

Assignment – Rodolfo Lerma

The first part of this analysis (as seen in the appendix) was to explore the data to get a better idea of what kind of data that I have available for this analysis and identify where some cleaning and/or formatting might be needed and get familiar with the dataset in general. After this section and the identification of the areas that might need improvement, I started working on cleaning and formatting the data. During this step I decided to normalize the numerical columns, change the format from the Boolean columns to numeric (0 & 1), drop missing values as they were around 15% of the training dataset (another alternative if could have been to replace the missing values for the median in case of numerical columns and the mode for categorical ones, which is what I did for the test data set as the missing values were around 5%), get one hot encodings for some of the categorical columns and to extract as much information as possible from the 2 columns that were formatted as a list of lists (VehFeats & VehHistory).

Once the data was clean and with the needed formatting, I decided to work on 2 different models: A **regression** for the prediction of the Dealer Listing Price and a **classifier** for the Vehicle Trim.

The first part of the **regression** model was to do a feature selection for the features that would be part of the model, for this section I went for an embedded method (Lasso) as it provides a way to minimize to zero those coefficients that contribute the least to the regression. For the Regression model itself I used the ElasticNet model from the sklearn – linear model package. I split the Training dataset into “Training” & “Validation” and selected multiple hyperparameters to select the best combination of values which I check by means of cross validation and looking at the r squared (it is important to notice that other measures of performance could have been used, such as mean squared error. This relates to the end use of the model). After this process the training and validation datasets were combined to train the selected model again and use this last model as the final one for the test data.

For the second model (**classifier**) I followed the same structure as before: Feature selection (used Lasso as well), did a comparison of 3 models looking for the best hyperparameters for each of the classifiers (3 legacy machine learning models) using the “RandomizedSearchCV” for Cross Validation. After looking at the performance of the validation dataset and verifying that no overfitting was present, I re-trained the model using the entire “Training dataset”.

I did not use deep learning because the dataset was too small for the noted application, and I believe interpretability was important for this model, but maybe an option to explore for the future. (A basic Deep Learning Section was added at the end of the analysis).

BOEING ASSIGNMENT

Applicant: Rodolfo Lerma

Problem Statement:

Include all the training dataset is information on used cars previously sold. Each row corresponds to one used car listing. The first column of the data contains a unique identifier for the listing. The next twenty-six columns contain information on parameters relevant to the transaction, with those parameters described in more detail in the appendix attached. Finally, the last two columns of the 'Training_dataset.csv' contain information on 'Vehicle_Trim' and 'Dealer_Listing_Price', which describe the trim of the vehicle involved in the sale, and the price at which the vehicle was listed by the dealer. Your challenge is to build one or more models, through whatever means you find most appropriate, capable of predicting **vehicle trim** and **dealer listing price** given the other twenty-six variables provided.

ABSTRACT:

The analysis is structured in the following way:

Analysis - Prediction (Price & Trim)

- **Data Exploration**
 - Training Dataset
 - Testing Dataset
 - General Exploration
 - Distribution of Target Variables
 - Distribution of Numerical Variables
 - Categorical Variables
 - Unique Values Training Dataset
- **Data Formatting (Training Dataset & Testing Dataset)**
 - Missing Values
 - Removing Columns from Dataset
 - Formatting Boolean Variables
 - Normalizing Numerical Variables
 - One Hot Encoding For Categorical Columns
 - Formatting of two List of Lists Columns (VehHistory & VehFeats)
 - VehFeats
 - VehHistory
 - Splitting the dataset into Target Variables and Features (Training Dataset)
- **Model 1: Prediction of Dealer Listing Price**
 - Feature Selection (LASSO) for Dealer Listing Price Variable
 - Model for Dealer Listing Price prediction
 - Training
 - Results/Predictions (Training & Validation dataset)
 - Re-train the model but now with all the datapoints (training + validation)
 - Prediction using the Test_data
- **Model 2: Prediction of Vehicle Trim**
 - Feature Selection (LASSO) for Vehicle Trim Variable
 - Different Models and Hyperparameters
 - Hyperparameter Settings for the Decision Tree: Model
 - Hyperparameter Settings for the Random Forest: Model
 - Hyperparameter Settings for the Logistic Regression: Model
 - Comparison & Final Model Selection
 - Final Model 2 Selection and Training (using training + validation data)
- **Export to Dataframe**

Analysis - Prediction (Price & Trim)

Data Exploration

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

In [2]: Test_df = pd.read_csv('Test_Dataset.csv')
Training_df = pd.read_csv('Training_Dataset.csv')
```

Training Dataset

