

Project 2: Avocado Price Prediction using Machine Learning

Autor: Barbara Jean

Date: 07/21/2024

output: pdf_document

```
In [36]: # handling warnings as errors
import warnings
warnings.filterwarnings("ignore")

In [2]: # Importing the necessary libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import statistics
import scipy.stats as st
from sklearn.model_selection import train_test_split
from sklearn.feature_selection import SelectKBest
from sklearn.metrics import r2_score, mean_squared_error
from sklearn.metrics import mean_squared_error
from math import sqrt
from sklearn.metrics import accuracy_score, confusion_matrix
from sklearn.metrics import confusion_matrix
from sklearn.ensemble import RandomForestRegressor
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from sklearn.preprocessing import OrdinalEncoder
from sklearn import preprocessing
from sklearn.metrics import accuracy_score
import joblib as job
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report
from sklearn.tree import DecisionTreeClassifier
from sklearn.tree import DecisionTreeRegressor
from pandas.plotting import autocorrelation_plot
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.tsa.stattools import adfuller
from pandas.plotting import autocorrelation_plot
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
import statsmodels.api as sm
from matplotlib import pyplot
from pandas.plotting import lag_plot
from statsmodels.tsa.ar_model import AutoReg
from pandas import DataFrame

In [3]: # Reading dataset
dataprd_read_csv("C:/Users/78bar/dsc_688/avocado.csv")
avocado=pd.DataFrame(data)
print("The loading of the dataset was successful.\n")

The loading of the dataset was successful.

In [4]: avocado.head() # Reading the first records by using the head() metho

Out[4]:
```

Unnamed: 0	Date	Average Price	TotalVolume	Plu4046	Plu4225	Plu4770	Total Bags	Small Bags	Large Bags	XLarge Bags	type	year	region
0	0 12/27/2015	1.33	64236.62	1036.74	54545.85	48.16	8696.87	8603.62	93.25	0.0	conventional	2015	Albany
1	1 12/20/2015	1.35	54876.98	674.28	44638.81	58.33	9505.56	9408.07	97.49	0.0	conventional	2015	Albany
2	2 12/13/2015	0.93	118220.22	794.70	109149.67	130.50	8145.35	8042.21	103.14	0.0	conventional	2015	Albany
3	3 12/6/2015	1.08	78992.15	1132.00	71976.41	72.58	5811.16	5677.40	133.76	0.0	conventional	2015	Albany
4	4 11/29/2015	1.28	51039.60	941.48	43838.39	75.78	6183.95	5986.26	197.69	0.0	conventional	2015	Albany

Data exploration

```
In [5]: avocado.shape # Returning a tuple representing the dimensionality of the DataFrame

Out[5]: (18249, 14)

In [6]: avocado.size # Returning an int representing the number of elements in this object

Out[6]: 255486

In [7]: avocado.info() # Printing a summary of the dataframe, index dtype and columns, non-null values

Out[7]:
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18249 entries, 0 to 18248
Data columns (total 14 columns):
 #   Column        Non-Null Count  Dtype
---  --
 0   Unnamed: 0    18249 non-null  int64
 1   Date          18249 non-null  object
 2   Average Price 18249 non-null  float64
 3   TotalVolume   18249 non-null  float64
 4   Plu4046       18249 non-null  float64
 5   Plu4225       18249 non-null  float64
 6   Plu4770       18249 non-null  float64
 7   Total Bags    18249 non-null  float64
 8   Small Bags    18249 non-null  float64
 9   Large Bags    18249 non-null  float64
10  XLarge Bags   18249 non-null  float64
11  type          18249 non-null  object
12  year          18249 non-null  int64
13  region        18249 non-null  object
dtypes: float64(9), int64(2), object(3)
memory usage: 1.9+ MB

In [8]: avocado.dtypes.unique()

Out[8]: array([dtype('int64'), dtype('O'), dtype('float64')], dtype=object)

In [9]: avocado.describe() # Looking at the statistical summary of the variables with describe()

Out[9]:
```

	Unnamed: 0	Average Price	TotalVolume	Plu4046	Plu4225	Plu4770	Total Bags	Small Bags	Large Bags	XLarge Bags	year
count	18249.000000	18248.000000	1.824900e+04	1.824900e+04	1.824900e+04	1.824900e+04	1.824900e+04	1.824900e+04	1.824900e+04	1.824900e+04	18249.000000
mean	2423232	1.409978	3.456450e+01	2.930064e+05	2.951546e+05	2.283974e+04	2.396392e+05	1.821947e+05	5.433960e+04	3105.428507	2016.147899
std	15.401405	0.402677	3.456450e+01	1.264980e+00	1.074641e+00	1.074641e+00	5.086640e+00	5.086640e+00	0.000000e+00	0.000000e+00	2015.000000
min	0.000000	0.440000	8.456000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	2015.000000
25%	10.000000	1.100000	1.068356e+04	8.540700e+02	3.008780e+03	0.000000e+00	5.086640e+00	2.648420e+03	1.274700e+02	0.000000	2015.000000
50%	24.000000	1.370000	1.073766e+05	6.540300e+03	2.906102e+04	1.849900e+02	1.974383e+04	2.636202e+04	2.647710e+03	0.000000	2016.000000
75%	38.000000	1.650000	4.329623e+05	1.110220e+05	1.502069e+05	6.243420e+03	1.107938e+05	8.333767e+04	2.202925e+04	132.500000	2017.000000
max	52.000000	3.260000	6.250656e+07	2.274362e+07	2.047076e+07	2.545496e+06	1.937313e+07	1.338459e+06	5.719097e+06	55169.650000	2018.000000

