



# Prediction of Average Prices of Avocado in USA

## Context

It is a well known fact that Millenials LOVE Avocado Toast. It's also a well known fact that all Millenials live in their parents basements.

Clearly, they aren't buying home because they are buying too much Avocado Toast!

But maybe there's hope... if a Millenial could find a city with cheap avocados, they could live out the Millenial American Dream.

## Content

The data represents weekly 2018 retail scan data for National retail volume (units) and price. Retail scan data comes directly from retailers' cash registers based on actual retail sales of Hass avocados. Starting in 2013, the data reflects an expanded, multi-outlet retail data set. Multi-outlet reporting includes an aggregation of the following channels: grocery, mass, club, drug, dollar and military. The Average Price (of avocados) in the data reflects a per unit (per avocado) cost, even when multiple units (avocados) are sold in bags. The Product Lookup codes (PLU's) in the table are only for Hass avocados. Other varieties of avocados (e.g. greenskins) are not included in this data.

## Columns in our dataset

- index
- Date: The date of the observation
- AveragePrice: The average price of a single avocado
- Total Volume: Total number of avocados sold
- 4046: Total number of avocados with PLU 4046 sold
- 4225: Total number of avocados with PLU 4225 sold
- 4770: Total number of avocados with PLU 4770 sold
- Total Bags
- Small Bags
- Large Bags
- XLarge Bags
- type: conventional or organic
- year: The year
- region: The city or region of the observation

## Step 1: Importing our Dataset

In [62]:

```
import pandas as pd
dataset = pd.read_csv("avocado.csv")
dataset.head()
```

Out[62]:

|   | Unnamed:<br>0 | Date       | AveragePrice | Total<br>Volume | 4046    | 4225      | 4770   | Total<br>Bags | Small<br>Bags | Large<br>Bags | XL<br>B |
|---|---------------|------------|--------------|-----------------|---------|-----------|--------|---------------|---------------|---------------|---------|
| 0 | 0             | 2015-12-27 | 1.33         | 64236.62        | 1036.74 | 54454.85  | 48.16  | 8696.87       | 8603.62       | 93.25         |         |
| 1 | 1             | 2015-12-20 | 1.35         | 54876.98        | 674.28  | 44638.81  | 58.33  | 9505.56       | 9408.07       | 97.49         |         |
| 2 | 2             | 2015-12-13 | 0.93         | 118220.22       | 794.70  | 109149.67 | 130.50 | 8145.35       | 8042.21       | 103.14        |         |
| 3 | 3             | 2015-12-06 | 1.08         | 78992.15        | 1132.00 | 71976.41  | 72.58  | 5811.16       | 5677.40       | 133.76        |         |
| 4 | 4             | 2015-11-29 | 1.28         | 51039.60        | 941.48  | 43838.39  | 75.78  | 6183.95       | 5986.26       | 197.69        |         |

We are using the pandas library to read the csv file named avocado. The file contains information on the following features:

- Date - The date of the observation
- AveragePrice - the average price of a single avocado
- Total Volume - Total number of avocados sold
- 4046 - Total number of avocados with PLU 4046 sold
- 4225 - Total number of avocados with PLU 4225 sold
- 4770 - Total number of avocados with PLU 4770 sold
- Total Bags
- Small Bags
- Large Bags
- XLarge Bags
- type - conventional or organic
- year - the year
- region - the city or region of the observation

In [63]:

```
dataset.tail()
```

Out[63]:

|       | Unnamed:<br>0 | Date       | AveragePrice | Total<br>Volume | 4046    | 4225    | 4770   | Total<br>Bags | Small<br>Bags | Large<br>Bags |
|-------|---------------|------------|--------------|-----------------|---------|---------|--------|---------------|---------------|---------------|
| 18244 | 7             | 2018-02-04 | 1.63         | 17074.83        | 2046.96 | 1529.20 | 0.00   | 13498.67      | 13066.82      | 431.85        |
| 18245 | 8             | 2018-01-28 | 1.71         | 13888.04        | 1191.70 | 3431.50 | 0.00   | 9264.84       | 8940.04       | 324.80        |
| 18246 | 9             | 2018-01-21 | 1.87         | 13766.76        | 1191.92 | 2452.79 | 727.94 | 9394.11       | 9351.80       | 42.31         |
| 18247 | 10            | 2018-01-14 | 1.93         | 16205.22        | 1527.63 | 2981.04 | 727.01 | 10969.54      | 10919.54      | 50.00         |
| 18248 | 11            | 2018-01-07 | 1.62         | 17489.58        | 2894.77 | 2356.13 | 224.53 | 12014.15      | 11988.14      | 26.01         |

◀ ▶

In [64]: `dataset.drop('Unnamed: 0', axis=1, inplace=True)`

**The Feature "Unnamed:0" is just a representation of the indexes, so it's useless to keep it, lets remove it !**

In [65]: `dataset.head()`

Out[65]:

|   | Date       | AveragePrice | Total<br>Volume | 4046    | 4225      | 4770   | Total<br>Bags | Small<br>Bags | Large<br>Bags | XLarge<br>Bags |        |
|---|------------|--------------|-----------------|---------|-----------|--------|---------------|---------------|---------------|----------------|--------|
| 0 | 2015-12-27 | 1.33         | 64236.62        | 1036.74 | 54454.85  | 48.16  | 8696.87       | 8603.62       | 93.25         | 0.0            | conven |
| 1 | 2015-12-20 | 1.35         | 54876.98        | 674.28  | 44638.81  | 58.33  | 9505.56       | 9408.07       | 97.49         | 0.0            | conven |
| 2 | 2015-12-13 | 0.93         | 118220.22       | 794.70  | 109149.67 | 130.50 | 8145.35       | 8042.21       | 103.14        | 0.0            | conven |
| 3 | 2015-12-06 | 1.08         | 78992.15        | 1132.00 | 71976.41  | 72.58  | 5811.16       | 5677.40       | 133.76        | 0.0            | conven |
| 4 | 2015-11-29 | 1.28         | 51039.60        | 941.48  | 43838.39  | 75.78  | 6183.95       | 5986.26       | 197.69        | 0.0            | conven |

◀ ▶

In [66]: `dataset.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  object
11  year            18249 non-null  int64
12  region          18249 non-null  object
dtypes: float64(9), int64(1), object(3)
memory usage: 1.8+ MB

```

**Well as a first observation we can see that we are lucky, we dont have any missing values (18249 complete data) and 13 columns. Now let's do some Feature Engineering on the Date Feature so we can be able to use the day and the month columns in building our machine learning model later.**

In [67]:

```

dataset['Date']=pd.to_datetime(dataset['Date'])
dataset['Month']=dataset['Date'].apply(lambda x:x.month)
dataset['Day']=dataset['Date'].apply(lambda x:x.day)
dataset.head()

```

Out[67]:

|   | Date       | AveragePrice | Total Volume | 4046    | 4225      | 4770   | Total Bags | Small Bags | Large Bags | XLarge Bags |        |
|---|------------|--------------|--------------|---------|-----------|--------|------------|------------|------------|-------------|--------|
| 0 | 2015-12-27 | 1.33         | 64236.62     | 1036.74 | 54454.85  | 48.16  | 8696.87    | 8603.62    | 93.25      | 0.0         | conven |
| 1 | 2015-12-20 | 1.35         | 54876.98     | 674.28  | 44638.81  | 58.33  | 9505.56    | 9408.07    | 97.49      | 0.0         | conven |
| 2 | 2015-12-13 | 0.93         | 118220.22    | 794.70  | 109149.67 | 130.50 | 8145.35    | 8042.21    | 103.14     | 0.0         | conven |
| 3 | 2015-12-06 | 1.08         | 78992.15     | 1132.00 | 71976.41  | 72.58  | 5811.16    | 5677.40    | 133.76     | 0.0         | conven |
| 4 | 2015-11-29 | 1.28         | 51039.60     | 941.48  | 43838.39  | 75.78  | 6183.95    | 5986.26    | 197.69     | 0.0         | conven |

Here we add two more columns month and day where 6 in Month implies the month "June" & 15 in Day implies the 15th day of a month.

## Step 2: Analysis of Average Prices

In [68]:

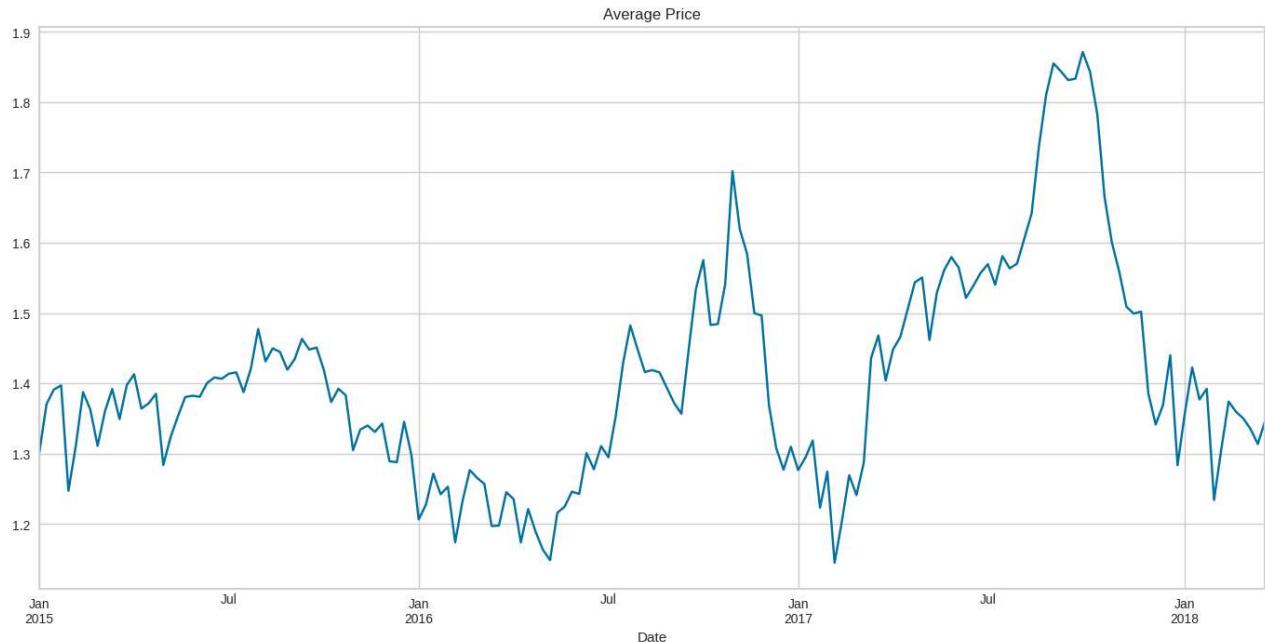
```

import matplotlib.pyplot as plt

byDate=dataset.groupby('Date').mean()
plt.figure(figsize=(17,8),dpi=100)
byDate['AveragePrice'].plot()
plt.title('Average Price')

```

Out[68]: Text(0.5, 1.0, 'Average Price')



Hence the plot shows the average price of avocado at various points of time

In [69]:

```
byMonth = dataset.groupby("Month").mean()
plt.figure(figsize=(17,8),dpi=100)
plt.plot(["Jan","Feb","Mar","Apr","May","Jun","Jul","Aug","Sept","Oct","Nov","Dec"],byM
plt.title('Average Price Per Month')
```

Out[69]: Text(0.5, 1.0, 'Average Price Per Month')



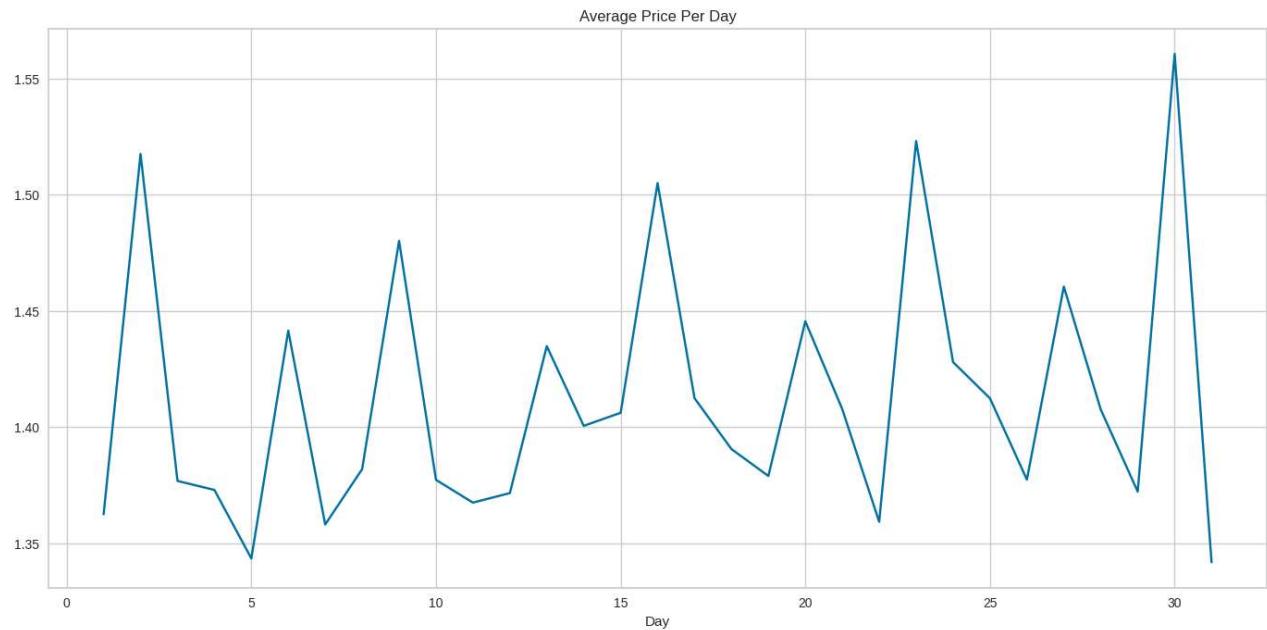
From the above graph plotted for average price of avocado per month we can observe that the price rises for a while in February to March then it falls in April and then the month of May witnesses a rise in the average price. This rise reaches its zenith in the month of October and henceforth it starts to fall.

In [70]:

```
byDay = dataset.groupby("Day").mean()
plt.figure(figsize=(17,8),dpi=100)
```

```
byDay['AveragePrice'].plot()
plt.title('Average Price Per Day')
```

Out[70]: Text(0.5, 1.0, 'Average Price Per Day')



**The above graph for average price per day implies that the price fluctuates in a similar manner at a regular interval.**

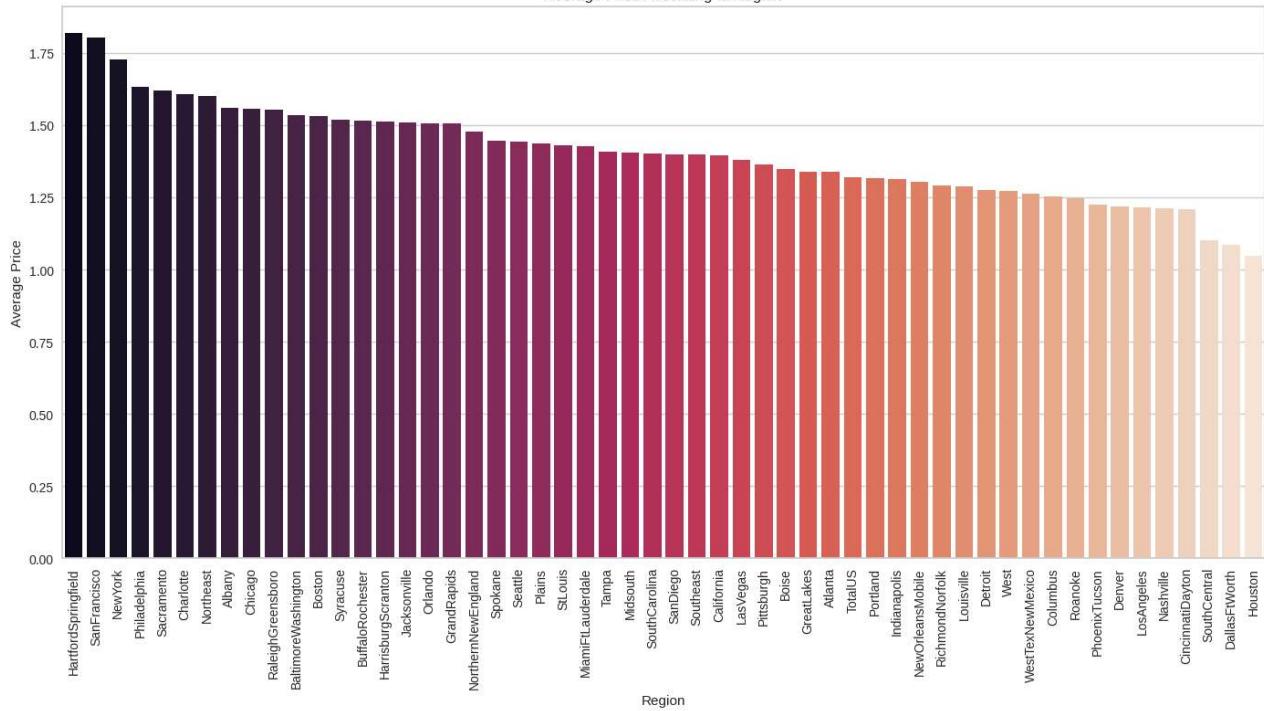
In [71]:

```
import seaborn as sns

byRegion=dataset.groupby('region').mean()
byRegion.sort_values(by=['AveragePrice'], ascending=False, inplace=True)
plt.figure(figsize=(17,8),dpi=100)
sns.barplot(x = byRegion.index,y=byRegion["AveragePrice"],data = byRegion,palette='rocket')
plt.xticks(rotation=90)
plt.xlabel('Region')
plt.ylabel('Average Price')
plt.title('Average Price According to Region')
```

Out[71]: Text(0.5, 1.0, 'Average Price According to Region')

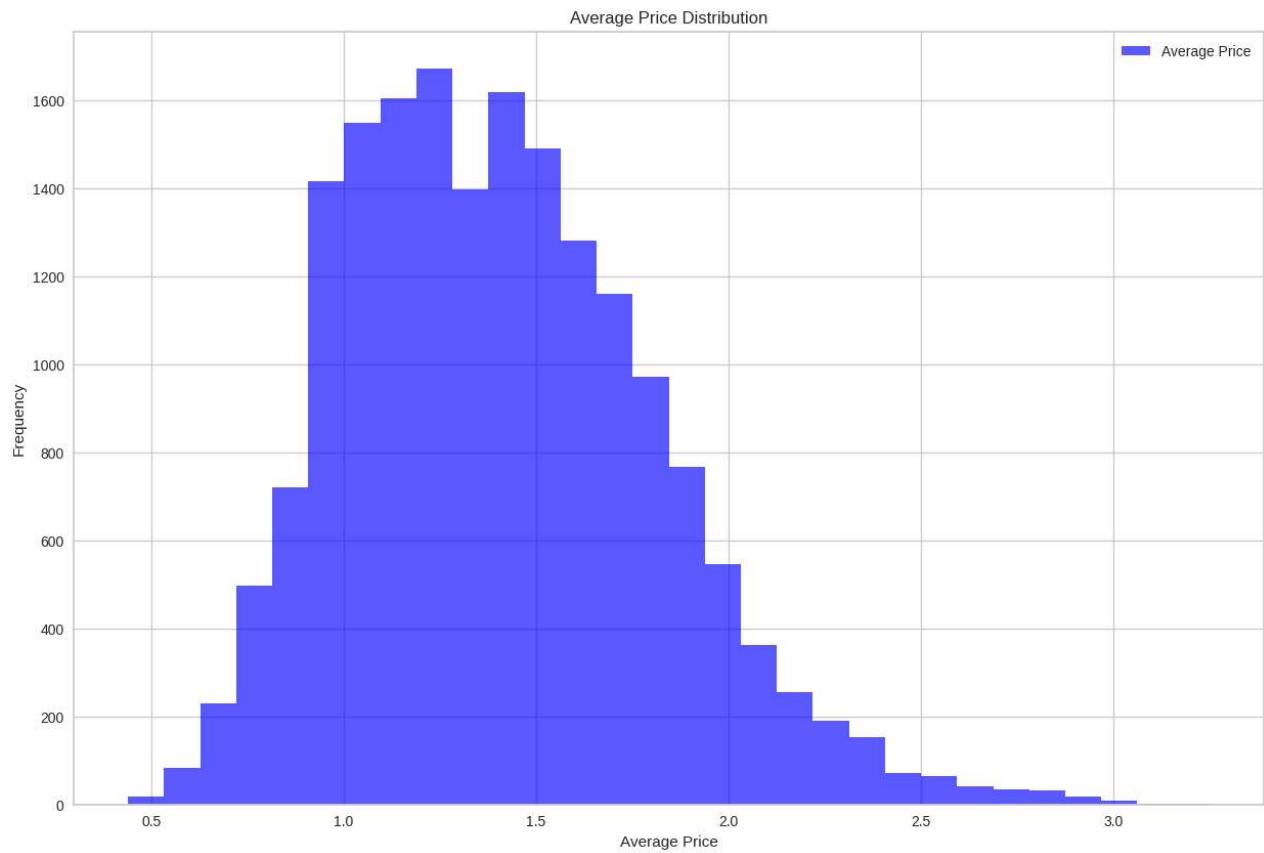
Average Price According to Region



The barplot shows the average price of avocado at various regions in a ascending order. Clearly Hartford Springfield, SanFrancisco, NewYork are the regions with the highest avocado prices.

In [72]:

```
plt.figure(figsize=(15,10),dpi=100)
dataset[ "AveragePrice" ].plot(kind="hist",color="blue",bins=30,grid=True,alpha=0.65,lab
plt.legend()
plt.xlabel("Average Price")
plt.title("Average Price Distribution")
plt.show()
```



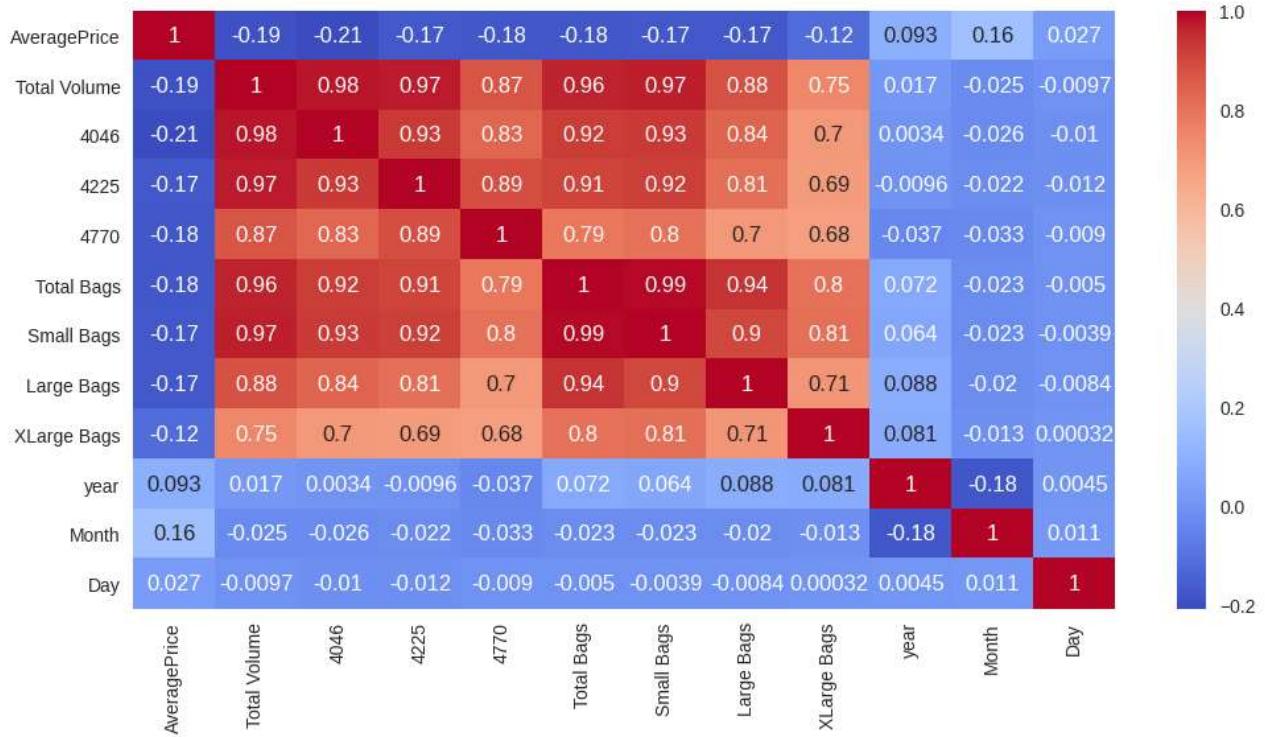
The above histogram for the average price of avocado suggests that its distribution is somewhat positively skewed.

In [73]:

```
import numpy as np

corr_df = dataset.corr(method='pearson')
plt.figure(figsize=(12,6),dpi=100)
sns.heatmap(corr_df,cmap='coolwarm',annot=True)
```

Out[73]: <AxesSubplot:>



**As we can from the heatmap above, all the Features are not correlated with the Average Price column, instead most of them are correlated with each other.**

In [74]:

```
sns.factorplot('AveragePrice','region',data=dataset,
               hue='year',
               aspect=0.8,
               height=15,
               palette='magma',
               join=False,
)
```

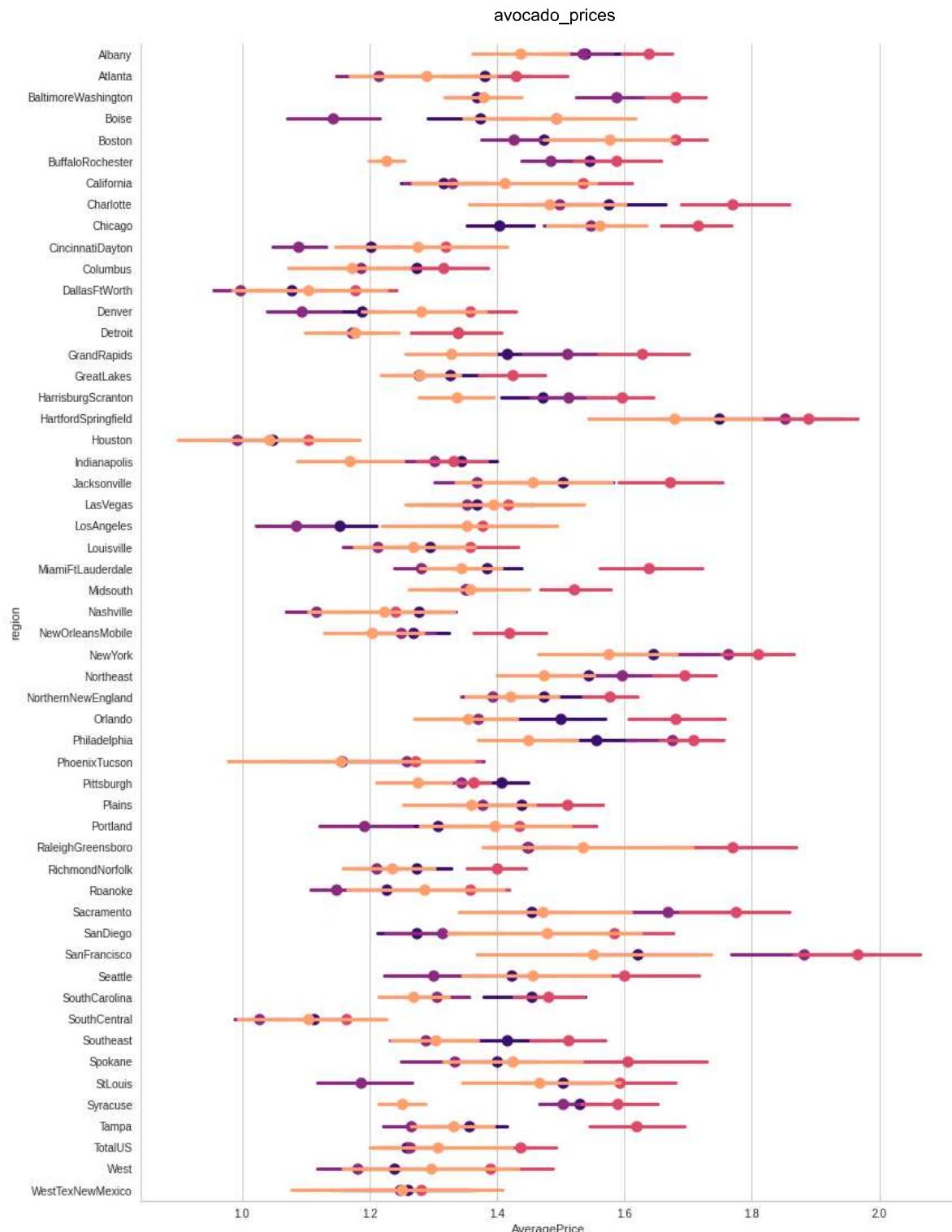
/opt/conda/lib/python3.8/site-packages/seaborn/categorical.py:3714: UserWarning: The `factorplot` function has been renamed to `catplot`. The original name will be removed in a future release. Please update your code. Note that the default `kind` in `factorplot`(`'point'`) has changed `'strip'` in `catplot`.

warnings.warn(msg)

/opt/conda/lib/python3.8/site-packages/seaborn/\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid position argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(

Out[74]: &lt;seaborn.axisgrid.FacetGrid at 0x7f67629b0970&gt;



**A factor plot is simply the same plot generated for different response and factor variables and arranged on a single page. The underlying plot generated can be any univariate or bivariate plot. The scatter plot is the most common application. The above plot is a factor plot of average avocado price for different regions classified by year.**

In [75]:

```
dataset_vif = dataset.copy()
dataset_vif.drop(columns=['Date', 'type', 'region'], inplace = True)

from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.tools.tools import add_constant

Xf = add_constant(dataset_vif)
```

```
pd.Series([variance_inflation_factor(Xf.values, i)
           for i in range(Xf.shape[1])],
           index=Xf.columns)
```

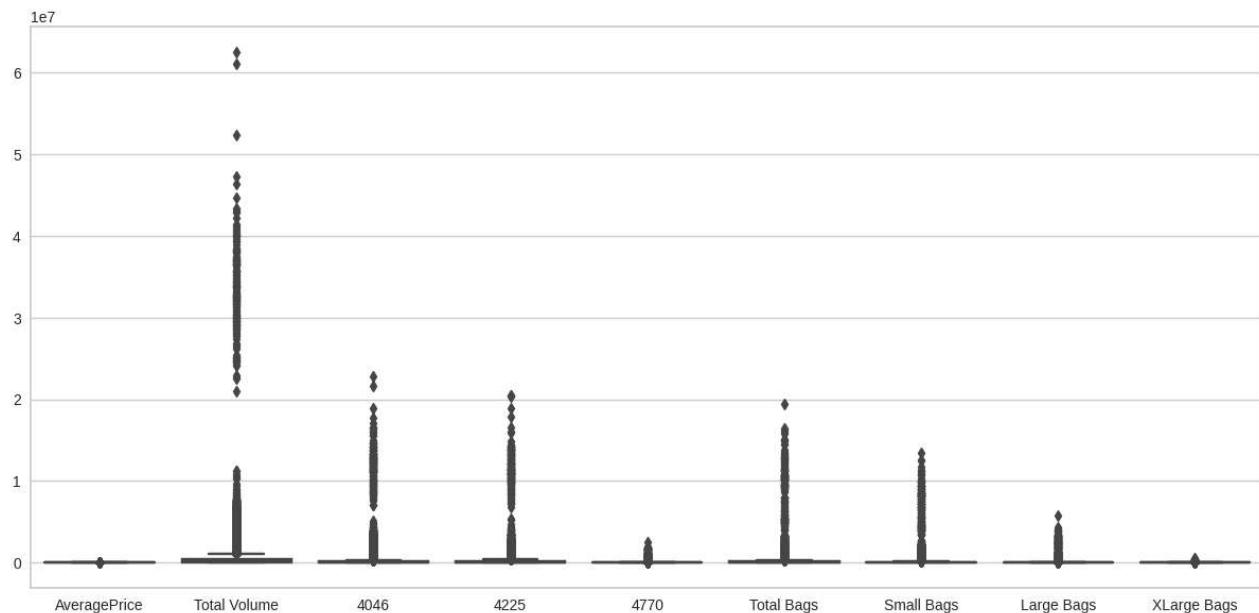
```
Out[75]: const      5.068485e+06
AveragePrice  1.099766e+00
Total Volume   4.918067e+09
4046          6.598339e+08
4225          5.978631e+08
4770          4.762133e+06
Total Bags    2.370316e+14
Small Bags    1.364727e+14
Large Bags    1.448103e+13
XLarge Bags   7.622174e+10
year          1.101665e+00
Month         1.071816e+00
Day           1.001467e+00
dtype: float64
```

The above code snippet calculates the variable inflation factor for the displayed variables.

## Step 3: Taking Care of the Outliers