```
In [3]: Training_df.head(5)
```

	ListingID	SellerCity	SellerPriv	SellerListSrc	SellerName	SellerRating	SellerRevCnt	SellerState	SellerZip	VehBodystyle	...	VehMake	V
0	3287	Warren	False	Inventory Command Center	Prime Motorz	5.0	32	MI	48091.0	SUV	...	Jeep	
1	3920	Fargo	False	Cadillac Certified Program	Gateway Chevrolet	4.8	1456	ND	58103.0	SUV	...	Cadillac	
2	4777	Waukeha	False	Jeep Certified Program	Chrysler Chrysler Ram Ram	4.8	1405	WI	53186.0	SUV	...	Jeep	
3	6242	Wentzville	False	Inventory Command Center	Century Dodge Jeep RAM	4.4	21	MO	63385.0	SUV	...	Jeep	
4	7108	Fayetteville	False	HomeNet Automotive	Superior Buick GMC of Fayetteville	3.7	74	AR	72703.0	SUV	...	Cadillac	

```
In [4]: Training_df.shape
```

Out[4]: (6298, 29)

It is possible to see that the Training Data Set contains 29 columns and 6298 data points. Out of those 29 columns, 2 are the 'target' variables.

Testing Dataset

```
In [5]: Test_df.head(5)
```

	ListingID	SellerCity	SellerPriv	SellerListSrc	SellerName	SellerRating	SellerRevCnt	SellerState	SellerZip	VehBodystyle	...	VehHistory	O
0	8622015	Seneca	False	HomeNet Automotive	Lake Keywee Chrysler Dodge Jeep Ram	2.5	59	SC	29678	SUV	...	1 Owner, Non-Personal Use Reported, Buyback Pr...	
1	8625693	Bedford	False	Inventory Command Center	North Coast Auto Mall	4.7	2116	OH	44146	SUV	...	1 Owner, Accident(s) Reported, Non-Personal Us...	
2	8625750	Webster	False	Jeep Certified Program	Marina Chrysler Dodge Jeep Mitsubishi RAM	3.9	46	NY	14580	SUV	...	1 Owner, Buyback Protection Eligible	
3	8626885	Louisville	False	Digital Motorworks (DM)	Oxmoor Ford Lincoln	4.5	1075	KY	40222	SUV	...	1 Owner, Buyback Protection Eligible	
4	8627430	Palmyra	False	Digital Motorworks (DM)	F.C. Kerbeck & Sons	4.6	162	NJ	8065	SUV	...	1 Owner, Non-Personal Use Reported, Buyback Pr...	

```
In [6]: Test_df.shape
```

Out[6]: (1000, 27)

It is possible to see that the Training Dataset contains 27 columns and 1000 data points, this since the 2 'target' variables were removed from the dataset.

General Exploration