Data cleaning

Convert strings to lower or proper case

```
In [18]: # Declaring clean_headers function to clean columns names in data set
def clean_headers(val):
    if isinstance(val, str):
        val="".join(char for char in val if char.isalnum() or char in (" ", "-"))
        val=val.strip().lower().replace(" ", "_")
        return val
    else:
        return val

In [11]: # Calling clean_headers function
avocado=avocado.rename(columns=clean_headers)
avocado.head()
```

```
Out[11]:
```

Unnamed: 0	date	average_price	totalvolume	plu4046	plu4225	plu4770	total_bags	small_bags	large_bags	xlarge_bags	type	year	region
0	0 12/27/2015	1.33	64236.62	1036.74	54545.85	48.16	8696.87	8603.62	93.25	0.0	conventional	2015	Albany
1	1 12/20/2015	1.35	54876.98	674.28	44638.81	58.33	9505.56	9408.07	97.49	0.0	conventional	2015	Albany
2	2 12/13/2015	0.93	118220.22	794.70	109149.67	130.50	8145.35	8042.21	103.14	0.0	conventional	2015	Albany
3	3 12/6/2015	1.08	78992.15	1132.00	71976.41	72.58	5811.16	5677.40	133.76	0.0	conventional	2015	Albany
4	4 11/29/2015	1.28	51039.60	941.48	43838.39	75.78	6183.95	5986.26	197.69	0.0	conventional	2015	Albany

Identify missing data

```
In [12]: avocado.isnull().any() # Checking missing data

Out[12]:
```

Unnamed_0	False
date	False
average_price	False
totalvolume	False
plu4046	False
plu4225	False
plu4770	False
total_bags	False
small_bags	False
large_bags	False
xlarge_bags	False
type	False
year	False
region	False
dtype: bool	

```
In [13]: avocado.isnull().sum() # Detecting missing value using integer

Out[13]:
```

Unnamed_0	0
date	0
average_price	0
totalvolume	0
plu4046	0
plu4225	0
plu4770	0
total_bags	0
small_bags	0
large_bags	0
xlarge_bags	0
type	0
year	0
region	0
dtype: int64	

Find outliers

Outliers for average_price variable

Outliers are the extreme values found in a dataset when exploring it. This indicates that the outlier data points deviate significantly from the predicted values, either being noticeably larger or smaller. Creating a plot is one of the quickest way to detect outliers

Using describe() gives descriptive statistics those that, excluding NaN values, summarize the central tendency, dispersion, and shape of the distribution of a dataset

