)
axes[1].set_title("Brand and Year")

#Brand and Mileage
sns.boxplot(ax=axes[2], x="brand", y="mileage", data=df1)
axes[2].set_title("Brand and Mileage")
```

Out [27]: Text(0.5, 1.0, 'Brand and Mileage')



Title Status

```
In [28]: fig, axes = plt.subplots(1, 3, figsize=(20, 5))
fig.suptitle('Visualizing title_status', fontsize=20)

#Title Status and Price
sns.boxplot(ax=axes[0], x="title_status", y="price", data=df1)
axes[0].set_title("title_status and Price")

#Title Status and Year
sns.boxplot(ax=axes[1], x="title_status", y="year", data=df1)
axes[1].set_title("title_status and Year")

#Title Status and Mileage
sns.boxplot(ax=axes[2], x="title_status", y="mileage", data=df1)
axes[2].set_title("title_status and Mileage")
```

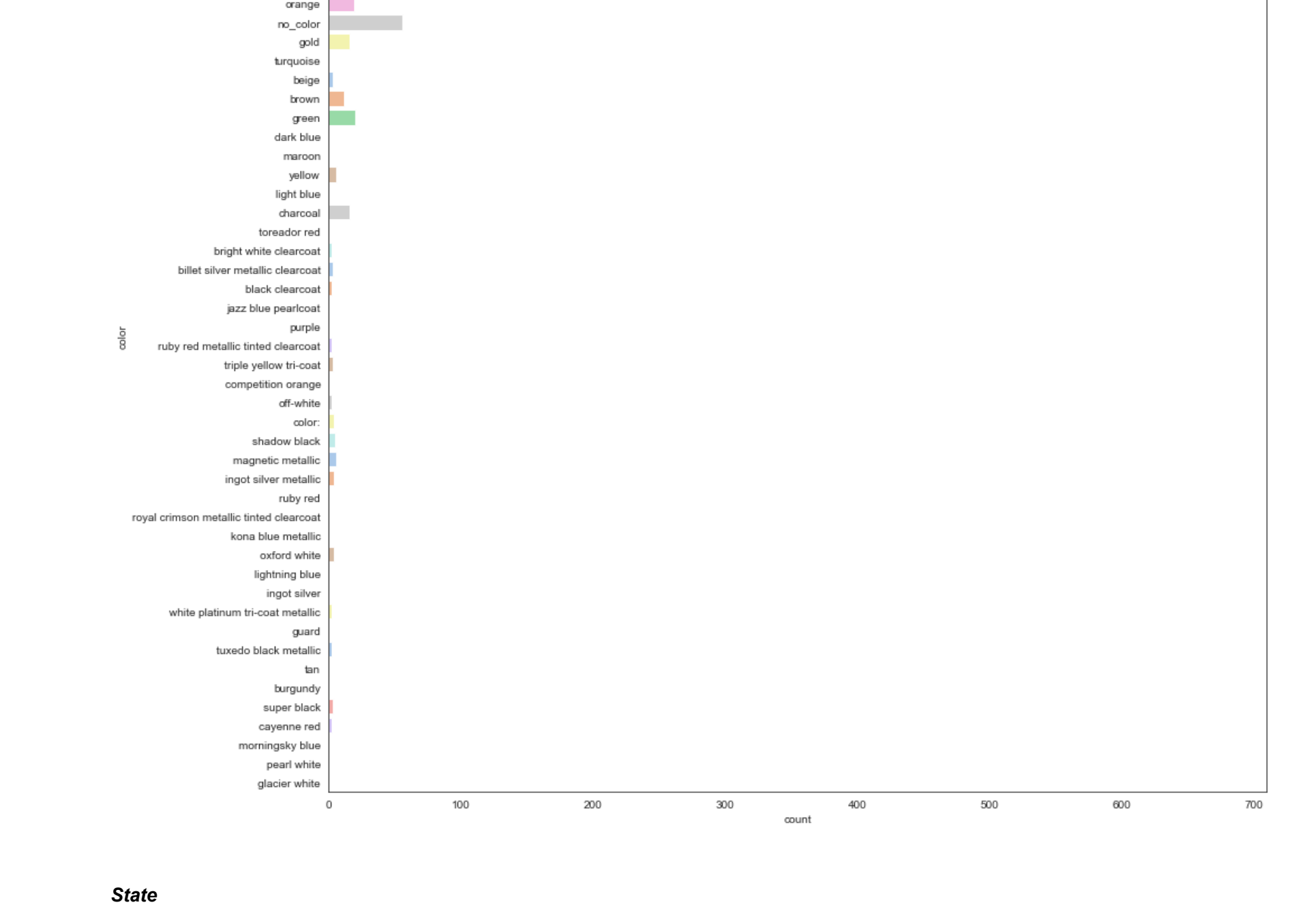
Out [28]: Text(0.5, 1.0, 'title_status and Mileage')



Model