```
In [76]: plt.figure(figsize=(15,7),dpi=100)
sns.boxplot(data = dataset[['
    'AveragePrice',
    'Total Volume',
    '4046',
    '4225',
    '4770',
    'Total Bags',
    'Small Bags',
    'Large Bags',
    'XLarge Bags']]])
```

```
Out[76]: <AxesSubplot:>
```



Clearly the boxplot indicates that all the variables contains outliers. Now we need to take care of the outliers.

In [77]:

```
dataset.drop(columns=["Date"], inplace=True)
dataset.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18249 entries, 0 to 18248
Data columns (total 14 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 
 12  Month          18249 non-null   int64  
 13  Day            18249 non-null   int64  
dtypes: float64(9), int64(3), object(2)
memory usage: 1.9+ MB
```

**Before we go on to taking care of the outliers we removed the "Date" variable from our dataset as it is useless now.**

In [78]:

```
import numpy as np
from numpy import percentile

columns = dataset.columns
for j in columns:
    if isinstance(dataset[j][0], str):
        continue
    else:
        #defining quartiles
        quartiles = percentile(dataset[j], [25,75])
        # calculate min/max
        lower_fence = quartiles[0] - (1.5*(quartiles[1]-quartiles[0]))
        upper_fence = quartiles[1] + (1.5*(quartiles[1]-quartiles[0]))
        dataset[j] = dataset[j].apply(lambda x: upper_fence if x > upper_fence else (lo
```

**In the following code snippet we have replaced the outliers higher than the upper whisker by the value of the upper whisker and the outliers lower than the lower whisker by the value of the lower whisker.**

In [79]:

```
dataset.head()
```

Out[79]:

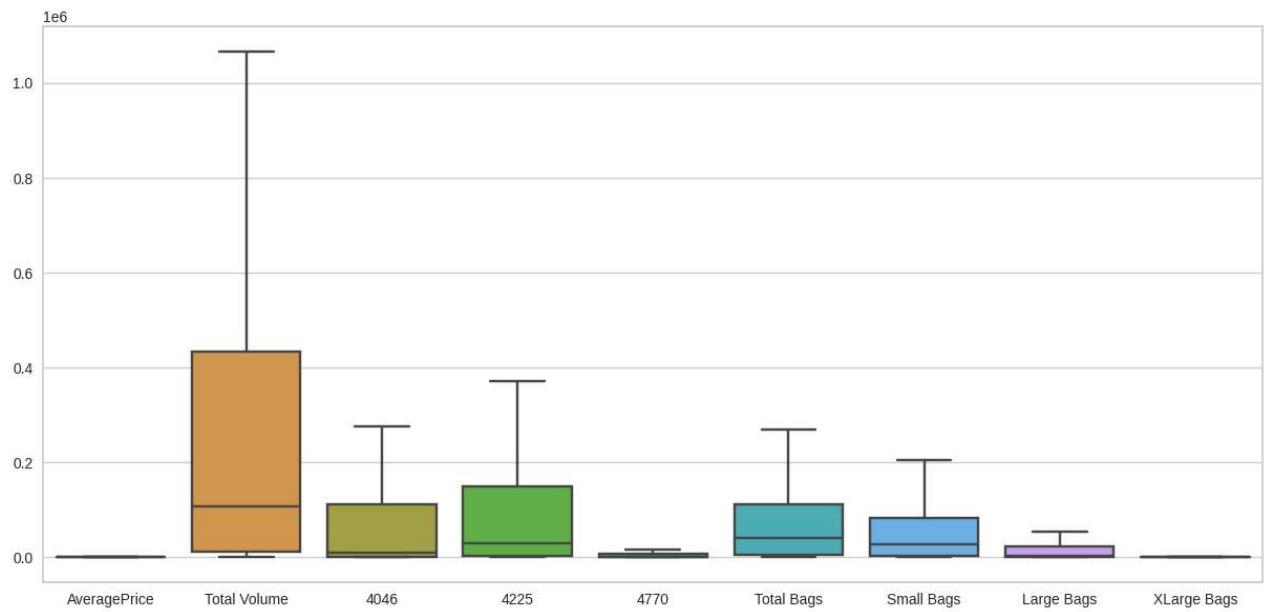
|   | AveragePrice | Total Volume | 4046    | 4225      | 4770   | Total Bags | Small Bags | Large Bags | XLarge Bags | type         |
|---|--------------|--------------|---------|-----------|--------|------------|------------|------------|-------------|--------------|
| 0 | 1.33         | 64236.62     | 1036.74 | 54454.85  | 48.16  | 8696.87    | 8603.62    | 93.25      | 0.0         | conventional |
| 1 | 1.35         | 54876.98     | 674.28  | 44638.81  | 58.33  | 9505.56    | 9408.07    | 97.49      | 0.0         | conventional |
| 2 | 0.93         | 118220.22    | 794.70  | 109149.67 | 130.50 | 8145.35    | 8042.21    | 103.14     | 0.0         | conventional |

| AveragePrice | Total Volume | 4046     | 4225    | 4770     | Total Bags | Small Bags | Large Bags | XLarge Bags |     | type         |
|--------------|--------------|----------|---------|----------|------------|------------|------------|-------------|-----|--------------|
| 3            | 1.08         | 78992.15 | 1132.00 | 71976.41 | 72.58      | 5811.16    | 5677.40    | 133.76      | 0.0 | conventional |
| 4            | 1.28         | 51039.60 | 941.48  | 43838.39 | 75.78      | 6183.95    | 5986.26    | 197.69      | 0.0 | conventional |

In [80]:

```
plt.figure(figsize=(15,7),dpi=100)
sns.boxplot(data = dataset[[
    'AveragePrice',
    'Total Volume',
    '4046',
    '4225',
    '4770',
    'Total Bags',
    'Small Bags',
    'Large Bags',
    'XLarge Bags']]])
```

Out[80]: &lt;AxesSubplot:&gt;



Now clearly our data is free from outliers. Now we can fit our data to appropriate models.

## Step 4: Taking Care of the Categorical Variables

Now since our data contains categorical variables like "type", "month" and "region" we apply one-hot encoding to our variables "region", "month" and apply label encoding in variable "type".

One hot encoding creates equal number of columns, with 1's and 0's, as the number of categories in a categorical variable a column for a specific category contains 1's where the category is present and 0's elsewhere.

**As for label encoding it assigns numerical value to the categories of a categorical variable in their alphabetical order, the indexing starts with 0.**

**OneHotEncoder in Python can encode a specific number of categories since for the variable 'region' we have crossed that threshold we have used pandas.get\_dummies instead. Had we use OneHotEncoder we would have eliminated one column to avoid dummy variable trap but here we have no use for that.**

In [81]:

```
dataset['region'] = pd.Categorical(dataset['region'])
dfDummies_region = pd.get_dummies(dataset['region'], prefix = 'region')
dfDummies_region
```

Out[81]:

|              | region_Albany | region_Atlanta | region_BaltimoreWashington | region_Boise | region_Boston | region |
|--------------|---------------|----------------|----------------------------|--------------|---------------|--------|
| <b>0</b>     | 1             | 0              |                            | 0            | 0             | 0      |
| <b>1</b>     | 1             | 0              |                            | 0            | 0             | 0      |
| <b>2</b>     | 1             | 0              |                            | 0            | 0             | 0      |
| <b>3</b>     | 1             | 0              |                            | 0            | 0             | 0      |
| <b>4</b>     | 1             | 0              |                            | 0            | 0             | 0      |
| ...          | ...           | ...            |                            | ...          | ...           | ...    |
| <b>18244</b> | 0             | 0              |                            | 0            | 0             | 0      |
| <b>18245</b> | 0             | 0              |                            | 0            | 0             | 0      |
| <b>18246</b> | 0             | 0              |                            | 0            | 0             | 0      |
| <b>18247</b> | 0             | 0              |                            | 0            | 0             | 0      |
| <b>18248</b> | 0             | 0              |                            | 0            | 0             | 0      |

18249 rows × 54 columns

In [82]:

```
dataset = pd.concat([dataset, dfDummies_region], axis=1)
dataset.drop(columns="region", inplace=True)
dataset
```

Out[82]:

|              | AveragePrice | Total Volume | 4046    | 4225      | 4770   | Total Bags | Small Bags | Large Bags | XLarge Bags |         |
|--------------|--------------|--------------|---------|-----------|--------|------------|------------|------------|-------------|---------|
| <b>0</b>     | 1.33         | 64236.62     | 1036.74 | 54454.85  | 48.16  | 8696.87    | 8603.62    | 93.25      | 0.0         | conveni |
| <b>1</b>     | 1.35         | 54876.98     | 674.28  | 44638.81  | 58.33  | 9505.56    | 9408.07    | 97.49      | 0.0         | conveni |
| <b>2</b>     | 0.93         | 118220.22    | 794.70  | 109149.67 | 130.50 | 8145.35    | 8042.21    | 103.14     | 0.0         | conveni |
| <b>3</b>     | 1.08         | 78992.15     | 1132.00 | 71976.41  | 72.58  | 5811.16    | 5677.40    | 133.76     | 0.0         | conveni |
| <b>4</b>     | 1.28         | 51039.60     | 941.48  | 43838.39  | 75.78  | 6183.95    | 5986.26    | 197.69     | 0.0         | conveni |
| ...          | ...          | ...          | ...     | ...       | ...    | ...        | ...        | ...        | ...         | ...     |
| <b>18244</b> | 1.63         | 17074.83     | 2046.96 | 1529.20   | 0.00   | 13498.67   | 13066.82   | 431.85     | 0.0         | or      |

|       | AveragePrice | Total Volume | 4046    | 4225    | 4770   | Total Bags | Small Bags | Large Bags | XLarge Bags |    |
|-------|--------------|--------------|---------|---------|--------|------------|------------|------------|-------------|----|
| 18245 | 1.71         | 13888.04     | 1191.70 | 3431.50 | 0.00   | 9264.84    | 8940.04    | 324.80     | 0.0         | or |
| 18246 | 1.87         | 13766.76     | 1191.92 | 2452.79 | 727.94 | 9394.11    | 9351.80    | 42.31      | 0.0         | or |
| 18247 | 1.93         | 16205.22     | 1527.63 | 2981.04 | 727.01 | 10969.54   | 10919.54   | 50.00      | 0.0         | or |
| 18248 | 1.62         | 17489.58     | 2894.77 | 2356.13 | 224.53 | 12014.15   | 11988.14   | 26.01      | 0.0         | or |

18249 rows × 67 columns

Adding the one hot encoded columns for region into our data and dropping the region column from our dataset.

In [83]:

```
dataset['Month'] = pd.Categorical(dataset['Month'])
dfDummies_month = pd.get_dummies(dataset['Month'], prefix = 'month')
dfDummies_month
```

Out[83]:

|       | month_1 | month_2 | month_3 | month_4 | month_5 | month_6 | month_7 | month_8 | month_9 | more |
|-------|---------|---------|---------|---------|---------|---------|---------|---------|---------|------|
| 0     | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| 1     | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| 2     | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| 3     | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| 4     | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| ...   | ...     | ...     | ...     | ...     | ...     | ...     | ...     | ...     | ...     | ...  |
| 18244 | 0       | 1       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| 18245 | 1       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| 18246 | 1       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| 18247 | 1       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |
| 18248 | 1       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0       | 0    |

18249 rows × 12 columns

Similarly applying one hot encoding on months.

In [84]:

```
dataset = pd.concat([dataset, dfDummies_month], axis=1)
dataset.drop(columns="Month", inplace=True)
dataset
```

Out[84]:

|   | AveragePrice | Total Volume | 4046    | 4225     | 4770  | Total Bags | Small Bags | Large Bags | XLarge Bags |         |
|---|--------------|--------------|---------|----------|-------|------------|------------|------------|-------------|---------|
| 0 | 1.33         | 64236.62     | 1036.74 | 54454.85 | 48.16 | 8696.87    | 8603.62    | 93.25      | 0.0         | conveni |

|       | AveragePrice | Total Volume | 4046    | 4225      | 4770   | Total Bags | Small Bags | Large Bags | XLarge Bags |              |
|-------|--------------|--------------|---------|-----------|--------|------------|------------|------------|-------------|--------------|
| 1     | 1.35         | 54876.98     | 674.28  | 44638.81  | 58.33  | 9505.56    | 9408.07    | 97.49      | 0.0         | conventional |
| 2     | 0.93         | 118220.22    | 794.70  | 109149.67 | 130.50 | 8145.35    | 8042.21    | 103.14     | 0.0         | conventional |
| 3     | 1.08         | 78992.15     | 1132.00 | 71976.41  | 72.58  | 5811.16    | 5677.40    | 133.76     | 0.0         | conventional |
| 4     | 1.28         | 51039.60     | 941.48  | 43838.39  | 75.78  | 6183.95    | 5986.26    | 197.69     | 0.0         | conventional |
| ...   | ...          | ...          | ...     | ...       | ...    | ...        | ...        | ...        | ...         | ...          |
| 18244 | 1.63         | 17074.83     | 2046.96 | 1529.20   | 0.00   | 13498.67   | 13066.82   | 431.85     | 0.0         | organic      |
| 18245 | 1.71         | 13888.04     | 1191.70 | 3431.50   | 0.00   | 9264.84    | 8940.04    | 324.80     | 0.0         | organic      |
| 18246 | 1.87         | 13766.76     | 1191.92 | 2452.79   | 727.94 | 9394.11    | 9351.80    | 42.31      | 0.0         | organic      |
| 18247 | 1.93         | 16205.22     | 1527.63 | 2981.04   | 727.01 | 10969.54   | 10919.54   | 50.00      | 0.0         | organic      |
| 18248 | 1.62         | 17489.58     | 2894.77 | 2356.13   | 224.53 | 12014.15   | 11988.14   | 26.01      | 0.0         | organic      |

18249 rows × 78 columns

Adding the one hot encoded columns for Month into our data and dropping the Month column from our dataset.

In [85]:

```
from sklearn import preprocessing

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

Out[85]:

|       | AveragePrice | Total Volume | 4046    | 4225      | 4770   | Total Bags | Small Bags | Large Bags | XLarge Bags | type |
|-------|--------------|--------------|---------|-----------|--------|------------|------------|------------|-------------|------|
| 0     | 1.33         | 64236.62     | 1036.74 | 54454.85  | 48.16  | 8696.87    | 8603.62    | 93.25      | 0.0         | 0    |
| 1     | 1.35         | 54876.98     | 674.28  | 44638.81  | 58.33  | 9505.56    | 9408.07    | 97.49      | 0.0         | 0    |
| 2     | 0.93         | 118220.22    | 794.70  | 109149.67 | 130.50 | 8145.35    | 8042.21    | 103.14     | 0.0         | 0    |
| 3     | 1.08         | 78992.15     | 1132.00 | 71976.41  | 72.58  | 5811.16    | 5677.40    | 133.76     | 0.0         | 0    |
| 4     | 1.28         | 51039.60     | 941.48  | 43838.39  | 75.78  | 6183.95    | 5986.26    | 197.69     | 0.0         | 0    |
| ...   | ...          | ...          | ...     | ...       | ...    | ...        | ...        | ...        | ...         | ...  |
| 18244 | 1.63         | 17074.83     | 2046.96 | 1529.20   | 0.00   | 13498.67   | 13066.82   | 431.85     | 0.0         | 1    |
| 18245 | 1.71         | 13888.04     | 1191.70 | 3431.50   | 0.00   | 9264.84    | 8940.04    | 324.80     | 0.0         | 1    |
| 18246 | 1.87         | 13766.76     | 1191.92 | 2452.79   | 727.94 | 9394.11    | 9351.80    | 42.31      | 0.0         | 1    |
| 18247 | 1.93         | 16205.22     | 1527.63 | 2981.04   | 727.01 | 10969.54   | 10919.54   | 50.00      | 0.0         | 1    |
| 18248 | 1.62         | 17489.58     | 2894.77 | 2356.13   | 224.53 | 12014.15   | 11988.14   | 26.01      | 0.0         | 1    |

18249 rows × 78 columns

**Now label encoding on the variable "type"**

**Hence our preprocessing ends here!!!**

**Now its time that we fit multiple linear regression, decision tree regression and random forest regression onto our data.**

## Step 5: Model Fitting

In [86]:

```
dataset.head()
```

Out[86]:

|   | AveragePrice | Total Volume | 4046    | 4225      | 4770   | Total Bags | Small Bags | Large Bags | XLarge Bags | type | ... | mo  |
|---|--------------|--------------|---------|-----------|--------|------------|------------|------------|-------------|------|-----|-----|
| 0 | 1.33         | 64236.62     | 1036.74 | 54454.85  | 48.16  | 8696.87    | 8603.62    | 93.25      | 0.0         | 0    | 0   | ... |
| 1 | 1.35         | 54876.98     | 674.28  | 44638.81  | 58.33  | 9505.56    | 9408.07    | 97.49      | 0.0         | 0    | 0   | ... |
| 2 | 0.93         | 118220.22    | 794.70  | 109149.67 | 130.50 | 8145.35    | 8042.21    | 103.14     | 0.0         | 0    | 0   | ... |
| 3 | 1.08         | 78992.15     | 1132.00 | 71976.41  | 72.58  | 5811.16    | 5677.40    | 133.76     | 0.0         | 0    | 0   | ... |
| 4 | 1.28         | 51039.60     | 941.48  | 43838.39  | 75.78  | 6183.95    | 5986.26    | 197.69     | 0.0         | 0    | 0   | ... |

5 rows × 78 columns

**Having a look at our data after complete preprocessing.**

**Splitting our dataset to training and test numpy arrays with the names having their intended meaning. Where we are using 80% of our dataset for training and 20% of the data for testing.**

In [87]:

```
X=dataset.iloc[:,1:78]
y=dataset['AveragePrice']
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=50)
y_test = np.array(y_test,dtype = float)
```

**Normalizing our X\_train and X\_test using standard scaler**

In [88]:

```
from sklearn.preprocessing import StandardScaler
sc=StandardScaler()
X_train=sc.fit_transform(X_train)
X_test=sc.transform(X_test)
```

**The function regression\_results defined below calculates and prints the following features of a model: explained\_variance, r2, adjusted\_r2, MAE, MSE, RMSE. It accepts the original and predicted values as its arguments.**

In [89]:

```
import sklearn.metrics as metrics

def regression_results(y_true, y_pred):
    explained_variance=metrics.explained_variance_score(y_true, y_pred)
```

```

mean_absolute_error=metrics.mean_absolute_error(y_true, y_pred)
mse=metrics.mean_squared_error(y_true, y_pred)
r2=metrics.r2_score(y_true, y_pred)
adjusted_r2 = 1 - (1-r2)*(len(y_true)-1)/(len(y_true)-X_test.shape[1]-1)

print('Explained_variance: ', round(explained_variance,4))
print('R2: ', round(r2,4))
print('Adjusted_r2: ', round(adjusted_r2,4))
print('MAE: ', round(mean_absolute_error,4))
print('MSE: ', round(mse,4))
print('RMSE: ', round(np.sqrt(mse),4))

```

**Below is a function to find the accuracy of each model on the basis of K-fold cross validation.**

In [90]:

```

from sklearn.model_selection import cross_val_score
def model_accuracy(model,X_train=X_train,y_train=y_train):
    accuracies = cross_val_score(estimator = model, X = X_train, y = y_train, cv = 10)
    print("Accuracy: {:.2f} %".format(accuracies.mean()*100))
    print("Standard Deviation: {:.2f} %".format(accuracies.std()*100))

```

## Fitting Multiple Linear Regression Model

The following code snippet fits the multiple linear regression model on `X_train` and `y_train` and predicts the values for `X_test` and stores it in `y_pred`. It also prints the outputs of the functions defined above. Hence giving us a useful summary for the multiple linear regression model.

In [91]:

```

from sklearn.linear_model import LinearRegression
import statsmodels.api as sm

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:  0.6665
R2:  0.6664
Adjusted_r2:  0.6592
MAE:  0.1779
MSE:  0.0541
RMSE:  0.2327
Accuracy: 64.15 %
Standard Deviation: 1.48 %

```

In [99]:

```

from yellowbrick.regressor import ResidualsPlot

fig = plt.figure(figsize=(16, 12),dpi=100)
visualizer = ResidualsPlot(regressor, hist = True, qqplot = False)
visualizer.score(X_test, y_test) # Evaluate the model on the test data
visualizer.show() # Finalize and render the figure

```



```
Out[99]: <AxesSubplot:title={'center':'Residuals for LinearRegression Model'}, xlabel='Predicted Value', ylabel='Residuals'>
```

