

In [ ]:

# REGRESSION PROJECT\_Avocado Data Analysis

## Comparison of all Regression Models\_ Avocado Data Analysis

```
In [55]: #display image using python
from IPython.display import Image
url = 'E:\Data Science & AI\Dataset files\Fruit.jpg'
Image(url,height=200,width=900)
```

```
<>:3: SyntaxWarning: invalid escape sequence '\D'
<>:3: SyntaxWarning: invalid escape sequence '\D'
C:\Users\roy62\AppData\Local\Temp\ipykernel_13772\81040572.py:3: SyntaxWarning: i
nvalid escape sequence '\D'
url = 'E:\Data Science & AI\Dataset files\Fruit.jpg'
```

Out[55]:



```
In [57]: #importing libraries

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
import warnings
warnings.filterwarnings('ignore')
```

```
In [58]: #importing the dataset
data = pd.read_csv(r'E:\Data Science & AI\Dataset files\avocado.csv',index_col=0)
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 18249 entries, 0 to 11
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: 1.9+ MB
```

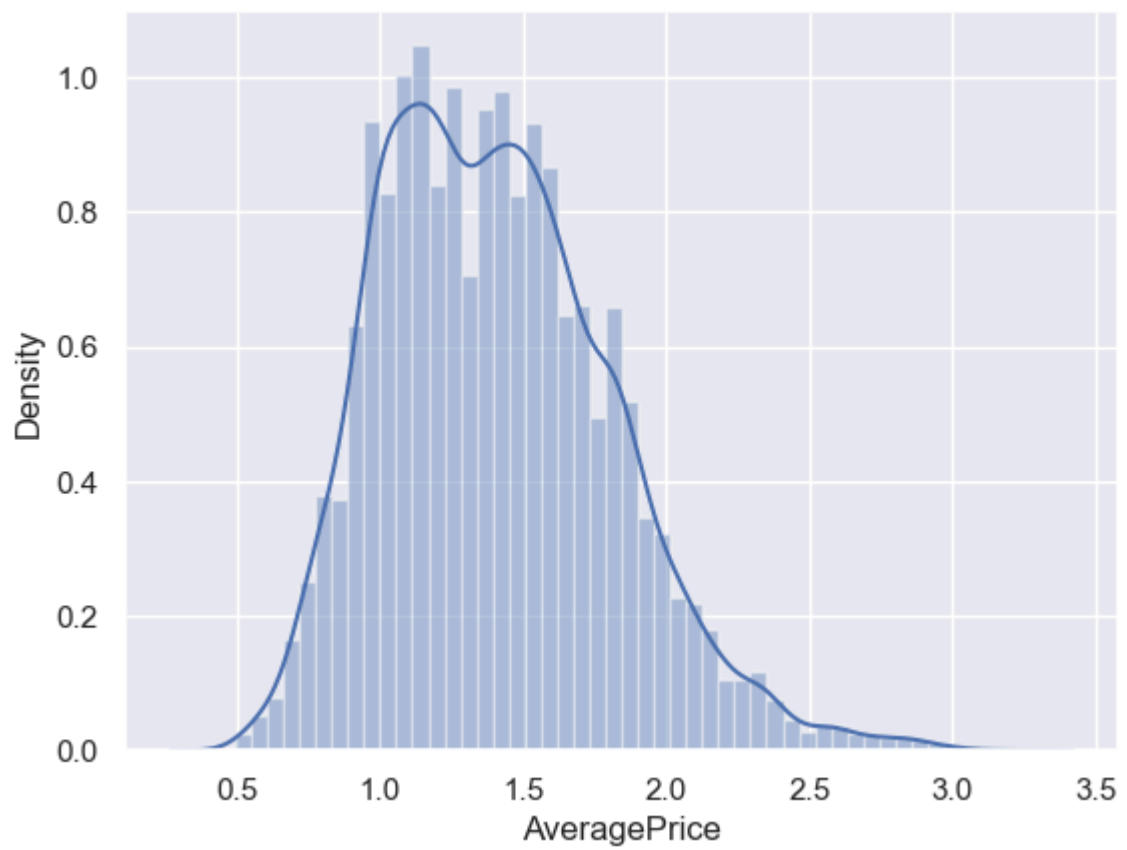
```
In [59]: data.head(3)
```

```
Out[59]:
```