```
In [29]: f, ax = plt.subplots(figsize=(20, 15))
sns.countplot(y="model", data=df1, palette="pastel")
```

Out [29]: <matplotlib.axes._subplots.AxesSubplot at 0x1ee795b0d0>



Color

```
In [30]: f, ax = plt.subplots(figsize=(15, 15))
sns.countplot(y="color", data=df1, palette="pastel")
```

Out [30]: <matplotlib.axes._subplots.AxesSubplot at 0x1ee7954f10>



State

```
In [31]: f, ax = plt.subplots(figsize=(15, 15))
sns.countplot(y="state", data=df1, palette="pastel")
```

Out [31]: <matplotlib.axes._subplots.AxesSubplot at 0x1ee7e1fa30>



Data Preparation

Separate the Categorical and Numeric Variables

```
In [32]: df1_categorical = df1[['brand', 'title_status', 'model', 'color', 'state']]
df1_categorical
```

Out [32]:

	brand	title_status	model	color	state
0	ford	clean vehicle	se	silver	tennessee
1	dodge	clean vehicle	mpv	silver	georgia
2	ford	clean vehicle	car	blue	virginia
3	chevrolet	clean vehicle	1500	red	texas
4	dodge	clean vehicle	mpv	white	florida
...
2309	nissan	clean vehicle	versa	red	california
2310	nissan	clean vehicle	versa	silver	florida
2311	nissan	clean vehicle	versa	silver	florida
2312	nissan	clean vehicle	versa	black	florida
2313	nissan	clean vehicle	versa	silver	florida
...
2314	nissan	clean vehicle	versa	red	california

2314 rows x 5 columns

```
In [33]: df1_numerical = df1[['price', 'year', 'mileage']]
df1_numerical
```

Out [33]:

	price	year	mileage
0	2899	10	190552.0
1	5350	3	30590.0
2	25000	7	64146.0
3	27700	3	6654.0
4	5700	3	45561.0
...
2309	7800	2	23609.0
2310	9200	3	34553.0
2311	9200	3	31594.0
2312	9200	3	32557.0
2313	9200	3	31371.0
...
2314	9200	3	31371.0

2314 rows x 3 columns

OneHot Encode the Categorical Variables

By default, the encoder derives the categories based on the unique values in each feature.

The method below one_hot encodes the categorical features to keep the prefixes for each feature when creating binary variables for categorical predictor levels.

```
In [34]: def onehot_encode(df, columns_with_prefixes):
df = df.copy()
for column, prefix in columns_with_prefixes:
dummies = pd.get_dummies(df[column], prefix=prefix)
df = pd.concat([df, dummies], axis=1)
df = df.drop(column, axis=1)
return df
```

In [35]: df1_categorical = onehot_encode(
df1_categorical,
columns_with_prefixes=[
('brand', 'br_'),
('model', 'md_'),
('color', 'cl_'),
('state', 'st_'),
('title_status', 'ts_')
])

In [36]: df1_categorical.describe(include="all")

Out [36]:

	br_chevrolet	br_dodge	br_ford	br_gmc	br_nissan	md_1500	md_2500	md_2500hd	md_3500	md_acadia	...	st_texas	st_ut
count	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	...	2314.000	2314.000
mean	0.127	0.187	0.534	0.018	0.135	0.017	0.003	4.322e-04	0.002	4.322e-04	...	-0.308	-0.062
std	0.333	0.390	0.499	0.134	0.342	0.129	0.059	2.079e-02	0.042	2.079e-02	...	-0.308	-0.062
min	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000e+00	0.000	0.000e+00	...	-0.308	-0.062
25%	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000e+00	0.000	0.000e+00	...	-0.308	-0.062
50%	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000e+00	0.000	0.000e+00	...	-0.308	-0.062
75%	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000e+00	0.000	0.000e+00	...	-0.308	-0.062
max	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000e+00	1.000	1.000e+00	...	-0.308	-0.062

8 rows x 188 columns