```
In [7]: Training_df.dtypes
```

Out[7]:

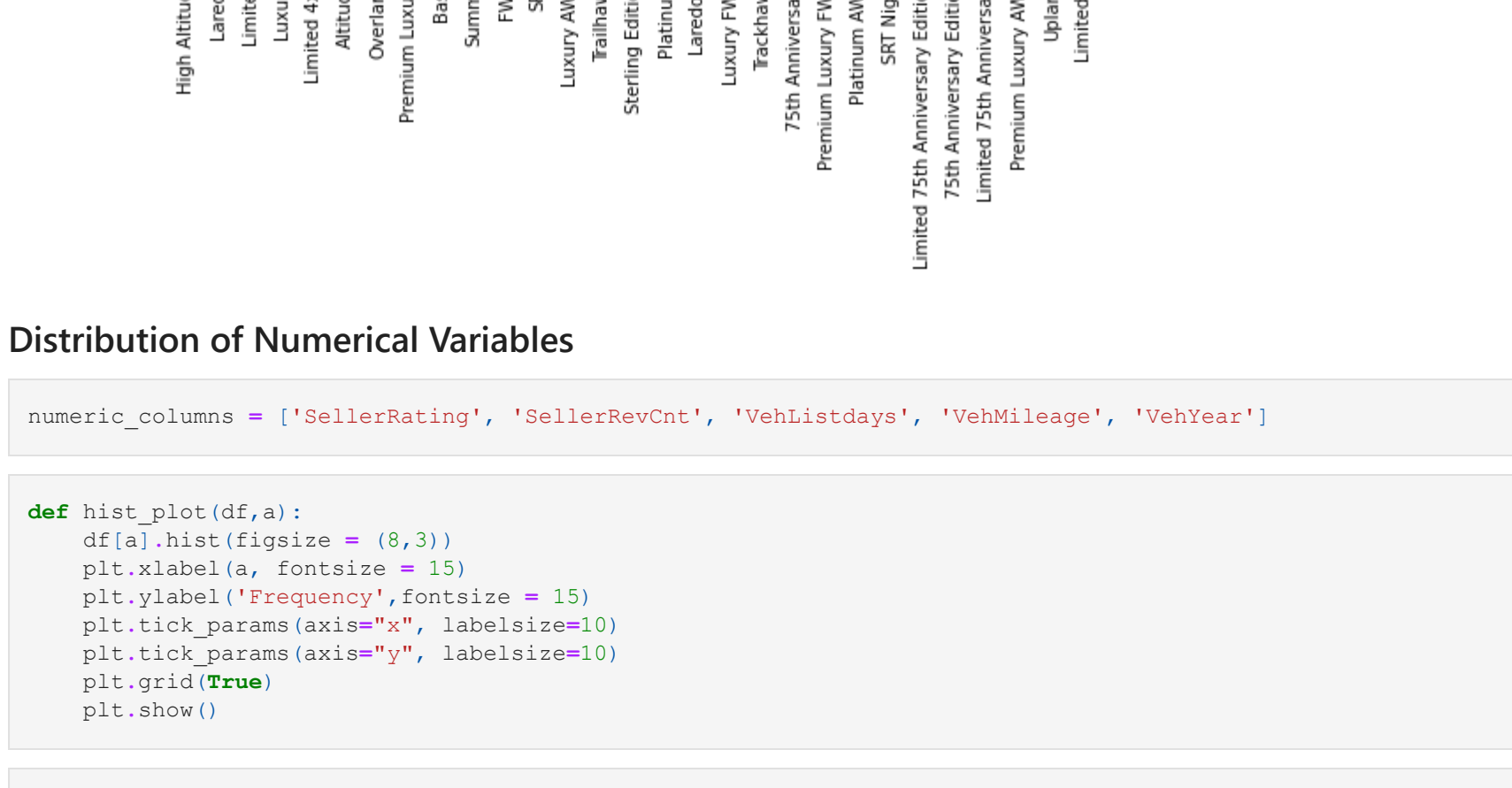
```
ListingID      int64
SellerCity     object
SellerPriv     bool
SellerListSrc  object
SellerName     object
SellerRating   float64
SellerRevCnt   int64
SellerState    object
SellerZip      float64
VehBodystyle   object
VehColorExt    object
VehColorInt    object
VehDriveTrain  object
VehEngine      object
VehFeats       object
VehHistory     object
VehType        object
VehYear        float64
VehMake        object
VehMileage     float64
VehModel       object
VehPriceLabel  object
VehSellerNotes object
VehTransmission object
VehType        int64
Vehicle_Trim   object
Dealer_Listing_Price float64
dtype: object
```

Looking at the data, the type distribution is as follows:

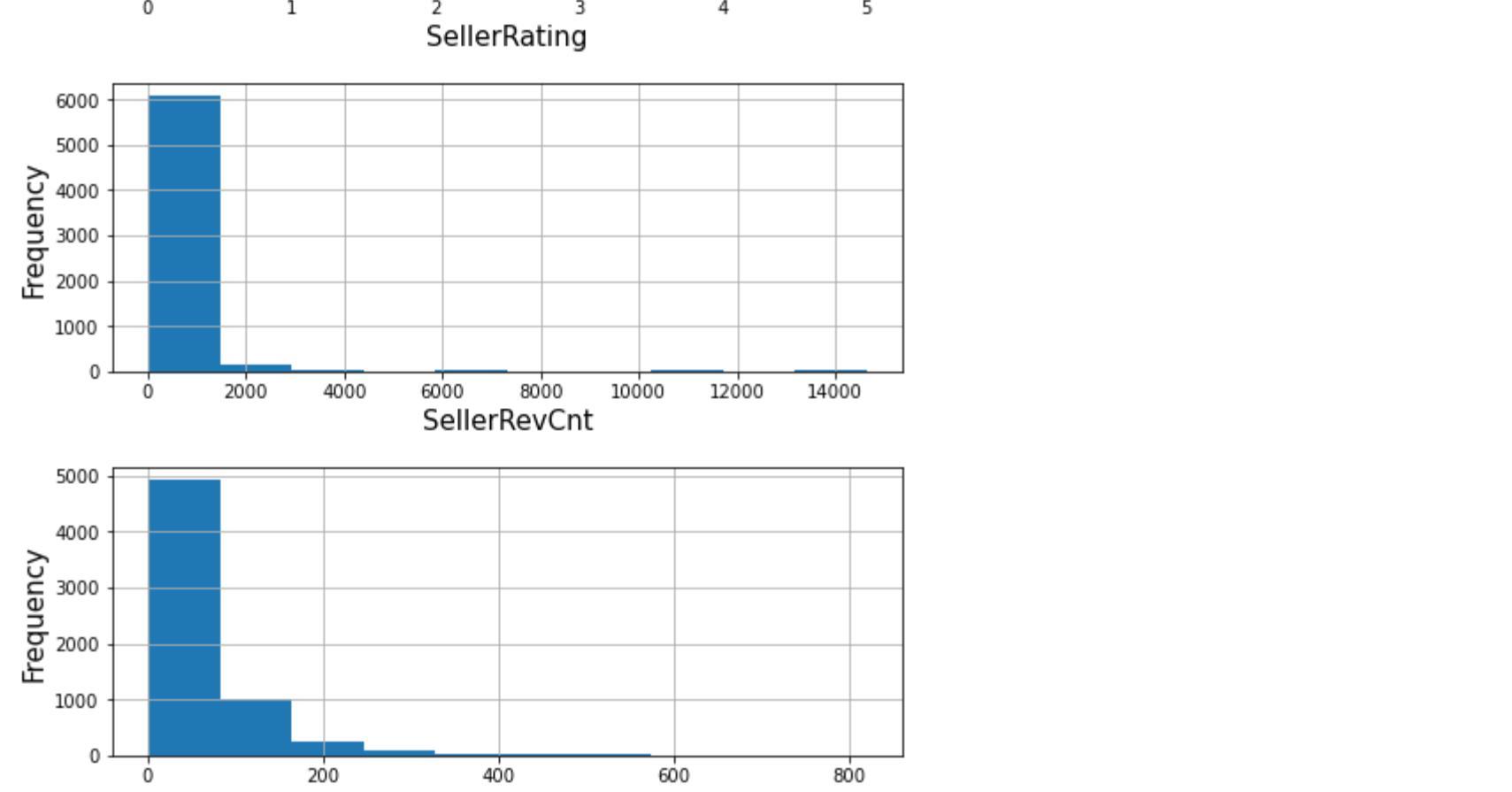
- 8 Numerical Variables (6 float64 & 2 int64)
- 2 Boolean Variables (bool)
- 19 Categorical Variables (object)

Distribution of the Target Variable (Price & Vehicle Trim)