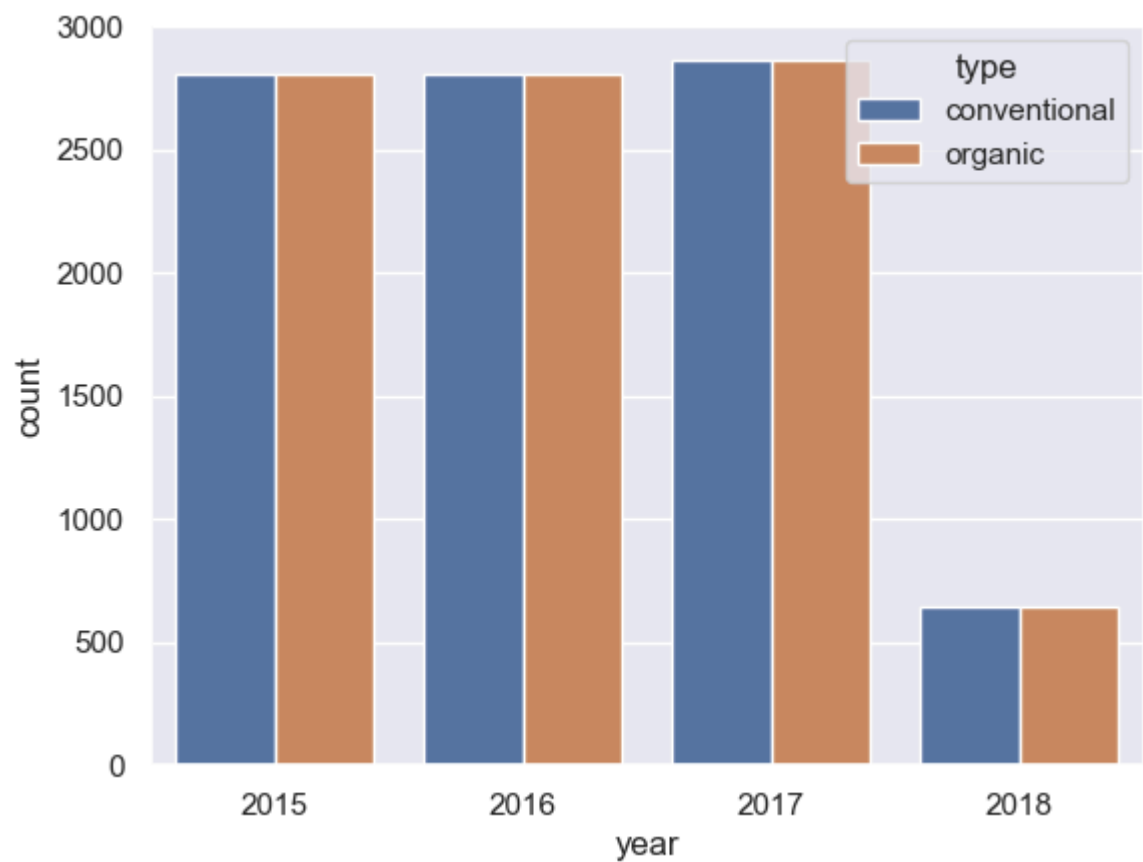
	Date	AveragePrice	Total Volume	4046	4225	4770	Total Bags	Small Bags	Large Bags
0	2015-12-27	1.33	64236.62	1036.74	54454.85	48.16	8696.87	8603.62	93.25
1	2015-12-20	1.35	54876.98	674.28	44638.81	58.33	9505.56	9408.07	97.49
2	2015-12-13	0.93	118220.22	794.70	109149.67	130.50	8145.35	8042.21	103.14



```
In [60]: sns.distplot(data['AveragePrice']);
```