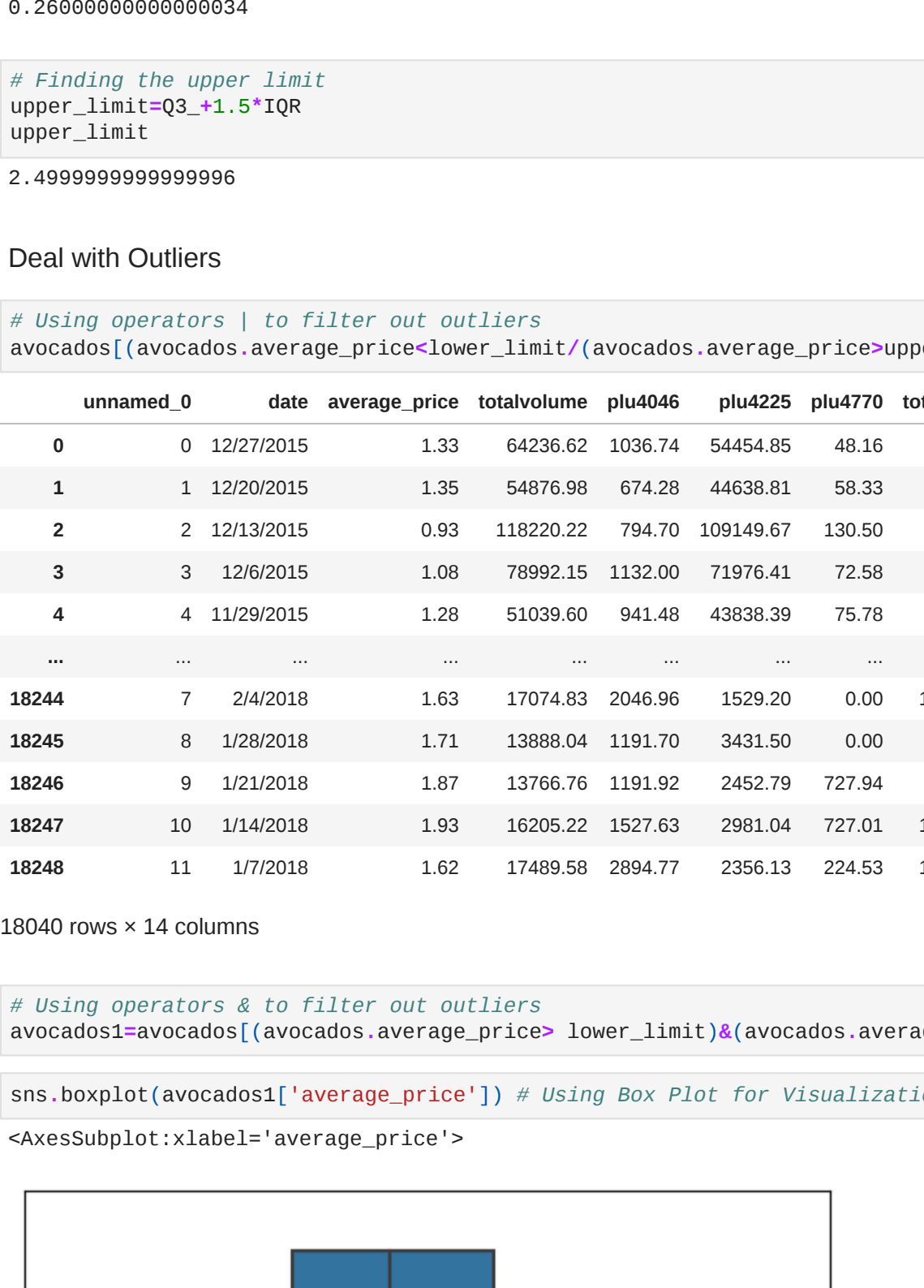
```
In [14]: avocado.average_price.describe() # Looking at descriptive statistics for 'average_price'

Out[14]:
```

count	18249.000000
mean	1.409978
std	0.402677
min	0.440000
25%	1.100000
50%	1.370000
75%	1.650000
max	3.260000
Name: average_price, dtype: float64	

```
In [15]: sns.boxplot(avocado['average_price']) # Using Box Plot for Visualization

Out[15]:
```



```
In [16]: # Checking the highest values where the condition is False
avocado.where(avocado.average_price>=max(avocado.average_price))

Out[16]:
```

Name: average_price, dtype: float64	
14125	3.25

```
In [17]: # Finding the quartile (Q1)
Q1 = avocado.average_price.quantile(0.25)
Q1

Out[17]: 1.1

In [18]: # Finding the upper quartile (Q3)
Q3 = avocado.average_price.quantile(0.75)
Q3

Out[18]: 1.66

In [19]: # Finding Interquartile range(IQR)
IQR= Q3-Q1
IQR

Out[19]: 0.55999999999999998
```

```
In [20]: # Finding the lower limit
lower_limit=Q1-1.5*IQR
lower_limit

Out[20]: 0.26000000000000034

In [21]: # Finding the upper limit
upper_limit=Q3+1.5*IQR
upper_limit

Out[21]: 2.4999999999999996
```

Deal with Outliers

```
In [22]: # Using operators > to filter out outliers
avocado[avocado.average_price>lower_limit][avocado.average_price<upper_limit]]

Out[22]:
```

Unnamed: 0	date	average_price	totalvolume	plu4046	plu4225	plu4770	total_bags	small_bags	large_bags	xlarge_bags	type	year	region
0	0 12/27/2015	1.33	64236.62	1036.74	54545.85	48.16	8696.87	8603.62	93.25	0.0	conventional	2015	Albany
1	1 12/20/2015	1.35	54876.98	674.28	44638.81	58.33	9505.56	9408.07	97.49	0.0	conventional	2015	Albany
2	2 12/13/2015	0.93	118220.22	794.70	109149.67	130.50	8145.35	8042.21	103.14	0.0	conventional	2015	Albany
3	3 12/6/2015	1.08	78992.15	1132.00	71976.41	72.58	5811.16	5677.40	133.76	0.0	conventional	2015	Albany
4	4 11/29/2015	1.28	51039.60	941.48	43838.39	75.78	6183.95	5986.26	197.69	0.0	conventional	2015	Albany