## Fitting Random Forest Regression Model

**The following code snippet fits the random forest regression model on X\_train and y\_train and predicts the values for X\_test and stores it in y\_pred\_rf. It also prints the outputs of the functions defined above. Hence giving us a useful summary for the random forest regression model.**

```
In [93]: 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: 0.9059

R2: 0.9058

Adjusted\_r2: 0.9038

MAE: 0.089

MSE: 0.0153

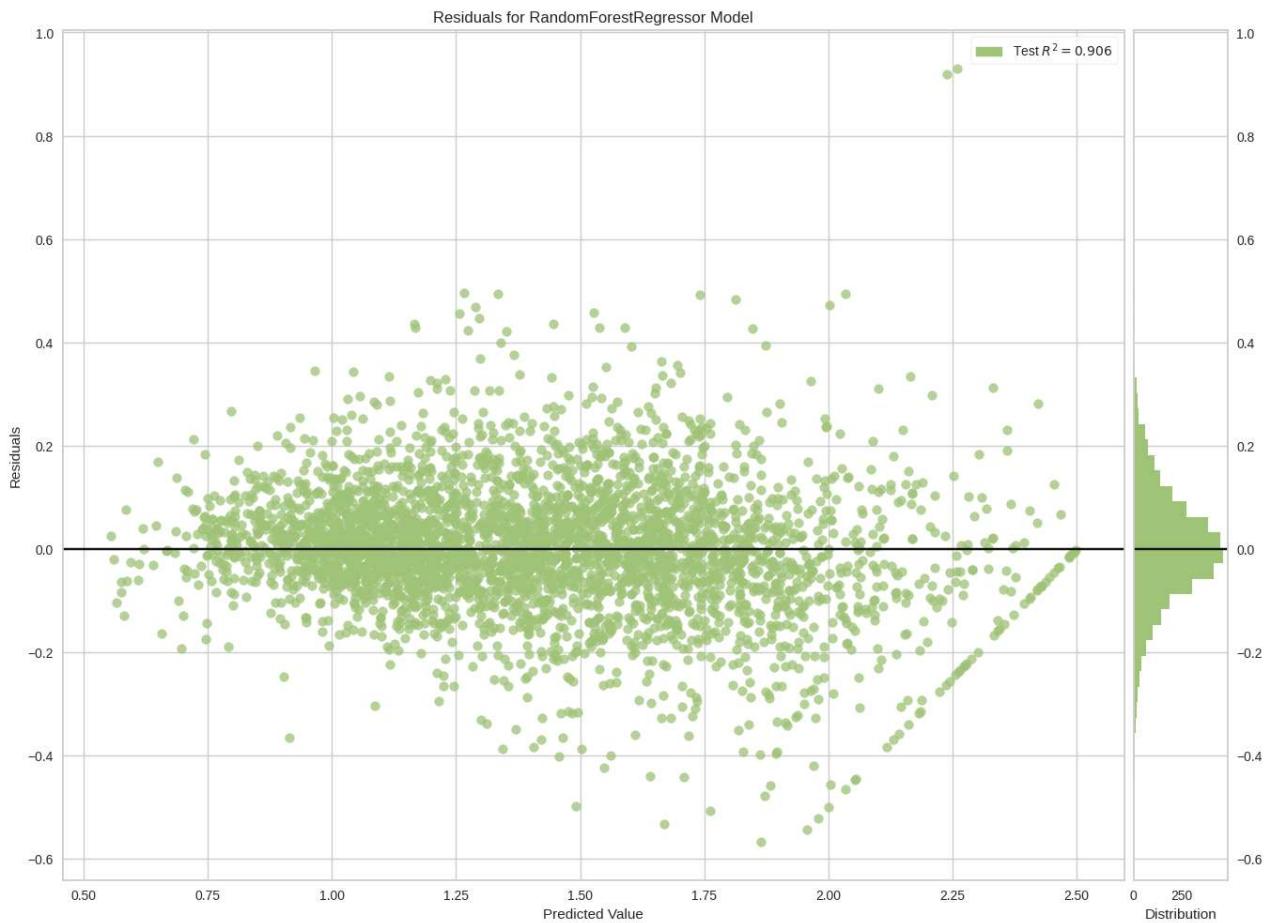
RMSE: 0.1236

Accuracy: 88.53 %

Standard Deviation: 0.52 %

```
In [100...]: fig = plt.figure(figsize=(16, 12),dpi=100)
visualizer = ResidualsPlot(rand_regressor, hist = True, qqplot = False)
```

```
visualizer.score(X_test, y_test) # Evaluate the model on the test data
visualizer.show() # Finalize and render the figure
```



```
Out[100... <AxesSubplot:title={'center':'Residuals for RandomForestRegressor Model'}, xlabel='Predicted Value', ylabel='Residuals'>
```

## Fitting Decision Tree Regression Model

**The following code snippet fits the decision tree regression model on X\_train and y\_train and predicts the values for X\_test and stores it in y\_pred\_dt. It also prints the outputs of the functions defined above. Hence giving us a useful summary for the decision tree regression model.**

In [96]:

```
from sklearn.tree import DecisionTreeRegressor

decision_tree=DecisionTreeRegressor(criterion='mse', splitter='random', random_state=10)
decision_tree.fit(X_train, y_train)
y_pred_dt = decision_tree.predict(X_test)
regression_results(y_test,y_pred_dt)
model_accuracy(decision_tree)
```

Explained\_variance: 0.8343

R2: 0.834

Adjusted\_r2: 0.8304

MAE: 0.1089

MSE: 0.0269

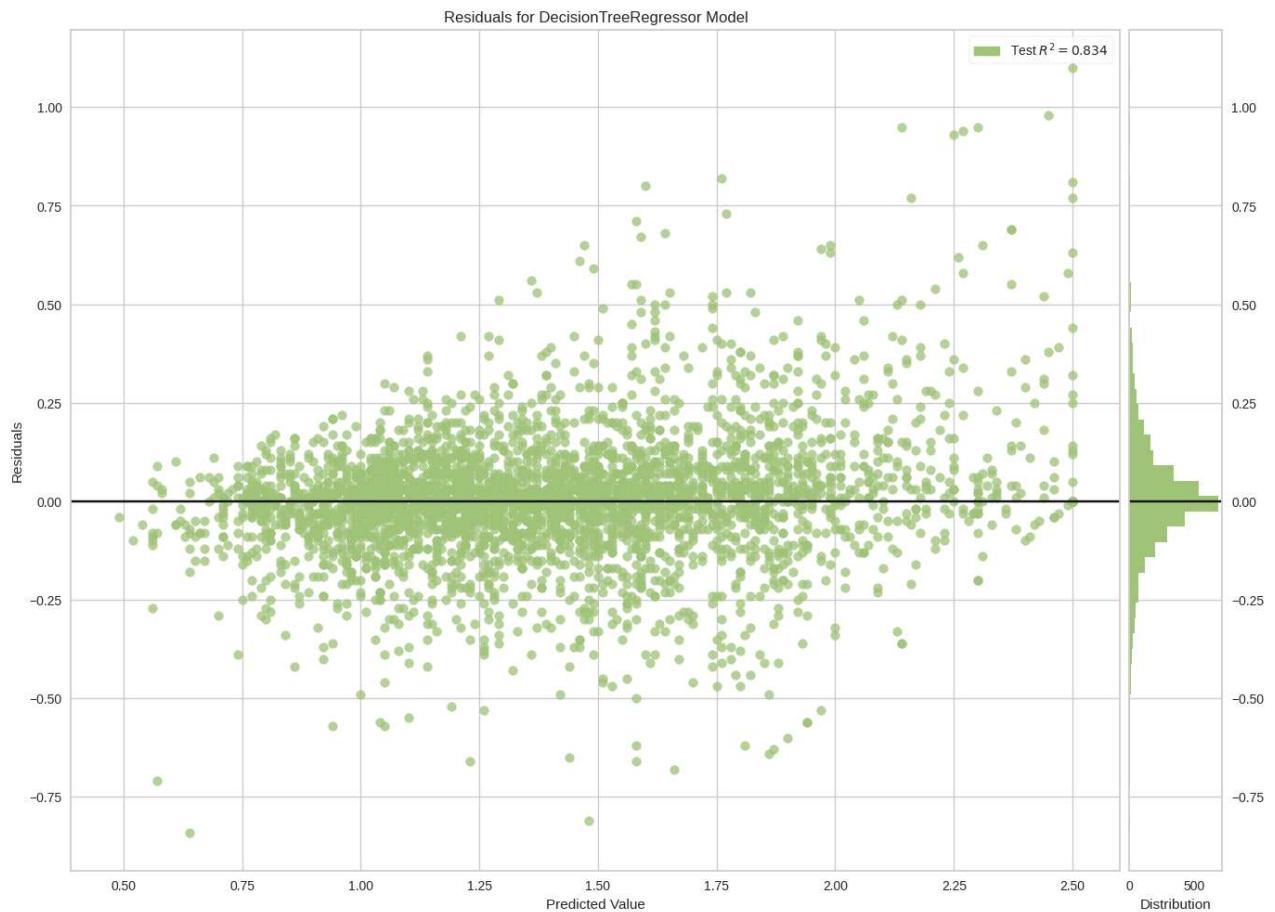
RMSE: 0.1642

Accuracy: 80.41 %

Standard Deviation: 1.71 %

In [101...]

```
fig = plt.figure(figsize=(16, 12), dpi=100)
visualizer = ResidualsPlot(decision_tree, hist = True, qqplot = False)
visualizer.score(X_test, y_test) # Evaluate the model on the test data
visualizer.show() # Finalize and render the figure
```



Out[101...]: &lt;AxesSubplot:title={'center':'Residuals for DecisionTreeRegressor Model'}, xlabel='Predicted Value', ylabel='Residuals'&gt;

**As our conclusion we proclaim that, using k-fold cross validation as the basis for model selection we declare random forest model as the best suited model for our purpose of predicting average avocado prices**