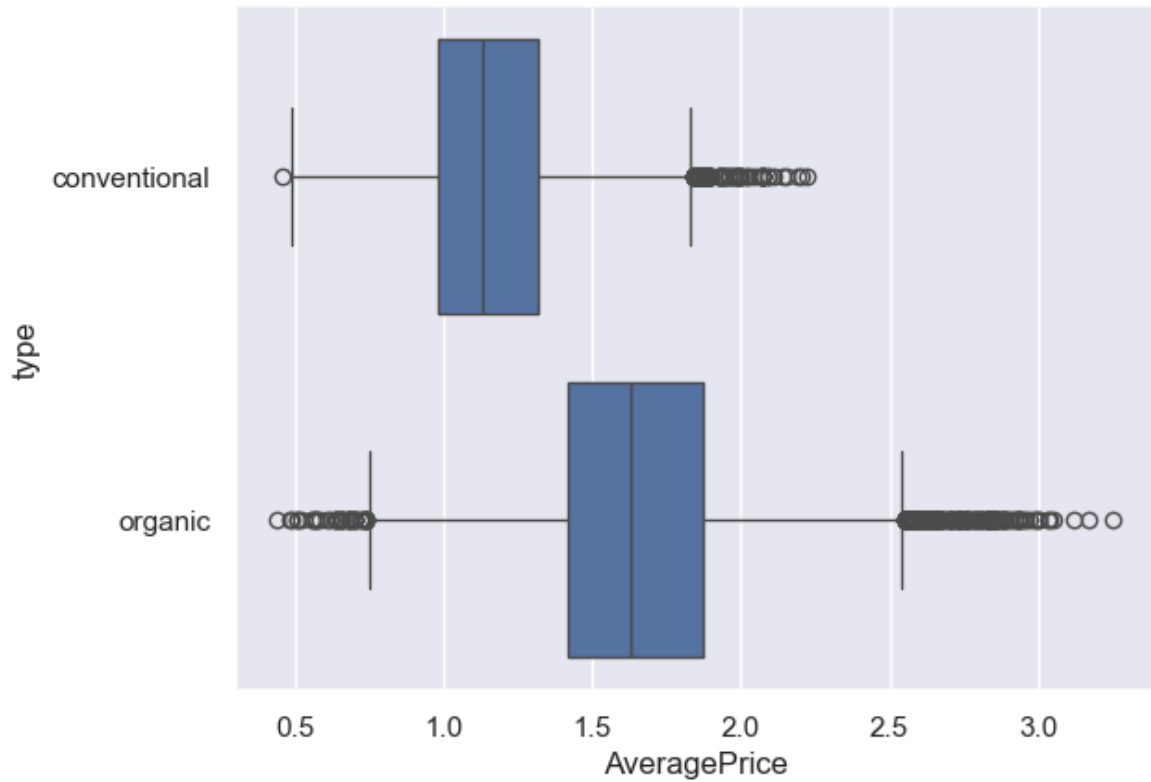
```
In [61]: sns.countplot(x='year',data=data,hue='type');
```