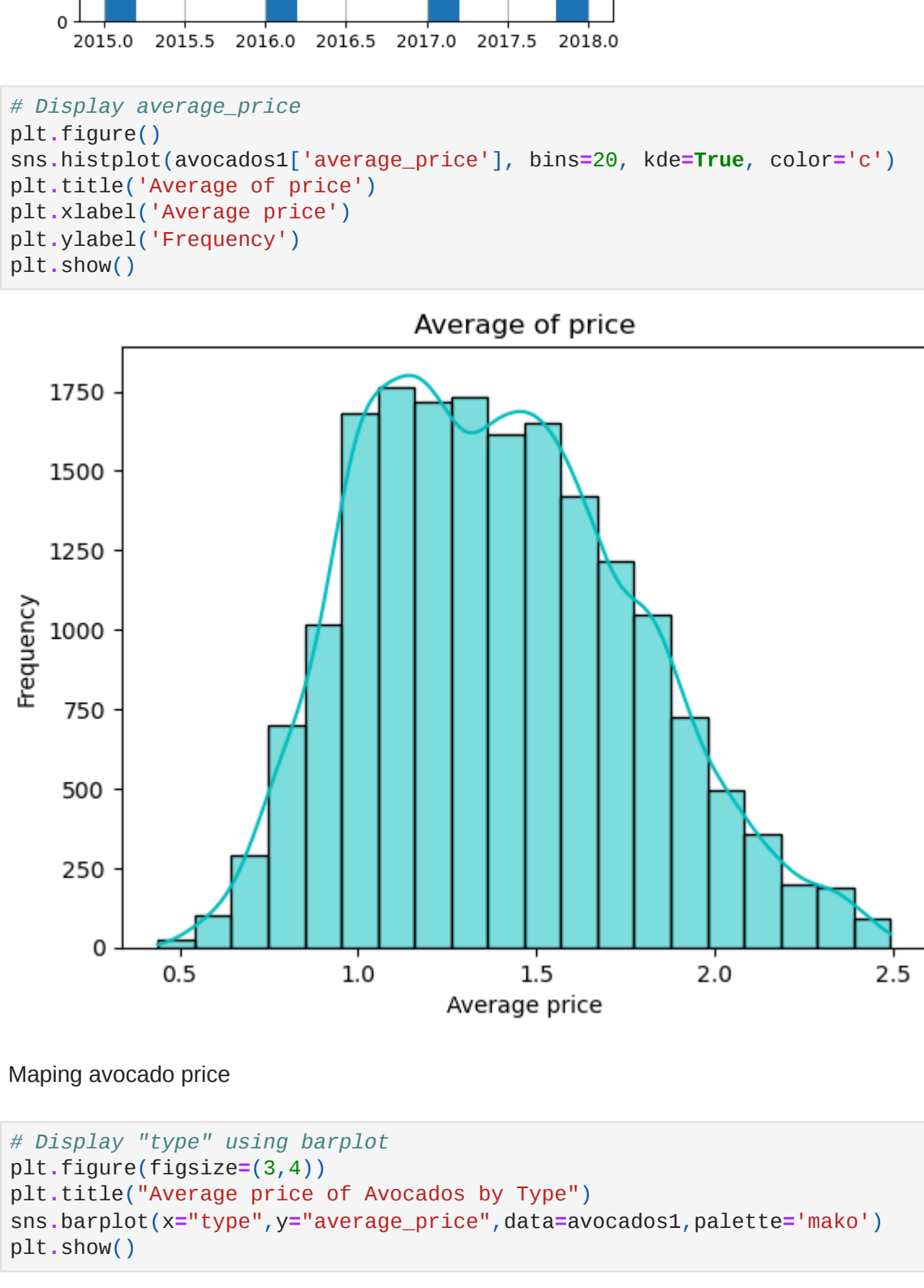
18244	7	2/20/2018	1.63	17074.83	2046.96	1529.20	0.00	13498.67	13065.82	431.85	0.0	organic	2018	WestTexasNewMexico
18245	8	1/28/2018	1.71	13888.04	1191.70	3431.50	0.00	9264.84	8940.04	324.80	0.0	organic	2018	WestTexasNewMexico
18246	9	1/21/2018	1.87	13766.76	1191.92	2452.79	727.84	9394.11	9351.80	42.31	0.0	organic	2018	WestTexasNewMexico
18247	10	1/14/2018	1.93	16205.22	1527.63	2981.04	727.01	10969.64	10919.54	50.00	0.0	organic	2018	WestTexasNewMexico
18248	11	1/7/2018	1.62	17489.58	2864.77	2356.13	224.53	12014.15	11968.14	26.01	0.0	organic	2018	WestTexasNewMexico

10240 rows x 14 columns

```
In [23]: # Using operators & to filter out outliers
avocado=avocado[avocado.average_price>lower_limit][avocado.average_price<upper_limit]]

In [24]: sns.boxplot(avocado['average_price']) # Using Box Plot for Visualization without outliers

Out[24]:
```



Find duplicates

The duplicated() method returns a Series with True and False values to show which rows in the DataFrame are duplicated and which are not. By default, if all the values in a row are the same, duplicated() considers the entire row as a duplicate. The method also considers the first occurrence of a row as unique, so it will always return False for the initial row, as a duplicate row won't appear until a subsequent occurrence.

```
In [25]: avocado.duplicated() # Finding duplicate value in the data

Out[25]:
```

