

```
[196]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib inline
import warnings
warnings.filterwarnings("ignore")

In [196]: file=pd.read_csv("avocado.csv.zip")

In [197]: file

Out[197]: Unnamed: 0      Date      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      year      region
0      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
1      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
...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...
18247      10      2018-01-14      1.93      16205.22      1527.63      2981.04      727.01      10969.54      10919.54      50.00      0.0      organic      2018      WestTexNewMexico
18248      11      2018-01-07      1.62      17489.58      2894.77      2356.13      224.53      12014.15      11988.14      26.01      0.0      organic      2018      WestTexNewMexico

18249 rows x 14 columns

In [198]: file.head(5)

Out[198]: Unnamed: 0      Date      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      year      region
0      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
1      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
2      1      2015-12-13      0.93      118220.22      794.70      109149.67      130.50      8145.35      8042.21      103.14      0.0      conventional      2015      Albany
3      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
4      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

In [199]: file.tail(5)

Out[199]: Unnamed: 0      Date      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      year      region
18244      7      2018-02-04      1.63      17074.83      2046.90      1529.20      0.00      15488.67      13066.12      431.85      0.0      organic      2018      WestTexNewMexico
18245      8      2018-01-26      1.71      13888.84      1181.70      3481.90      0.00      9264.84      8840.14      324.80      0.0      organic      2018      WestTexNewMexico
18246      9      2018-01-21      1.87      13765.76      1181.92      2462.78      727.94      8294.11      9261.90      42.31      0.0      organic      2018      WestTexNewMexico
18247      10      2018-01-14      1.93      16205.22      1527.63      2981.04      727.01      10969.54      10919.54      50.00      0.0      organic      2018      WestTexNewMexico
18248      11      2018-01-07      1.62      17489.58      2894.77      2356.13      224.53      12014.15      11988.14      26.01      0.0      organic      2018      WestTexNewMexico

In [210]: # checking data shape and size
file.shape

Out[210]: (18249, 14)

In [211]: # checking null values
file.isnull().sum()

Out[211]: Unnamed: 0      0
Date      0
year      0
region      0
Length: 14, dtype: int64

In [212]: file.dtypes

Out[212]: Unnamed: 0      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      year
count      18249.000000      18249.000000      1.824900e+04      1.824900e+04      1.824900e+04      1.824900e+04      1.824900e+04      1.824900e+04      18249.000000      18249.000000
mean      24.232232      1.405978      8.506440e+05      2.930084e+05      2.951546e+05      2.283974e+04      2.396392e+05      1.821947e+05      5.433809e+04      3106.426507      2016.147899
...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...
79%      38.000000      1.660000      4.329623e+05      1.110202e+05      1.502006e+05      6.243420e+03      1.107834e+05      8.333767e+04      2.202925e+04      132.500000      2017.000000
max      52.000000      3.260000      6.250566e+07      2.274362e+07      2.047057e+07      2.546439e+06      1.937313e+07      1.338459e+07      5.719097e+06      551695.650000      2018.000000

8 rows x 11 columns

In [213]: # dropping unnamed:0 because it is repeating the index number

In [214]: file.drop('Unnamed: 0', axis =1, inplace =True)

In [215]: file.head(5)

Out[215]:      Date      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      year      region
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
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
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
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
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

In [261]: # analysis of average price

In [281]: file.AveragePrice.plot(kind='hist' , figsize =(15,10), color ="pink",)
plt.show()

In [281]: plt.figure(figsize =(15,10))
plt.plot(file.AveragePrice)
plt.show()

In [282]: from pandas.plotting import scatter_matrix

In [283]: scatter_matrix(file , figsize=(15,15) ,diagonal="kde" ,color="r")
plt.show()

In [284]: # indicating number of regions
plt.figure(figsize =(50,50))
sns.relplot( x ="region", y ="year" , data =file)
plt.xticks(rotation = 50)
plt.show()

<Figure size 3600x3600 with 0 Axes>

In [285]: file.head()

Out[285]:      Date      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      year      region