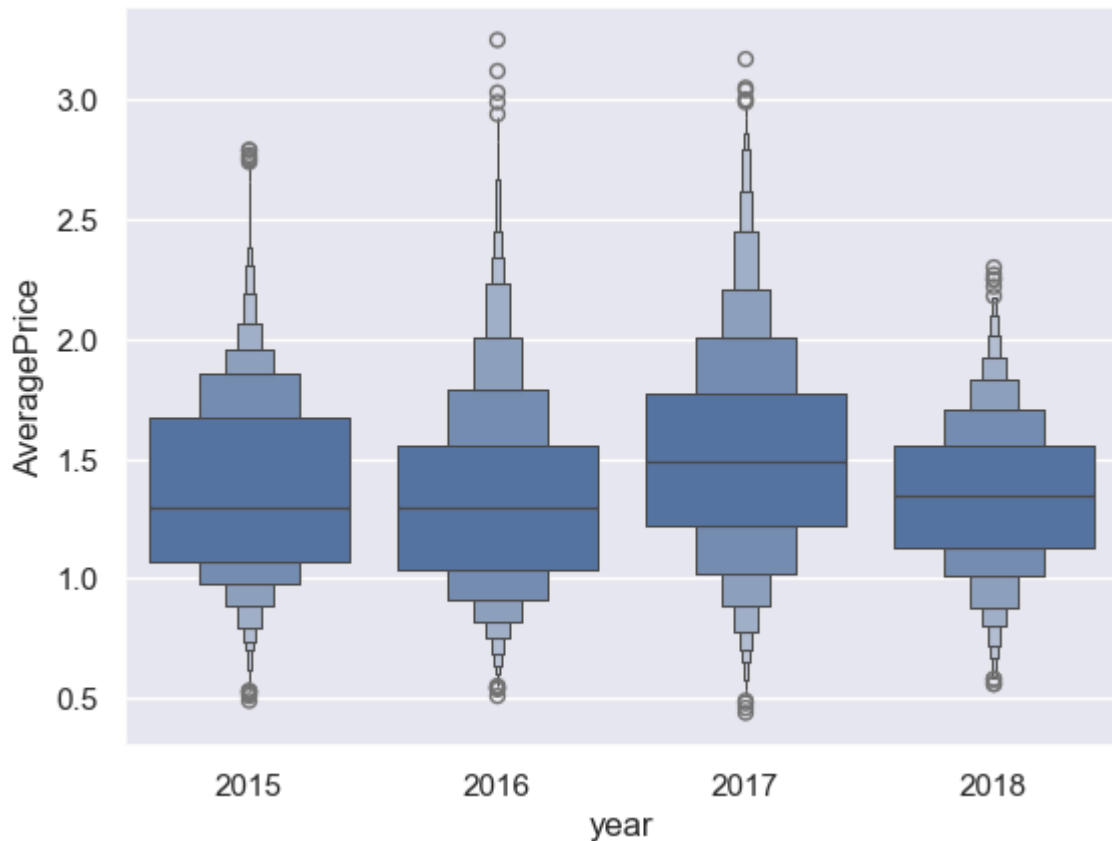
```
In [62]: data.year.value_counts()
```

```
Out[62]: year
2017    5722
2016    5616
2015    5615
2018    1296
Name: count, dtype: int64
```

```
In [63]: sns.boxplot(y="type", x="AveragePrice", data=data);
```



```
In [64]: data.year=data.year.apply(str)
sns.boxenplot(x="year", y="AveragePrice", data=data);
```