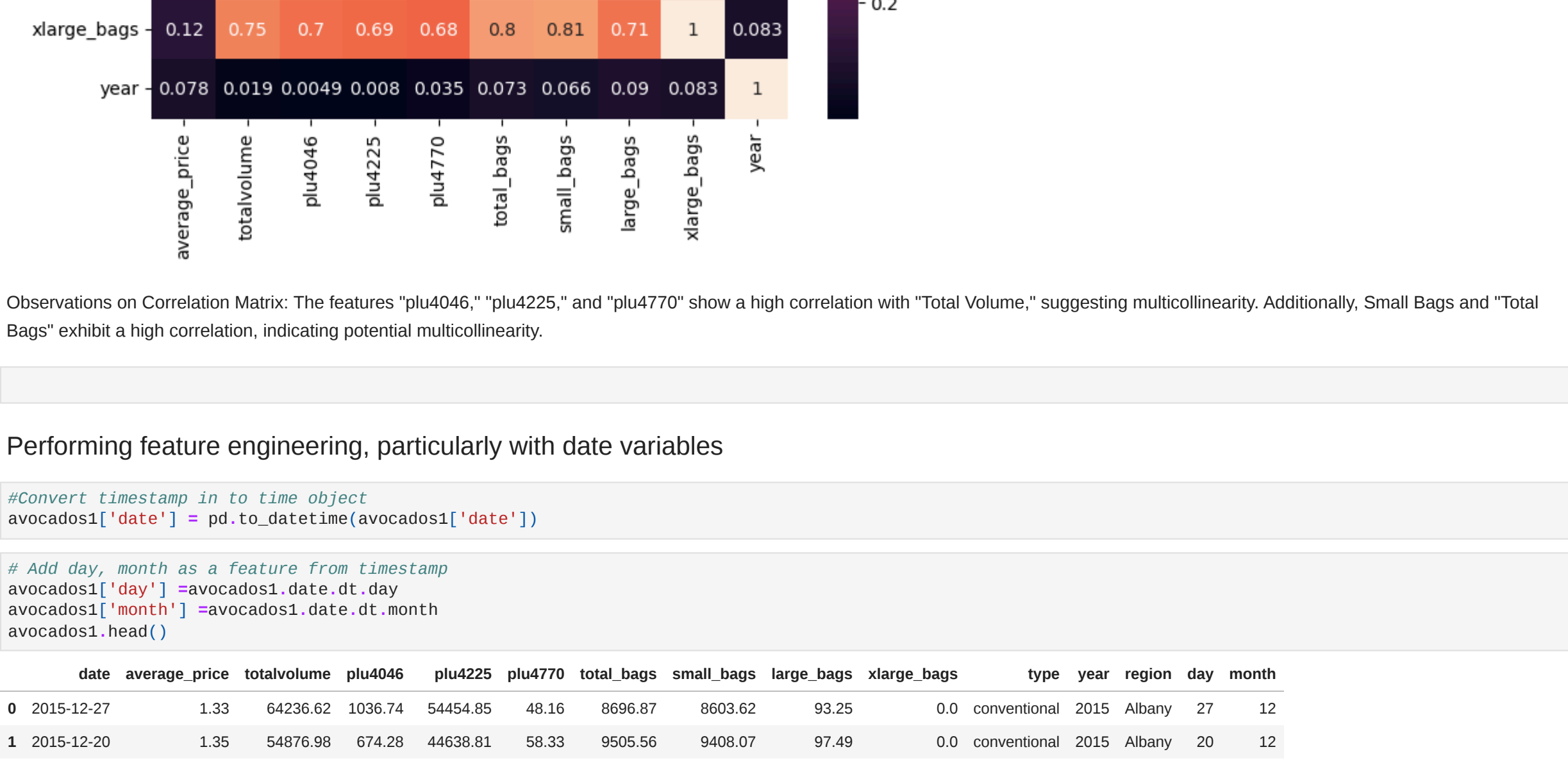
0	False
1	False
2	False
3	False
4	False
...	...
18244	False
18245	False
18246	False
18247	False
18248	False
Length: 18848, dtype: bool	

```
In [26]: # Drop the index column if it's not needed
avocado1= avocado.d.drop(columns="Unnamed: 0")

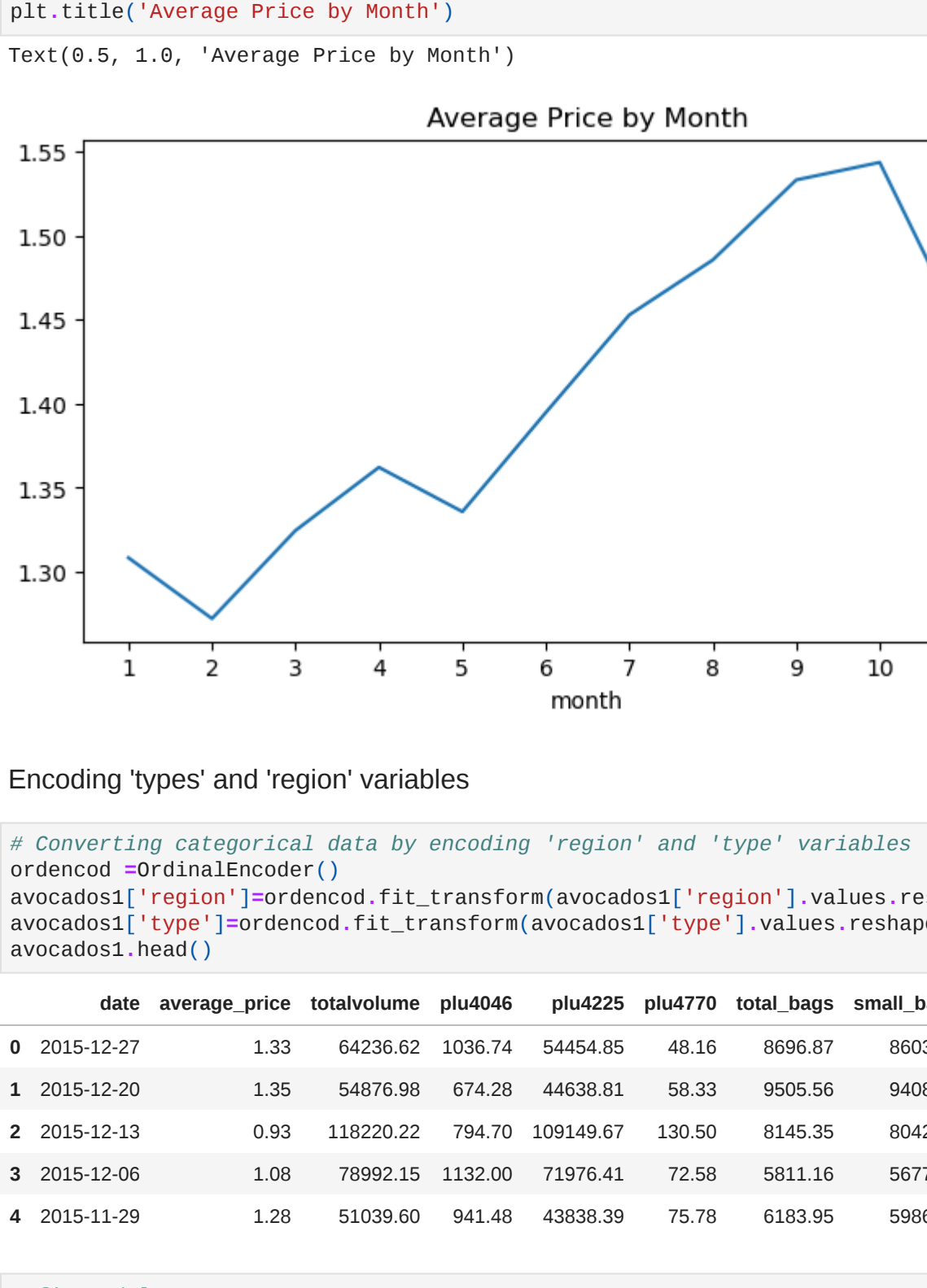
Data visualization
```

```
In [27]: # Hist plots the values and their frequencies as a bar graph
avocado1.hist(bins=15, figsize=(20,10))

Out[27]:
```

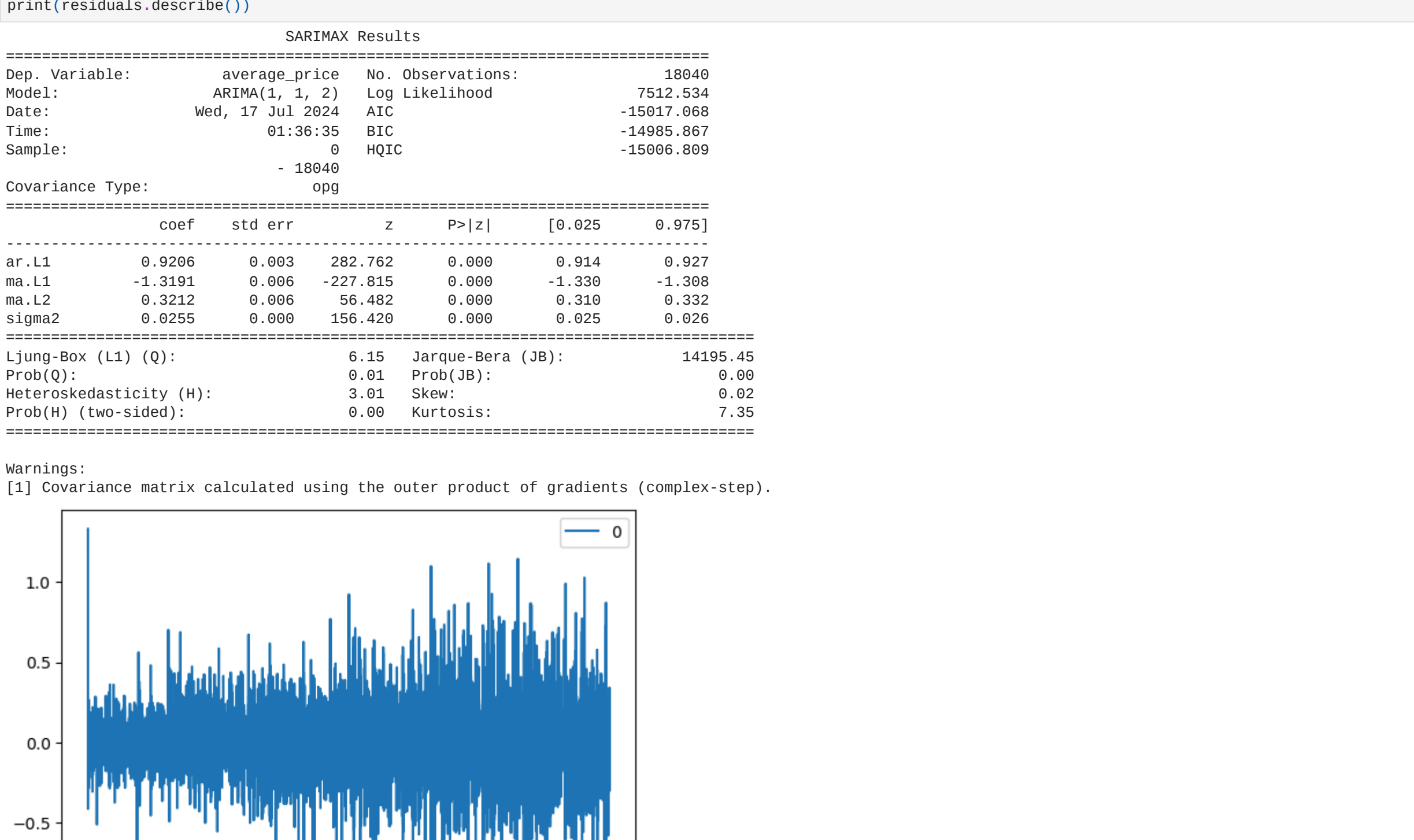


```
In [28]: # Display average price
plt.figure()
x= avocado1['average_price'].values
plt.figure(figsize=(3,4))
plt.title('Average price of Avocado by Type')
plt.xlabel('Average price')
plt.ylabel('Frequency')
plt.show()
```



```
In [29]: # Display 'type' using barplot
plt.figure(figsize=(3,4))
plt.title('Average price of Avocado by Type')
plt.xlabel('type', y='average_price', data=avocado1.palette('nako'))
plt.show()

Out[29]:
```



Observations on Correlation Matrix: The features 'plu4046' and 'plu4770' show a high correlation with 'Total Volume,' suggesting multicollinearity. Additionally, 'Small Bags' and 'Total Bags' exhibit a high correlation, indicating potential multicollinearity.

Performing feature engineering, particularly with date variables

```
In [31]: #Convert timestamp to time object
avocado1['date'] = pd.to_datetime(avocado1['date'])

In [32]: # Add day, month as a feature from timestamp
avocado1['day']=avocado1.date.dt.day
avocado1['month']=avocado1.date.dt.month
avocado1.head()
```

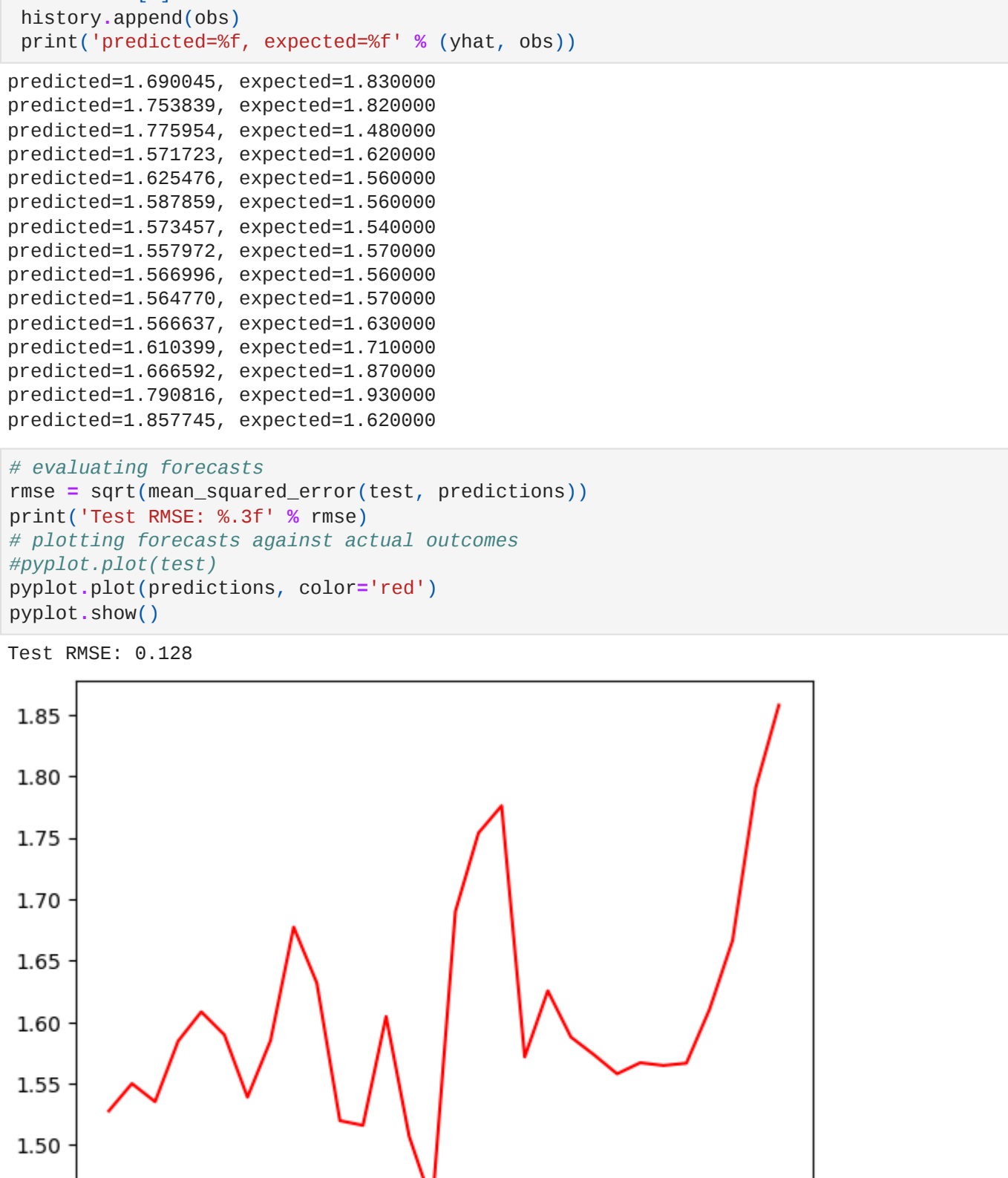
```
Out[32]:
```

date	average_price	totalvolume	plu4046	plu4225	plu4770	total_bags	small_bags	large_bags	xlarge_bags	type	year	region	day	month
0 2015-12-27	1.33	64236.62	1036.74	54545.85	48.16	8696.87	8603.62	93.25	0.0	conventional	2015	Albany	27	12
1 2015-12-20	1.35	54876.98	674.28	44638.81	58.33	9505.56	9408.07	97.49	0.0	conventional	2015	Albany	20	12
2 2015-12-13	0.93	118220.22	794.70	109149.67	130.50	8145.35	8042.21	103.14	0.0	conventional	2015	Albany	13	12
3 2015-12-06	1.08	78992.15	1132.00	71976.41	72.58	5811.16	5677.40	133.76	0.0	conventional	2015	Albany	6	12
4 2015-11-29	1.28	51039.60	941.48	43838.39	75.78	6183.95	5986.26	197.69	0.0	conventional	2015	Albany	29	11

Creating time-series graphs for visual analysis

```
In [33]: datagroup=avocado1.groupby('month').mean()
fig, ax = plt.subplots(figsize=(8,4))
ax.xaxis.set(ticksrange(0,35)) # manually set x-ticks
datagroup['average_price'].plot(x=avocado1.month)
plt.title('Average Price by Month')

Out[33]:
```



Encoding types and region variables

```
In [47]: # Converting categorical data by encoding 'region' and 'type' variables
ordencod=OrdinalEncoder()
avocado1['region']=ordencod.fit_transform(avocado1['region']).values.reshape(-1,1)
avocado1['type']=ordencod.fit_transform(avocado1['type']).values.reshape(-1,1)
avocado1.head()
```

```
Out[47]:
```

date	average_price	totalvolume	plu4046	plu4225	plu4770	total_bags	small_bags	large_bags	xlarge_bags	type	region	day	month	
0 2015-12-27	1.33	64236.62	1036.74	54545.85	48.16	8696.87	8603.62	93.25	0.0	0.0	2015	0	27	12
1 2015-12-20	1.35	54876.98	674.28	44638.81	58.33	9505.56	9408.07	97.49	0.0	0.0	2015	0	20	12
2 2015-12-13	0.93	118220.22	794.70	109149.67	130.50	8145.35	8042.21	103.14	0.0	0.0	2015	0	13	12
3 2015-12-06	1.08	78992.15	1132.00	71976.41	72.58	5811.16	5677.40	133.76	0.0	0.0	2015	0	6	12
4 2015-11-29	1.28	51039.60	941.48	43838.39	75.78	6183.95	5986.26	197.69						