```
In [37]: df1_categorical
```

Out [37]:

	br_chevrolet	br_dodge	br_ford	br_gmc	br_nissan	md_1500	md_2500	md_2500hd	md_3500	md_acadia	...	st_texas	st_utah
0	0	0	1	0	0	0	0	0	0	0	...	0	0
1	0	1	0	0	0	0	0	0	0	0	...	0	0
2	0	0	1	0	0	0	0	0	0	0	...	0	0
3	1	0	0	0	0	1	0	0	0	0	...	0	0
4	0	1	0	0	0	0	0	0	0	0	...	1	0
...
2309	0	0	0	0	1	0	0	0	0	0	...	0	0
2310	0	0	0	0	1	0	0	0	0	0	...	0	0
2311	0	0	0	0	1	0	0	0	0	0	...	0	0
2312	0	0	0	0	1	0	0	0	0	0	...	0	0
2313	0	0	0	0	1	0	0	0	0	0	...	0	0
...
2314	0	0	0	0	1	0	0	0	0	0	...	0	0

2314 rows x 188 columns

Concatenate the numerical and categorical dataframe

```
In [38]: df_2 = pd.concat([df1_numerical, df1_categorical], axis=1)
df_2
```

Out [38]:

	price	year	mileage	br_chevrolet	br_dodge	br_ford	br_gmc	br_nissan	md_1500	md_2500	...	st_texas	st_utah	st_vermont
0	2899	10	190552.0	0	0	1	0	0	0	0	...	-0.308	-0.062	-0.029
1	5350	3	30590.0	0	1	0	0	0	0	0	...	-0.308	-0.062	-0.029
2	25000	7	64146.0	0	0	1	0	0	0	0	...	-0.308	-0.062	-0.029
3	27700	3	6654.0	1	0	0	0	0	1	0	...	-0.308	-0.062	-0.029
4	5700	3	45561.0	0	1	0	0	0	0	0	...	1	0	0
...
2309	7800	2	23609.0	0	0	0	0	1	0	0	...	-0.308	-0.062	-0.029
2310	9200	3	34553.0	0	0	0	0	1	0	0	...	-0.308	-0.062	-0.029
2311	9200	3	31594.0	0	0	0	0	1	0	0	...	-0.308	-0.062	-0.029
2312	9200	3	32557.0	0	0	0	0	1	0	0	...	-0.308	-0.062	-0.029
2313	9200	3	31371.0	0	0	0	0	1	0	0	...	-0.308	-0.062	-0.029
...
2314	9200	3	31371.0	0	0	0	0	1	0	0	...	-0.308	-0.062	-0.029

2314 rows x 191 columns

Scaling the input

Standardize features by removing the mean and scaling to unit variance

```
In [39]: scaler = StandardScaler()
df_2 = pd.DataFrame(scaler.fit_transform(df_2.drop(['price'], axis=1)), columns=df_2.drop(['price'], axis=1).columns)
df_2
```

Out [39]:

	year	mileage	br_chevrolet	br_dodge	br_ford	br_gmc	br_nissan	md_1500	md_2500	...	st_texas	st_utah	st_vt
count	1746	2769	-0.381	-0.479	0.935	-0.136	-0.395	-0.131	-0.059	-0.021	...	-0.308	-0.062
mean	0.127	0.187	0.534	0.018	0.135	0.017	0.003	4.322e-04	0.002	4.322e-04	...	-0.308	-0.062
std	0.333	0.390	0.499	0.134	0.342	0.129	0.059	2.079e-02	0.042	2.079e-02	...	-0.308	-0.062
min	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000e+00	0.000	0.000e+00	...	-0.308	-0.062
25%	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000e+00	0.000	0.000e+00	...	-0.308	-0.062
50%	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000e+00	0.000	0.000e+00	...	-0.308	-0.062
75%	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000e+00	0.000	0.000e+00	...	-0.308	-0.062
max	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000e+00	1.000	1.000e+00	...	-0.308	-0.062

8 rows x 188 columns