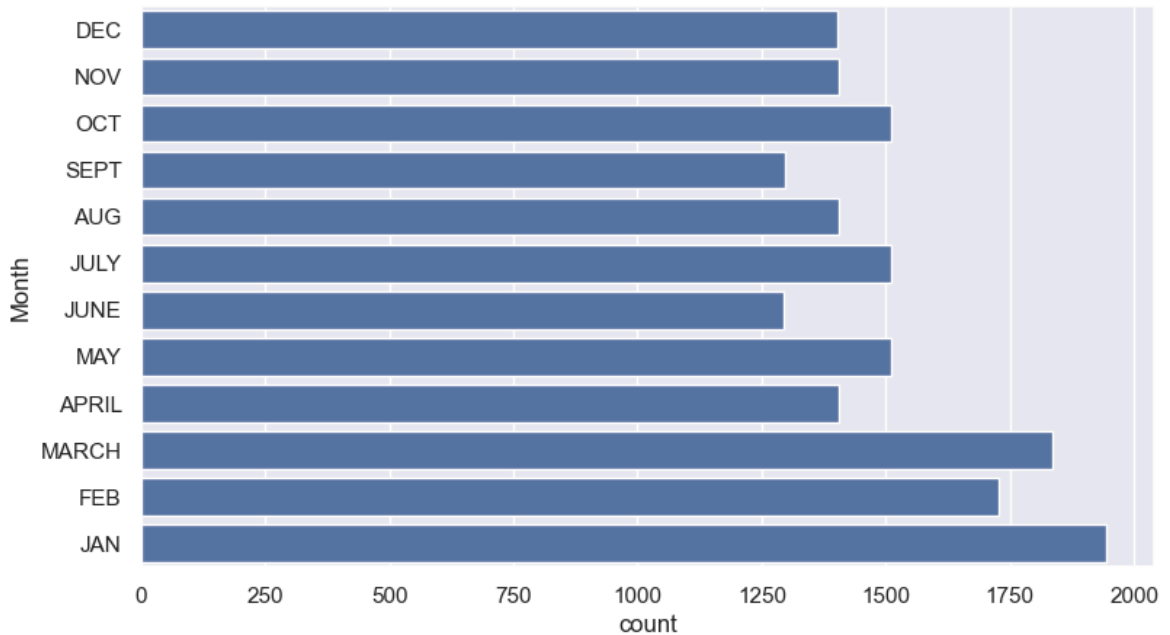
```
In [65]: # Dealing with categorical features

data['type'] = data['type'].map({'conventional':0, 'organic':1})

# Extracting month from date column.
data.Date = data.Date.apply(pd.to_datetime)
data['Month'] = data['Date'].apply(lambda x: x.month)
data.drop('Date', axis=1, inplace=True)
data.Month = data.Month.map({1:'JAN', 2:'FEB', 3:'MARCH', 4:'APRIL', 5:'MAY', 6:'JUNE'})
```

```
In [66]: plt.figure(figsize=(9,5))
sns.countplot(data['Month'])
plt.title('Monthwise Distribution of Sales', fontdict={'fontsize':25});
```

## Monthwise Distribution of Sales



```
In [70]: # Preparing data for ML models

# Creating dummy variables
dummies = pd.get_dummies(data[['year','region','Month']],drop_first=True)
df_dummies = pd.concat([data[['Total Volume', '4046', '4225', '4770', 'Total Bags',
                             'Small Bags', 'Large Bags', 'XLarge Bags', 'type']],dummies],axis=1)
target = data['AveragePrice']

# Splitting data into training and test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(df_dummies,target,test_size=

# Standardizing the data
cols_to_std = ['Total Volume', '4046', '4225', '4770', 'Total Bags', 'Small Bags']
from sklearn.preprocessing import StandardScaler
scaler=StandardScaler()
scaler.fit(X_train[cols_to_std])
X_train[cols_to_std] = scaler.transform(X_train[cols_to_std])
X_test[cols_to_std] = scaler.transform(X_test[cols_to_std])
```

```
In [72]: #importing ML models from scikit-learn
from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.svm import SVR
from sklearn.neighbors import KNeighborsRegressor
from sklearn.metrics import mean_absolute_error,mean_squared_error,r2_score
```

```
In [73]: #to save time all models can be applied once using for loop
regressors = {
    'Linear Regression' : LinearRegression(),
    'Decision Tree' : DecisionTreeRegressor(),
    'Random Forest' : RandomForestRegressor(),
    'Support Vector Machines' : SVR(gamma=1),
    'K-nearest Neighbors' : KNeighborsRegressor(n_neighbors=1),
}
results=pd.DataFrame(columns=['MAE','MSE','R2-score'])
```

```
for method, func in regressors.items():
    model = func.fit(X_train, y_train)
    pred = model.predict(X_test)
    results.loc[method] = [np.round(mean_absolute_error(y_test, pred), 3),
                           np.round(mean_squared_error(y_test, pred), 3),
                           np.round(r2_score(y_test, pred), 3)
                           ]
```

In [76]: *# Deep Neural Network*

```
# Splitting train set into training and validation sets.
X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, test_size=0.20)
```

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]: