

```
In [3]: import pandas as pd
import matplotlib
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
matplotlib inline

import plotly.offline as py
import plotly.graph_objs as go
from plotly.offline import init_notebook_mode
init_notebook_mode(connected=True)
from plotly import tools

import warnings
warnings.filterwarnings("ignore")
warnings.filterwarnings("ignore",category=DeprecationWarning)

In [4]: df= pd.read_csv("https://raw.githubusercontent.com/insalid2018/Term-2/master/Projects/avocado.csv")

In [5]: df.shape
Out[5]: (18249, 14)

In [6]: df.columns
Out[6]: Index([ 'Unnamed: 0', 'Date', 'AveragePrice', 'Total Volume', '4046', '4225',
          '4770', 'Total Bags', 'Small Bags', 'Large Bags', 'XLarge Bags', 'type',
          'year', 'region'],
          dtype='object')

In [7]: df.head()
Out[7]:
   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  54454.85  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      2  2015-12-13      0.93    11820.22    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 [8]: df.tail()
Out[8]:
   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    2046.96  1529.20      0.00    13498.67    13066.82    431.85      0.0  organic    2018  WestTexNewMexico
18245      8  2018-01-28      1.71    13888.04    1191.70  3431.50      0.00    9264.84    8940.04    324.80      0.0  organic    2018  WestTexNewMexico
18246      9  2018-01-21      1.87    13766.76    1191.92  2452.79  727.94    9394.11    9351.80    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 [9]: df.info()
<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   AveragePrice  18249 non-null    float64
3   Total Volume  18249 non-null    float64
4   4046         18249 non-null    float64
5   4225         18249 non-null    float64
6   4770         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(10), int64(2), object(3)
memory usage: 1.5+ MB

In [10]: df.describe()
Out[10]:
   Unnamed: 0  AveragePrice  Total Volume  4046  4225  4770  Total Bags  Small Bags  Large Bags  XLarge Bags  type  year
count  18249.000000  18249.000000  18249000e+00  18249000e+04  1.824900e+04  1.824900e+04  18249.000000  18249.000000
mean    24.232232    1.405877  8.506440e+05  2.930084e+05  2.951548e+05  2.283974e+04  2.396392e+05  2.439660e+05  17692.894652  0.939938
std    15.481045    0.402677  3.453345e+06  1.264989e+06  1.204120e+06  1.074641e+05  9.862424e+05  7.461785e+05  2.439660e+05  0.939938
min      0.000000    0.440000  8.456000e+01  0.000000e+00  0.000000e+00  0.000000e+00  0.000000e+00  0.000000e+00  0.000000e+00  0.000000
25%    10.000000    1.100000  1.083858e+01  8.540700e+02  3.008780e+03  5.088640e+03  2.849420e+03  1.274700e+02  0.000000  2015.000000
50%    24.000000    1.370000  1.073768e+05  8.645300e+03  2.906102e+04  1.849900e+02  3.974383e+04  2.636282e+04  2.647710e+03  0.000000
75%    38.000000    1.660000  4.329623e+05  1.110202e+05  1.502069e+05  6.243420e+03  1.107834e+05  8.333767e+04  2.202925e+04  132.500000
max    52.000000    3.250000  6.250565e+07  2.274362e+07  2.047057e+07  2.546439e+06  1.937313e+07  1.338459e+07  551693.650000  2018.000000

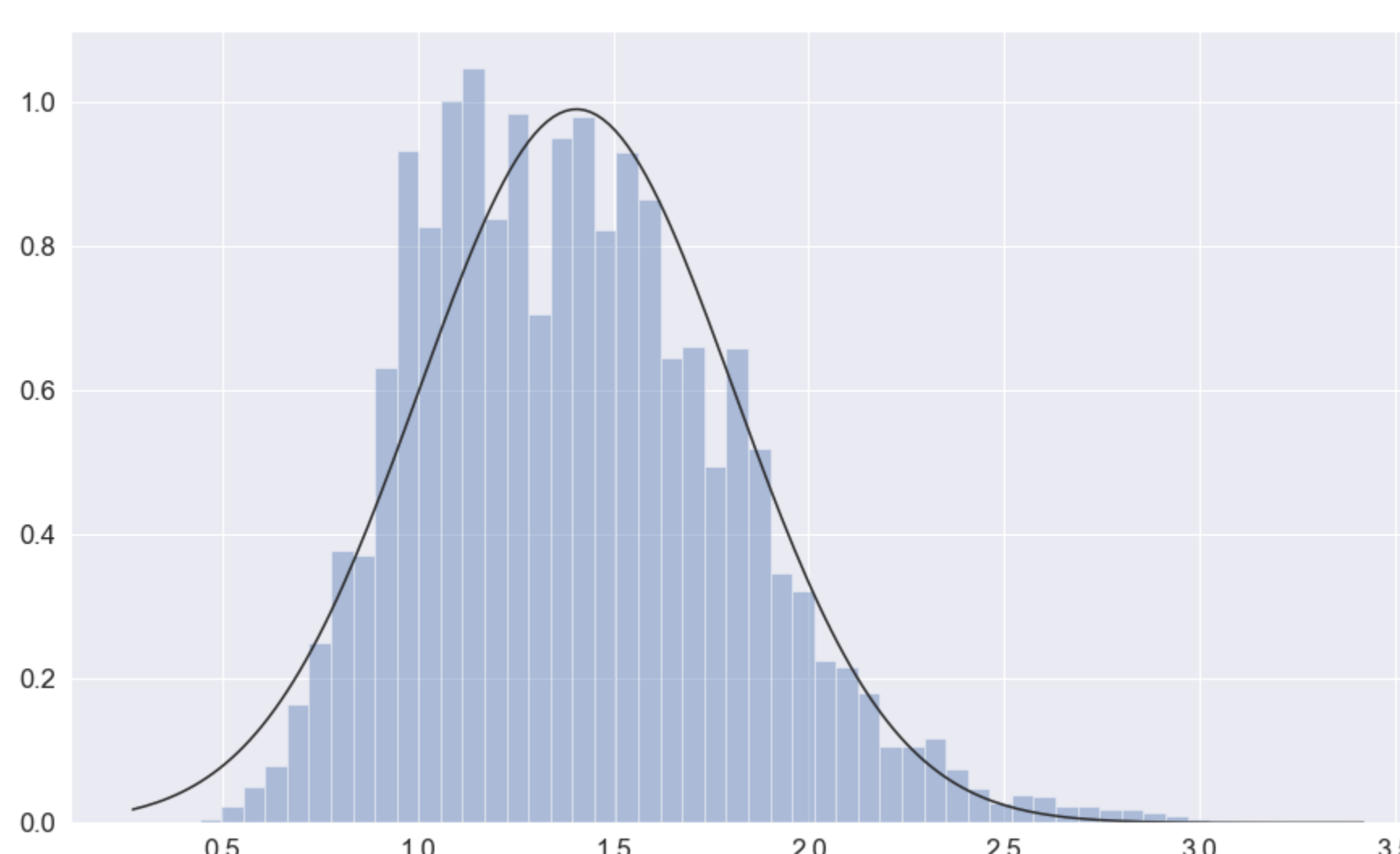
In [11]: df.isnull().sum()
Out[11]:
Unnamed: 0    0
Date          0
AveragePrice  0
Total Volume  0
4046          0
4225          0
4770          0
Total Bags    0
Small Bags    0
Large Bags    0
XLarge Bags   0
type          0
year          0
region        0
dtype: int64

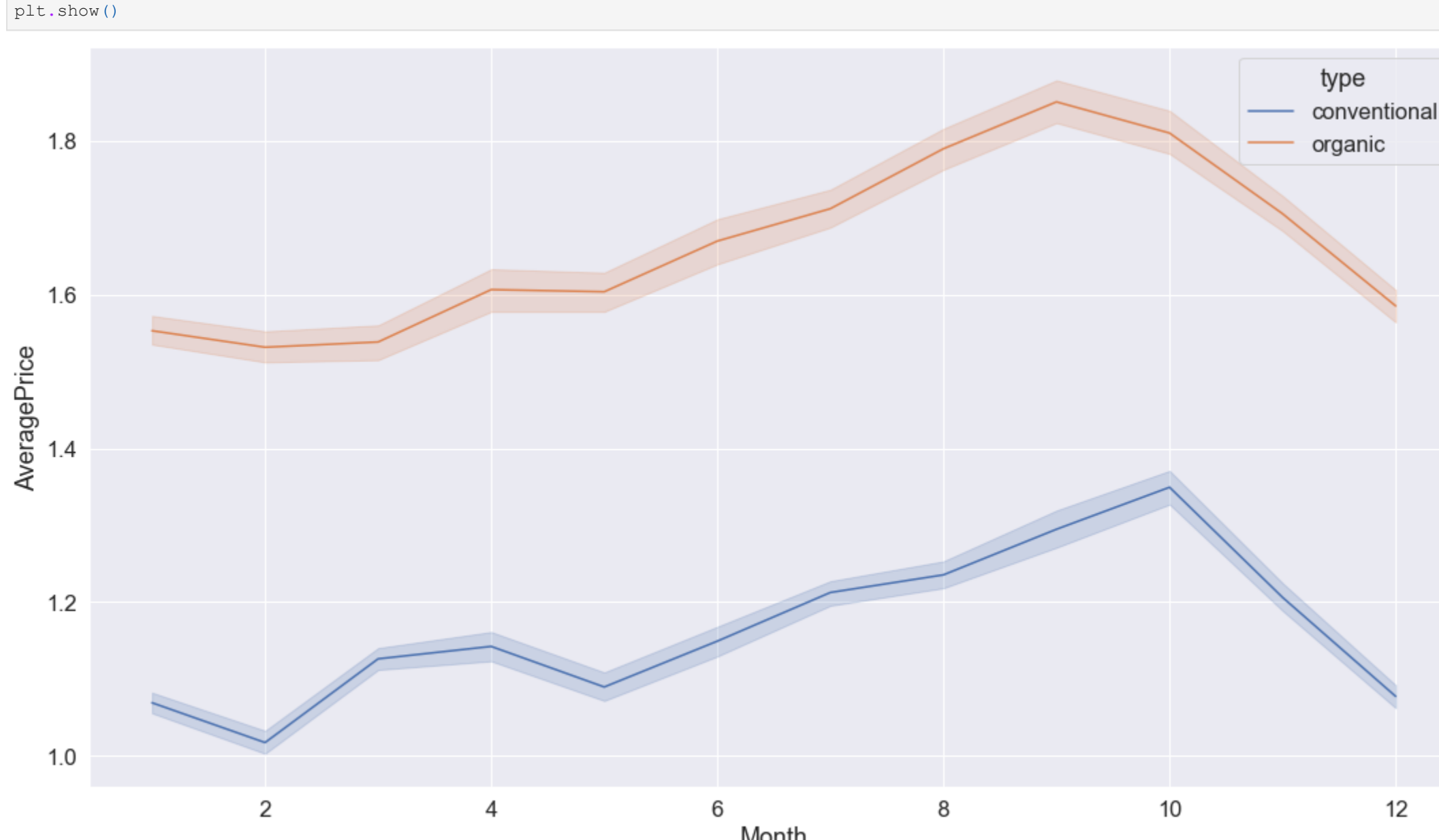
In [12]: df.drop('Unnamed: 0',axis=1,inplace=True)

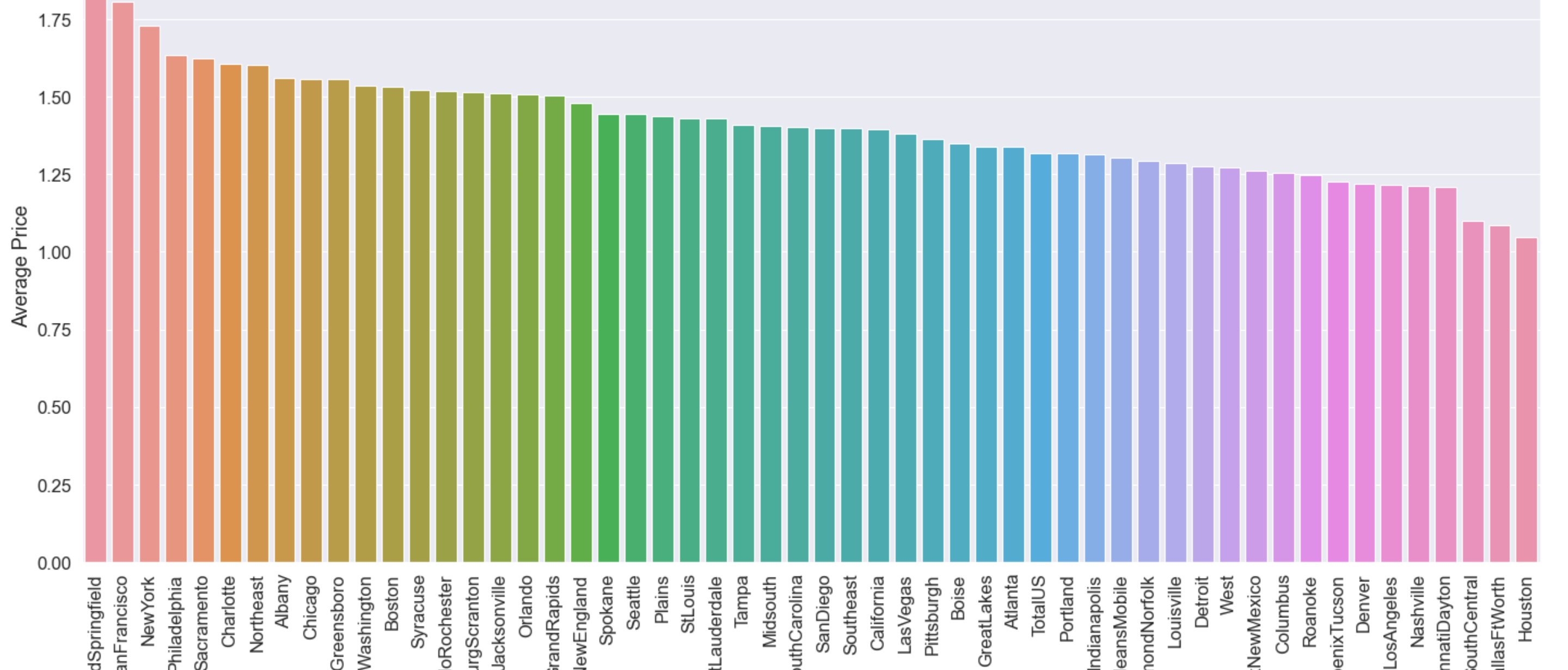
In [13]: df.head()
Out[13]:
   Date  AveragePrice  Total Volume  4046  4225  4770  Total Bags  Small Bags  Large Bags  XLarge Bags  type  year  region
0  2015-12-27      1.33    64236.62    1036.74  54454.85  48.16    8696.87    8603.62    93.25      0.0  conventional  2015  Albany
1  2015-12-20      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    11820.22    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 [14]: df['Date']=pd.to_datetime(df['Date'])
df['Month']=df['Date'].apply(lambda x:x.month)
df['Day']=df['Date'].apply(lambda x:x.day)

In [15]: df.head()
Out[15]:
   Date  AveragePrice  Total Volume  4046  4225  4770  Total Bags  Small Bags  Large Bags  XLarge Bags  type  year  region  Month  Day
0  2015-12-27      1.33    64236.62    1036.74  54454.85  48.16    8696.87    8603.62    93.25      0.0  conventional  2015  Albany    12    27
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    12    20
2  2015-12-13      0.93    11820.22    794.70    109149.67  130.50    8145.35    8042.21    103.14      0.0  conventional  2015  Albany    12    13
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    12     6
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    11    29

In [18]: sns.set(font_scale=1.5)
from scipy.stats import norm
fig, ax = plt.subplots(figsize=(15, 9))
sns.distplot(df['AveragePrice'], kde=False, fit=norm)
<AxesSubplot: xlabel='AveragePrice'>
Out[18]:
<AxesSubplot: xlabel='AveragePrice'>


In [19]: plt.figure(figsize=(18,10))
sns.lineplot(x='Month', y='AveragePrice', hue='type', data=df)
plt.show()


In [20]: region_list=list(df.region.unique())
average_price=[]
for i in region_list:
    x=df[df.region==i]
    region_average=sum(x.AveragePrice)/len(x)
    average_price.append(region_average)
df=pd.DataFrame({'region_list':region_list,'average_price':average_price})
new_index=df.average_price.sort_values(ascending=False).index.values
sorted_data=df.ix[new_index]
plt.figure(figsize=(24,10))
ax=sns.barplot(x=sorted_data.region_list,y=sorted_data.average_price)
plt.xticks(rotation=90)
plt.xlabel('Region')
plt.ylabel('Average Price')
plt.title('Average Price of Avocado According to Region')
Out[20]:
Text(0.5, 1.0, 'Average Price of Avocado According to Region')


In [22]: plt.figure(figsize=(25,20))
sns.heatmap(df.corr(),cmap='coolwarm',annot=True)
Out[22]:
<AxesSubplot:
AveragePrice
Total Volume
4046
4225
4770
Total Bags
Small Bags
Large Bags
XLarge Bags
year
Month
Day
AveragePrice
Total Volume
4046
4225
4770
Total Bags
Small Bags
Large Bags
XLarge Bags
year
Month
Day
1
-0.19
-0.21
-0.17
-0.18
-0.18
-0.07
-0.17
-0.12
0.093
0.16
0.027
-0.19
1
0.98
0.97
0.87
0.96
0.97
0.88
0.75
0.017
-0.025
-0.0097
-0.21
0.98
1
0.93
0.83
0.92
0.93
0.84
0.7
0.0034
-0.026
-0.01
-0.17
0.97
0.93
1
0.89
0.91
0.92
0.81
0.69
-0.0096
-0.022
-0.012
-0.17
0.87
0.83
0.89
1
0.79
0.8
0.7
0.68
-0.037
-0.033
-0.009
-0.18
0.96
0.92
0.91
0.79
1
0.99
0.94
0.8
0.072
-0.023
-0.005
-0.17
0.97
0.93
0.92
0.8
0.99
1
0.9
0.81
0.064
-0.023
-0.0039
-0.17
0.88
0.84
0.81
0.7
0.94
0.9
1
0.71
0.088
-0.02
-0.0084
-0.12
0.75
0.7
0.69
0.68
0.8
0.81
0.71
1
0.081
-0.013
0.00032
year
0.093
0.017
0.0034
-0.0096
-0.037
0.072
0.064
0.088
0.081
1
-0.18
0.0045
Month
0.16
-0.025
-0.026
-0.022
-0.033
-0.023
-0.023
-0.02
-0.013
-0.18
1
0.011
Day
0.027
-0.0097
-0.01
-0.012
-0.009
-0.005
-0.0039
-0.0084
0.00032
0.0045
0.011
1
AveragePrice
Total Volume
4046
4225
4770
Total Bags
Small Bags
Large Bags
XLarge Bags
year
Month
Day
1
-0.19
-0.21
-0.17
-0.18
-0.18
-0.07
-0.17
-0.12
0.093
0.16
0.027
-0.19
1
0.98
0.97
0.87
0.96
0.97
0.88
0.75
0.017
-0.025
-0.0097
-0.21
0.98
1
0.93
0.83
0.92
0.93
0.84
0.7
0.0034
-0.026
-0.01
-0.17
0.97
0.93
1
0.89
0.91
0.92
0.81
0.69
-0.0096
-0.022
-0.012
-0.17
0.87
0.83
0.89
1
0.79
0.8
0.7
0.68
-0.037
-0.033
-0.009
-0.18
0.96
0.92
0.91
0.79
1
0.99
0.94
0.8
0.072
-0.023
-0.005
-0.17
0.97
0.93
0.92
0.8
0.99
1
0.9
0.81
0.064
-0.023
-0.0039
-0.17
0.88
0.84
0.81
0.7
0.94
0.9
1
0.71
0.088
-0.02
-0.0084
-0.12
0.75
0.7
0.69
0.68
0.8
0.81
0.71
1
0.081
-0.013
0.00032
year
0.093
0.017
0.0034
-0.0096
-0.037
0.072
0.064
0.088
0.081
1
-0.18
0.0045
Month
0.16
-0.025
-0.026
-0.022
-0.033
-0.023
-0.023
-0.02
-0.013
-0.18
1
0.011
Day
0.027
-0.0097
-0.01
-0.012
-0.009
-0.005
-0.0039
-0.0084
0.00032
0.0045
0.011
1
```