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
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
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
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
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

In [286]: file.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18249 entries, 0 to 18248
Data columns (total 13 columns):
#      Column      Non-Null Count      Dtype
---      ---      ---
0      Date      18249 non-null      object
1      AveragePrice      18249 non-null      float64
2      Total Volume      18249 non-null      float64
3      4046      18249 non-null      float64
4      4225      18249 non-null      float64
5      4770      18249 non-null      float64
6      Total Bags      18249 non-null      float64
7      Small Bags      18249 non-null      float64
8      Large Bags      18249 non-null      float64
9      XLarge Bags      18249 non-null      float64
10     type      18249 non-null      int8
11     year      18249 non-null      int64
12     region      18249 non-null      int8
dtypes: float64(9), int64(1), int8(2), object(1)
memory usage: 1.6+ MB

In [288]: plt.figure(figsize=(5,5))
plt.plot(file['Date'], file['AveragePrice'])
plt.show()

In [289]: plt.figure(figsize=(10,5),dpi=80)
sns.boxplot(data = file[[
'AveragePrice',
'Total Volume',
'4046',
'4225',
'4770',
'Total Bags',
'Small Bags',
'Large Bags',
'XLarge Bags']])
plt.show()

In [290]: # now we can see that data is free from outliers
plt.figure(figsize=(10,5),dpi=80)
sns.boxplot(data = file[[
'AveragePrice',
'Total Volume',
'4046',
'4225',
'4770',
'Total Bags',
'Small Bags',
'Large Bags',
'XLarge Bags']])
plt.show()

In [291]: file.drop(columns=["Date"],inplace=True)
file.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18249 entries, 0 to 18248
Data columns (total 12 columns):
#      Column      Non-Null Count      Dtype
---      ---      ---
0      AveragePrice      18249 non-null      float64
1      Total Volume      18249 non-null      float64
2      4046      18249 non-null      float64
3      4225      18249 non-null      float64
4      4770      18249 non-null      float64
5      Total Bags      18249 non-null      float64
6      Small Bags      18249 non-null      float64
7      Large Bags      18249 non-null      float64
8      XLarge Bags      18249 non-null      float64
9      type      18249 non-null      object
10     year      18249 non-null      int64
11     region      18249 non-null      object
dtypes: float64(9), int64(1), object(2)
memory usage: 1.7+ MB

In [292]: from numpy import percentile

columns = file.columns
for i in columns:
    if isinstance(file[i][0], str):
        continue
    else:
        #finding quartiles
        quartiles = percentile(file[i], [20,60])
        # calculate min/max
        lower_fence = quartiles[0] - (1.5*(quartiles[1]-quartiles[0]))
        upper_fence = quartiles[1] + (1.5*(quartiles[1]-quartiles[0]))
        file[i] = file[i].apply(lambda x: upper_fence if x > upper_fence else (lower_fence if x < lower_fence else x))

In [293]: file.head()

Out[293]:      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      year      region
0      1.33      64236.62      1036.74      54545.85      48.16      8696.87      8603.62      93.25      0.0      conventional      2015      Albany
1      1.35      54876.98      674.28      44638.81      58.33      9505.56      9408.07      97.49      0.0      conventional      2015      Albany
2      0.93      118220.22      794.70      109149.67      130.50      8145.35      8042.21      103.14      0.0      conventional      2015      Albany
3      1.08      78992.15      1132.00      71976.41      72.58      5811.16      5677.40      133.76      0.0      conventional      2015      Albany
4      1.28      51039.60      941.48      43838.39      75.78      6183.95      5986.26      197.69      0.0      conventional      2015      Albany

In [294]: # now we can see that data is free from outliers
plt.figure(figsize=(10,5),dpi=80)
sns.boxplot(data = file[[
'AveragePrice',
'Total Volume',
'4046',
'4225',
'4770',
'Total Bags',
'Small Bags',
'Large Bags',
'XLarge Bags']])
plt.show()

In [295]: file['region'] = pd.Categorical(file['region'])
dummies_region = pd.get_dummies(file['region'], prefix = 'region')
dummies_region

Out[295]:      region_Albania      region_Atlanta      region_BaltimoreWashington      region_Boston      region_Boston      region_BuffaloRochester      region_California      region_Charlotte      region_Chicago      region_CincinnatiDayton      region_SouthCarolina      region_South
0      0      1      0      0      0      0      0      0      0      0      0      0
1      1      1      0      0      0      0      0      0      0      0      0      0
...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...