```
In [8]: Training_df['Dealer_Listing_Price'].hist(figsize = (10,4))
plt.xlabel('Dealer Price', fontsize = 20)
plt.ylabel('Frequency', fontsize = 20)
plt.title('Dealer Listing Price', fontsize = 20)
plt.tick_params(axis='x', labelsize=10, rotation=90)
plt.tick_params(axis='y', labelsize=15)
plt.grid(True)
plt.show()
```



```
In [9]: Training_df['Vehicle_Trim'].hist(figsize = (10,4))
plt.xlabel('Frequency', fontsize = 20)
plt.ylabel('Frequency', fontsize = 20)
plt.title('Vehicle Trim', fontsize = 20)
plt.tick_params(axis='x', labelsize=10, rotation=90)
plt.tick_params(axis='y', labelsize=15)
plt.grid(True)
plt.show()
```

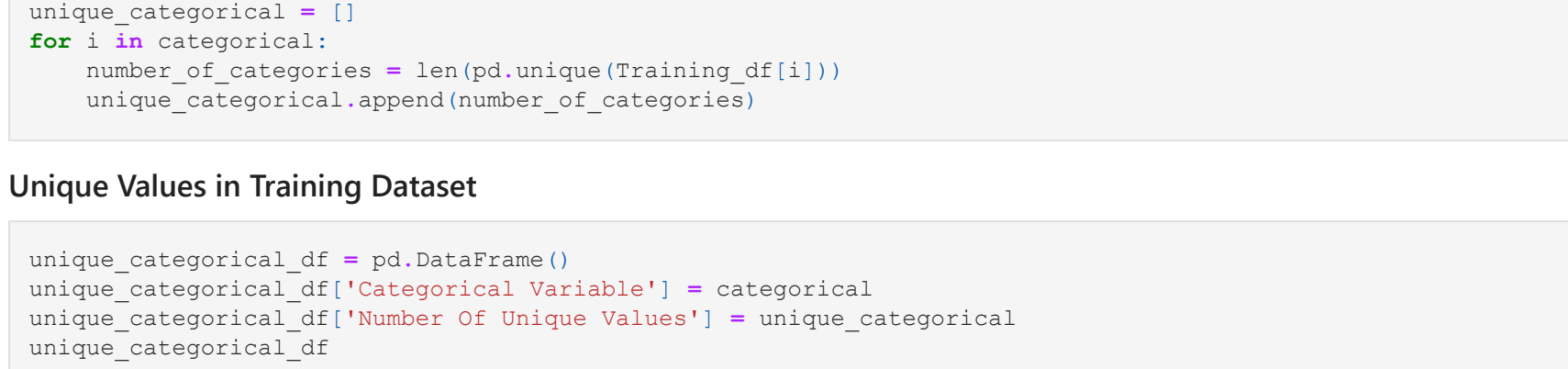
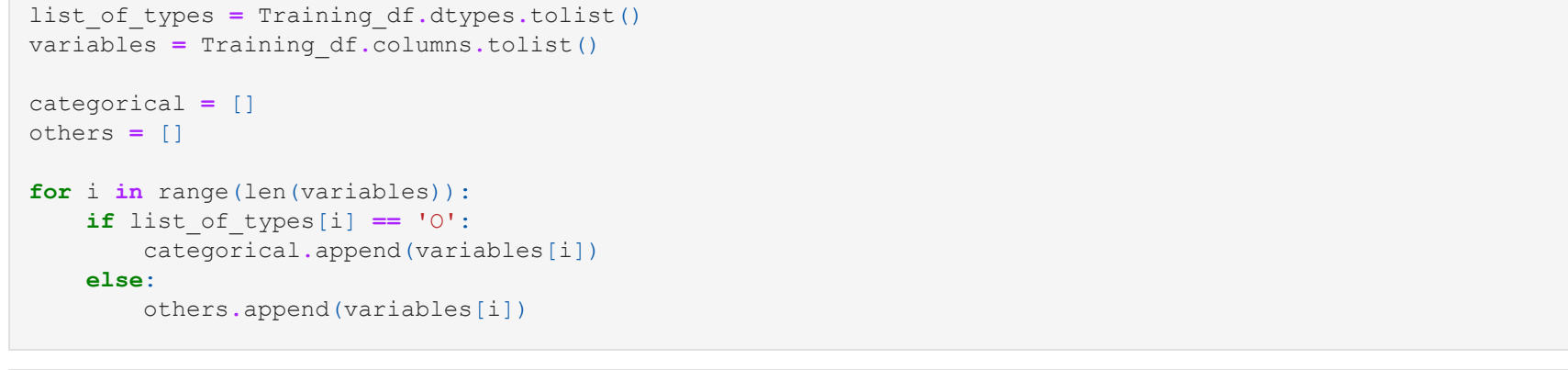
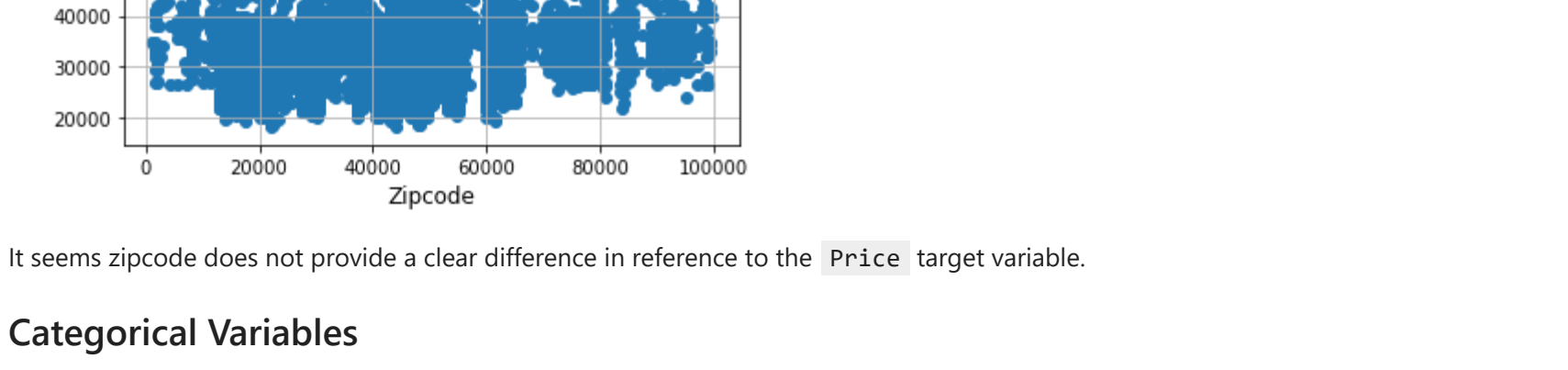


Distribution of Numerical Variables

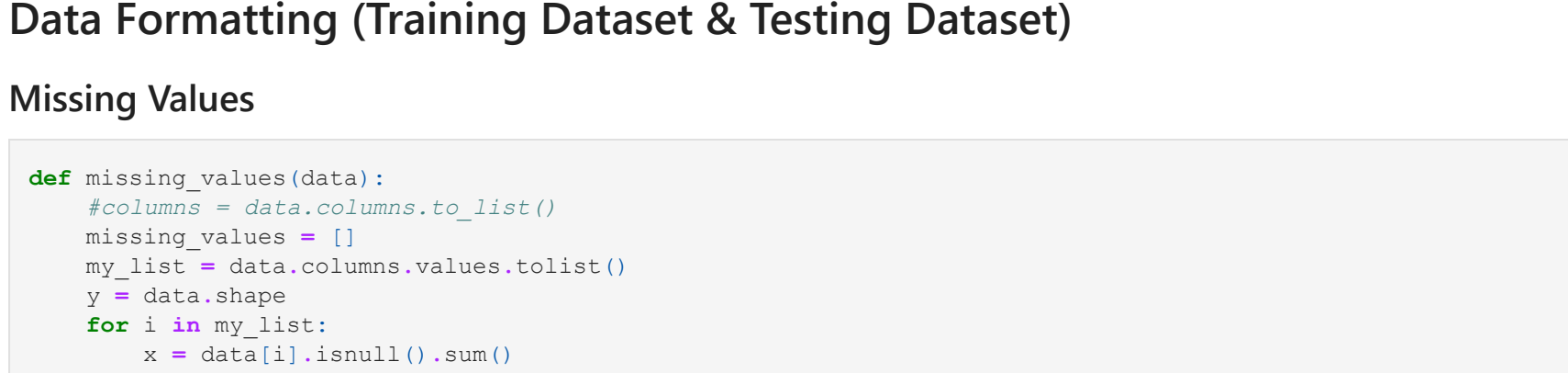
```
In [10]: numeric_columns = ['SellerRating', 'SellerRevCnt', 'VehListdays', 'VehMileage', 'VehYear']
```

```
In [11]: df = Training_df[numeric_columns]
df.hist(figsize = (8,3))
plt.xlabel('Frequency', fontsize = 15)
plt.ylabel('Frequency', fontsize = 15)
plt.tick_params(axis='x', labelsize=10)
plt.tick_params(axis='y', labelsize=10)
plt.grid(True)
plt.show()
```

```
In [12]: for i in numeric_columns:
    hist_plot(Training_df, i)
```



```
In [13]: plt.scatter(Training_df['SellerZip'], Training_df['Dealer_Listing_Price'])
plt.title('Price Vs Zipcode', fontsize = 15)
plt.xlabel('Zipcode', fontsize = 12)
plt.ylabel('Price', fontsize = 12)
plt.grid()
plt.show()
```



It seems zipcode does not provide a clear difference in reference to the 'Price' target variable.

Categorical Variables

```
In [14]: list_of_types = Training_df.dtypes.tolist()
variables = Training_df.columns.tolist()

categorical = []
others = []

for i in range(len(variables)):
    if list_of_types[i] == 'object':
        categorical.append(variables[i])
    else:
        others.append(variables[i])
```

```
In [15]: unique_categorical = []
for i in categorical:
    unique_of_categories = len(pd.unique(Training_df[i]))
    unique_categorical.append(unique_of_categories)
```

Unique Values in Training Dataset

```
In [16]: unique_categorical_df = pd.DataFrame()
unique_categorical_df['Categorical Variable'] = categorical
unique_categorical_df['Number Of Unique Values'] = unique_categorical
unique_categorical_df
```

	Categorical Variable	Number Of Unique Values
0	SellerCity	1318
1	SellerListSrc	9
2	SellerName	2452
3	SellerState	50
4	VehBodystyle	1
5	VehColorExt	170
6	VehColorInt	107
7	VehDriveTrain	21
8	VehEngine	97
9	VehFeats	844
10	VehFuel	5
11	VehHistory	34
12	VehMake	2
13	VehModel	2
14	VehPriceLabel	4
15	VehSellerNotes	4921
16	VehType	1
17	VehTransmission	34
18	Vehicle_Trim	30

Here we can see that some categories/columns have values that are very unique as 'VehSellerNotes' (which makes sense as these are notes from the seller) and other categories have very unique as 'VehType' & 'VehBodyStyle' with just one value.

Data Formatting (Training Dataset & Testing Dataset)

Missing Values

```
In [17]: def missing_values(data):
    #columns = data.columns.tolist()
    missing_values = []
    my_list = data.columns.values.tolist()
    for i in my_list:
        x = data[i].isnull().sum()
        columns_missing.append(x)
        missing_values.append(x)
    a = sorted(columns_missing.items(), key=lambda x: x[1], reverse = True)
    missing_columns_names = []
    for key, value in columns_missing.items():
        if value > 0:
            x = key
            missing_columns_names.append(x)
    return columns_missing, missing_columns_names
```

```
In [18]: #Training Data
dic_missing, missing_columns_names = missing_values(Training_df)

#Testing Data
dic_missing_test, missing_columns_names_test = missing_values(Test_df)
```