```
In [40]: df_2 = pd.concat([df1_numerical['price'], df_2], axis = 1)
df_2
```

Out [40]:

	price	year	mileage	br_chevrolet	br_dodge	br_ford	br_gmc	br_nissan	md_1500	md_2500	...	st_texas	st_utah	st_vermont
0	2899	1.746	2.769	-0.381	-0.479	0.935	-0.136	-0.395	-0.131	-0.059	...	-0.308	-0.062	-0.029
1	5350	-0.380	-0.225	-0.381	2.087	-1.070	-0.136	-0.395	-0.131	-0.059	...	-0.308	-0.062	-0.029
2	25000	0.844	0.882	-0.381	-0.479	0.935	-0.136	-0.395	-0.131	-0.059	...	-0.308	-0.062	-0.029
3	27700	-0.380	-0.078	2.626	-0.479	-1.070	-0.136	-0.395	7.638	-0.059	...	-0.308	-0.062	-0.029
4	5700	-0.380	-0.107	-0.381	2.087	-1.070	-0.136	-0.395	-0.131	-0.059	...	3.251	-0.062	-0.029
...
2309	7800	-0.681	-0.542	-0.381	-0.479	-1.070	-0.136	2.533	-0.131	-0.059	...	-0.308	-0.062	-0.029
2310	9200	-0.360	-0.325	-0.381	-0.479	-1.070	-0.136	2.533	-0.131	-0.059	...	-0.308	-0.062	-0.029
2311	9200	-0.360	-0.364	-0.381	-0.479	-1.070	-0.136	2.533	-0.131	-0.059	...	-0.308	-0.062	-0.029
2312	9200	-0.360	-0.365	-0.381	-0.479	-1.070	-0.136	2.533	-0.131	-0.059	...	-0.308	-0.062	-0.029
2313	9200	-0.360	-0.388	-0.381	-0.479	-1.070	-0.136	2.533	-0.131	-0.059	...	-0.308	-0.062	-0.029
...
2314	9200	-0.360	-0.388	-0.381	-0.479	-1.070	-0.136	2.533	-0.131	-0.059	...	-0.308	-0.062	-0.029

2314 rows x 191 columns

df_2.describe()

Out [41]:

	price	year	mileage	br_chevrolet	br_dodge	br_ford	br_gmc	br_nissan	md_1500	md_2500	...	st_texas
count	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	2314.000	...	2314.000
mean	19058.407	6.939e+16	2.953e-17	1.915e-15	-0.175e-15	4.535e-16	-2.942e-15	3.702e-15	-2.665e-15	-8.020e-15	...	-0.308
std	11944.640	1.000e+00	1.000e+00	1.000e+00	1.000e+00	1.000e+00	1.000e+00	1.000e+00	1.000e+00	1.000e+00	...	1.000e+00
min	0.000	-0.617e-01	-1.010e+00	-3.808e-01	-4.791e-01	-1.070e+00	-1.360e-01	-3.948e-01	-1.309e-01	-5.890e-02	...	-0.307e-01
25%	10607.500	-6.608e-01	-5.850e-01	-3.808e-01	-4.791e-01	-1.070e+00	-1.360e-01	-3.948e-01	-1.309e-01	-5.890e-02	...	-0.307e-01
50%	17175.000	-3.599e-01	-3.089e-01	-3.808e-01	-4.791e-01	9.347e-01	-1.360e-01	-3.948e-01	-1.309e-01	-5.890e-02	...	-0.307e-01
75%	25900.000	2.419e-01	2.344e-01	-3.808e-01	-4.791e-01	9.347e-01	-1.360e-01	-3.948e-01	-1.309e-01	-5.890e-02	...	-0.307e-01
max	74000.000	2.138e+01	1.882e+01	2.626e+00	2.087e+00	9.347e-01	7.355e+00	2.533e+00	7.638e+00	1.698e+01	...	3.251e+00

8 rows x 191 columns