18247      0      0      0      0      0      0      0      0      0      0      0      0
18248      0      0      0      0      0      0      0      0      0      0      0      0

18249 rows x 54 columns

In [296]: dataset = pd.concat([file, dummies_region], axis=1)
dataset.drop(columns="region",inplace=True)
dataset

Out[296]:      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      region_SouthCarolina      region_SouthCentral      region_Southeast      region_Spokane      region_StLouis      region_Syracuse      region_Tampa
0      1.33      64236.62      1036.74      54545.85      48.16      8696.87      8603.62      93.25      0.0      conventional      ...      0      0      0      0      0      0
1      1.35      54876.98      674.28      44638.81      58.33      9505.56      9408.07      97.49      0.0      conventional      ...      0      0      0      0      0      0
...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...
18247      1.93      16205.22      1527.63      2981.04      727.01      10969.54      10919.54      50.00      organic      ...      0      0      0      0      0      0
18248      1.62      17489.58      2894.77      2356.13      224.53      12014.15      11988.14      26.01      organic      ...      0      0      0      0      0      0

18249 rows x 65 columns

In [296]: dataset

Out[296]:      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      region_SouthCarolina      region_SouthCentral      region_Southeast      region_Spokane      region_StLouis      region_Syracuse      region_Tampa
0      1.33      64236.62      1036.74      54545.85      48.16      8696.87      8603.62      93.25      0.0      conventional      ...      0      0      0      0      0      0
1      1.35      54876.98      674.28      44638.81      58.33      9505.56      9408.07      97.49      0.0      conventional      ...      0      0      0      0      0      0
...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...
18247      1.93      16205.22      1527.63      2981.04      727.01      10969.54      10919.54      50.00      organic      ...      0      0      0      0      0      0
18248      1.62      17489.58      2894.77      2356.13      224.53      12014.15      11988.14      26.01      organic      ...      0      0      0      0      0      0

18249 rows x 65 columns

In [297]: from sklearn import preprocessing

label_encoder = preprocessing.LabelEncoder()
dataset[['type']] = label_encoder.fit_transform(dataset['type'])
dataset

Out[297]:      AveragePrice      Total Volume      4046      4225      4770      Total Bags      Small Bags      Large Bags      XLarge Bags      type      region_SouthCarolina      region_SouthCentral      region_Southeast      region_Spokane      region_StLouis      region_Syracuse      region_Tampa
0      0      1.33      64236.62      1036.74      54545.85      48.16      8696.87      8603.62      93.25      0.0      0      ...      0      0      0      0      0
1      1      1.35      54876.98      674.28      44638.81      58.33      9505.56      9408.07      97.49      0.0      0      ...      0      0      0      0      0
...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...
18247      10     1.93      16205.22      1527.63      2981.04      727.01      10969.54      10919.54      50.00      1      ...      0      0      0      0      0      0
18248      11     1.62      17489.58      2894.77      2356.13      224.53      12014.15      11988.14      26.01      1      ...      0      0      0      0      0      0

18249 rows x 65 columns

In [298]: # fitting multiple linear regression model

regressor=LinearRegression()
regressor.fit(X=train,y=train)
y_pred = regressor.predict(X=test)
regression_results(y_test,y_pred)
model_accuracy(regressor)

Explained variance: 0.6294
R2: 0.6293
Adjusted_r2: 0.6227
MSE: 0.1799
RMSE: 0.4247
Accuracy: 60.14 %
Standard Deviation: 1.55 %

In [299]: # fitting random forest regression model

from sklearn.ensemble import RandomForestRegressor

rand_regressor = RandomForestRegressor()
rand_regressor.fit(X=train,y=train)
y_pred_rf = rand_regressor.predict(X=test)
regression_results(y_test,y_pred_rf)
model_accuracy(rand_regressor)

Explained variance: 0.8504
R2: 0.8503
Adjusted_r2: 0.8476
MSE: 0.1264
RMSE: 0.3557
Accuracy: 83.48 %
Standard Deviation: 0.81 %

In [300]: # splitting data into X and y
X = dataset.drop(['AveragePrice'],axis=1)
y = dataset['AveragePrice']

In [300]: # splitting data into training and testing dataset
X_train,X_test,y_train,y_test = train_test_split(X,
                                                    test_size=0.5,
                                                    random_state=10)

In [300]: print('training set:',X_train.shape,' :: ', 'samples ::',y_train.shape[0])
print('testing set:',X_test.shape,' :: ', 'samples ::',y_test.shape[0])

In [301]:

In [301]:

In [301]:

In [301]:

In [301]:

In [301]:

In [301]:
```