```
In [19]: missing_training_df = pd.DataFrame()
missing_training_df['Variable'] = Training_df.columns.values.tolist()
missing_training_df['Number Of Missing Variables'] = dic_missing_test_values()
missing_training_df
```

	Variable	Number Of Missing Variables
0	ListingID	0
1	SellerCity	0
2	SellerPriv	0
3	SellerListSrc	2
4	SellerName	0
5	SellerRating	0
6	SellerRevCnt	0
7	SellerState	0
8	SellerZip	2
9	VehBodystyle	0
10	VehCertified	0
11	VehColorExt	73
12	VehColorInt	728
13	VehDriveTrain	401
14	VehEngine	361
15	VehFeats	275
16	VehFuel	2
17	VehHistory	201
18	VehListdays	2
19	VehMake	0
20	VehMileage	2
21	VehModel	0
22	VehPriceLabel	38
23	VehSellerNotes	243
24	VehType	0
25	VehTransmission	197
26	VehYear	0
27	Vehicle_Trim	405
28	Dealer_Listing_Price	52

Looking at the Train Dataset we can see that there are some columns with more than 10% (VehColorInt) of missing data points.

```
In [20]: missing_test_df = pd.DataFrame()
missing_test_df['Variable'] = Test_df.columns.values.tolist()
missing_test_df['Number Of Missing Variables'] = dic_missing_test_values()
missing_test_df
```

	Variable	Number Of Missing Variables
0	ListingID	0
1	SellerCity	0
2	SellerPriv	0
3	SellerListSrc	0
4	SellerName	0
5	SellerRating	0
6	SellerRevCnt	0
7	SellerState	0
8	SellerZip	0
9	VehBodystyle	0
10	VehCertified	0
11	VehColorExt	7
12	VehColorInt	108
13	VehDriveTrain	64
14	VehEngine	58
15	VehFeats	37
16	VehFuel	0
17	VehHistory	27
18	VehListdays	0
19	VehMake	0
20	VehMileage	1
21	VehModel	0
22	VehPriceLabel	38
23	VehSellerNotes	41
24	VehType	0
25	VehTransmission	27
26	VehYear	0

```
In [21]: # For the training data the missing values will be removed
data = Training_df.copy()
data.dropna(inplace = True)
data.shape
data = data.reset_index()
```

If the missing data is dropped we will reduce the dataset by around 20%. Depending on the example and the size of the dataset other options might be replacing the values with the mean/median (if numerical) or mode (if categorical).

```
In [22]: # For the testing data the missing values will be filled with the mode and median
test_data = test_df.copy()
categorical_test_missing = ['VehColorExt', 'VehColorInt', 'VehDriveTrain', 'VehEngine', 'VehFeats', 'VehHistory', 'VehSellerNotes', 'VehTransmission']
numerical_test_missing = ['VehMileage']

# Replace missing values by median and mode
import statistics as st

for i in numerical_test_missing:
    test_data[i].fillna(test_data[i].median(), inplace=True)

for i in categorical_test_missing:
    test_data[i].fillna(st.mode(test_data[i]).mode[0], inplace=True)
```

Removing Columns from dataset

Based on the nature of the data some columns will be removed as they do not add to the prediction:

- ListingID (as is a unique value per entry)
- VehSellerNotes (need more processing to be used)
- SellerZip
- VehType (as it has just one value)
- VehBodystyle (as it has just one value)

```
In [23]: # For Training Dataset
data.drop(columns=['ListingID', 'VehSellerNotes', 'SellerZip', 'VehType', 'VehBodystyle', 'index'], inplace = True)
data_updated = data.columns.tolist()
```

```
In [24]: # For Testing Dataset
test_data.drop(columns=['VehSellerNotes', 'SellerZip', 'VehType', 'VehBodystyle'], inplace = True)
test_data_updated = test_data.columns.tolist()
```

Formatting Boolean Variables

```
In [25]: boolean_columns = ['SellerPriv', 'VehCertified']
def boolean_format(columns, df):
    for i in columns:
        df[i] = df[i].astype(int)
```

```
In [26]: #Training Dataset
boolean_format(boolean_columns, data)

#Test Dataset
boolean_format(boolean_columns, test_data)
```

Normalizing Numerical Variables

```
In [27]: def normalization_format(numeric_columns, df):
    from sklearn.preprocessing import StandardScaler
    scaler = StandardScaler()
    scaler.fit(df[numeric_columns])
    data_scaled = scaler.transform(df[numeric_columns])
    return df

In [28]: # Training Dataset
data_norm = normalization_format(numeric_columns, data)
data.drop(columns = numeric_columns, inplace = True)
data = data.join(data_norm)

# Test Dataset
test_data_norm = normalization_format(numeric_columns, test_data)
test_data.drop(columns = numeric_columns, inplace = True)
test_data = test_data.join(data_norm_test)
```

```
In [29]: # Special Columns
special_columns = ['VehHistory', 'VehFeats']
targets = ['Vehicle_Trim', 'Dealer_Listing_Price']

In [30]: # Categorical Variables:
data['Vehicle_Trim'] = data['Vehicle_Trim'].astype('category')
d = dict(enumerate(data['Vehicle_Trim'].cat.categories))
data['Vehicle_Trim_cat'] = data['Vehicle_Trim'].cat.codes
data.drop(columns = 'Vehicle_Trim', inplace = True)
```

```
In [31]: categorical_columns = [i for i in columns_updated if i not in numeric_columns]
categorical_columns = [i for i in categorical_columns if i not in boolean_columns]
categorical_columns = [i for i in categorical_columns if i not in special_columns]

In [32]: for i in categorical_columns:
    data[i] = data[i].str.lower() #Training
    test_data[i] = test_data[i].str.lower() #Test
```

One Hot Encoding For Categorical Columns

```
In [33]: def hot_encoding(data, categorical_columns):
    one_hot = pd.get_dummies(data[categorical_columns])
    return one_hot

In [34]: # Hot encoding for Training
one_hot = hot_encoding(data, categorical_columns)
data = data.join(one_hot)
data.drop(columns = categorical_columns, inplace = True)

In [35]: # Hot encoding for Test
one_hot_test = hot_encoding(test_data, categorical_columns)
test_data = test_data.join(one_hot_test)
test_data.drop(columns = categorical_columns, inplace = True)
```

Formatting of two List of Lists Columns (VehHistory & VehFeats)

There are 2 columns with a list of values that might be important for the analysis. For this reason they are going to be split and included as features (One Hot Encoding):

- VehHistory
- VehFeats

See examples below:

```
In [36]: data['VehHistory'].head()
```

Out[36]:

```
0      1 Owner, Buyback Protection Eligible
1      1 Owner, Non-Personal Use Reported, Buyback Pr...
2      1 Owner, Non-Personal Use Reported, Buyback Pr...
3      1 Owner, Non-Personal Use Reported, Buyback Pr...
4      1 Owner, Accident(s) Reported, Non-Personal Use...
Name: VehHistory, dtype: object
```

```
In [37]: data['VehFeats'].head()
```

Out[37]:

```
0      ['18 MHRRL amp; 8.4 RADIO GROUP-inc; Nav-Cap...
1      ['Android Auto', 'Antilock Brakes', 'Apple Car...
2      ['4-wheel Disc Brakes', 'ABS', 'Adjustable Ste...
3      ['1st and 2nd row curtain head airbags', '4-wh...
4      ['1st and 2nd row curtain head airbags', '4-wh...
Name: VehFeats, dtype: object
```

VehFeats

```
In [38]: def veh_features(data):
    x = pd.Series(data['VehFeats'])
    veh_list_df = data['VehFeats'].tolist()
    mySeries = pd.Series(split_df)
    letter_list = x.split(",")

    list_of_lists = []
    for i in range(len(mySeries)):
        x = mySeries[i][1:-1]
        new_string = x.replace(" ", "")
        new_string = new_string.replace("(", "")
        new_string = new_string.replace(")", "")
        new_string = new_string.replace("'", "")
        new_string = new_string.lower()
        letter_list.append(letter_list[i])

    from sklearn.preprocessing import MultiLabelBinarizer

    test = pd.Series(list_of_lists)
    mlb = MultiLabelBinarizer()
    onehot_VehFeats = pd.DataFrame(mlb.fit_transform(test),
                                   columns=mlb.classes_,
                                   index=test.index)

    return onehot_VehFeats
```

```
In [39]: # Training Dataset
onehot_VehFeats = veh_features(data)
data = data.join(onehot_VehFeats)
data.drop(columns = 'VehFeats', inplace = True)
```

```
# Test Dataset
onehot_VehFeats_test = veh_features(test_data)
test_data = test_data.join(onehot_VehFeats_test)
test_data.drop(columns = 'VehFeats', inplace = True)
```

VehHistory

```
In [40]: def veh_history(data):
    veh_list_df2 = data['VehHistory'].tolist()
    mySeries2 = pd.Series(split_df2)

    list_of_lists2 = []
    for i in range(len(mySeries2)):
        x = mySeries2[i]
        letter_list2 = x.split(",")
        list_of_lists2.append(letter_list2[i])

    from sklearn.preprocessing import MultiLabelBinarizer

    test2 = pd.Series(list_of_lists2)
    mlb2 = MultiLabelBinarizer()
    onehot_VehHistory = pd.DataFrame(mlb2.fit_transform(test2),
                                    columns=mlb2.classes_,
                                    index=test2.index)

    return onehot_VehHistory
```

```
In [41]: # Training Dataset
onehot_VehHistory = veh_history(data)
data = data.join(onehot_VehHistory)
data.drop(columns = 'VehHistory', inplace = True)
```

Splitting the dataset into Target Variables and Features (Training Dataset)

