

Avocado Dataset Analysis and ML Prediction</center>



Table of Contents

1. Problem Statement
2. Data Loading and Description
 - 3.1 Understanding the Dataset
 - 3.2 Profiling
 - 3.3 Preprocessing
4. Data Visualisation
 - 4.1 Which type of Avocados are more in demand (Conventional or Organic)?
 - 4.2 In which range Average price lies, what is distribution look like?
 - 4.3 How Average price is distributed over the months for Conventional and Organic Types?
 - 4.4 In which year and for which region was the Average price the highest?
 - 4.5 How dataset features are correlated with each other?
5. Feature Engineering for Model Building
 - 6.1 Average Price prediction using Linear Regression Regressor
 - 6.2 Average Price prediction using Decision Tree Regressor
 - 6.3 Average Price prediction using Random Forest Regressor
 - 6.4 Lets see final Actual Vs Predicted sample.
7. Conclusions

1. Problem Statement

- The notebooks explores the use to Pandas and EDA for analysis purpose.
- In this notebook we will try to predict the Avocado's Average Price based on different features. The features are different
 - Total Bags, Date, Type, Year, Region etc.
- Categorical: 'region', 'type'
- Date: 'Date'
- Numerical: 'Total Volume', '4046', '4225', '4770', 'Total Bags', 'Small Bags', 'Large Bags', 'XLarge Bags', 'year'
- Target: 'AveragePrice'

2. Data Loading and Description

This data provided by (INSAD), from the Hass Avocado Board website in May of 2018 & compiled into a single CSV. The dataset comprises of **18249** observations of **14** columns.

Below is a table showing names of all the columns and their description.

	Features	Description
Unnamed: 0	0	It's just a useless index feature that will be removed later
Total Volume		Total sales volume of avocados
4046		Total sales volume of Small/Medium Hass Avocado
4225		Total sales volume of Large Hass Avocado
4770		Total sales volume of Extra Large Hass Avocado
Total Bags		Total number of Bags sold
Small Bags		Total number of Small Bags sold
Large Bags		Total number of Large Bags sold
XLarge Bags		Total number of XLarge Bags sold

```
In [11]: import pandas as pd
import numpy as np
matplotlib.use("Agg", warn=False)
import matplotlib.pyplot as plt
import warnings
import seaborn as sns
import pandas_profiling
import pickle
import matplotlib

import plotly.offline as py
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 [23]: df = pd.read_csv("data/avocado.csv")
```

```
In [24]: df.head()
```

```
Out[24]:
```

Unnamed: 0	Date	AveragePrice	Total Volume	4046	4225	4770	Total Bags	Small Bags	Large Bags	XLarge Bags	type
0	2015-12-27	1.33	64236.62	1036.74	54454.85	48.16	8696.87	8603.62	93.25	0.0	conventional
1	2015-12-20	1.35	54876.98	674.28	44638.81	58.33	9505.56	9408.07	97.49	0.0	conventional
2	2015-12-13	0.93	118220.22	784.70	109149.67	130.50	8145.35	8042.21	103.14	0.0	conventional
3	2015-12-06	1.08	78992.15	1132.00	71976.41	72.58	5811.16	5677.40	133.76	0.0	conventional
4	2015-11-29	1.28	51039.60	941.48	43838.39	75.78	6183.95	5986.26	197.69	0.0	conventional

3. Data Cleaning

- Let's analyze and clean the data.

3.1 Understanding the Dataset

```
In [17]: df.shape
```

```
Out[17]: (18249, 14)
```

```
In [18]: df.columns # This will print the names of all columns.
```

```
Out[18]: Index(['Unnamed: 0', 'Date', 'AveragePrice', 'Total Volume', '4046', '4225', '4770', 'Total Bags', 'Small Bags', 'Large Bags', 'XLarge Bags', 'type', 'year', 'region'], dtype='object')
```

```
In [19]: df.head() # Will give you first 5 records
```

```
Out[19]:
```

Unnamed: 0	Date	AveragePrice	Total Volume	4046	4225	4770	Total Bags	Small Bags	Large Bags	XLarge Bags	type
0	2015-12-27	1.33	64236.62	1036.74	54454.85	48.16	8696.87	8603.62	93.25	0.0	conventional
1	2015-12-20	1.35	54876.98	674.28	44638.81	58.33	9505.56	9408.07	97.49	0.0	conventional
2	2015-12-13	0.93	118220.22	784.70	109149.67	130.50	8145.35	8042.21	103.14	0.0	conventional
3	2015-12-06	1.08	78992.15	1132.00	71976.41	72.58	5811.16	5677.40	133.76	0.0	conventional
4	2015-11-29	1.28	51039.60	941.48	43838.39	75.78	6183.95	5986.26	197.69	0.0	conventional

- The data set contains **18249** rows and **14** columns.

```
In [11]: df.info() # This will give Index, Datatype and Memory Information
```

```
Out[11]:
```

	Unnamed: 0	AveragePrice	Total Volume	4046	4225	4770	Total Bags	Small Bags
count	18249.000000	18249.000000	1.824900e+04	1.824900e+04	1.824900e+04	1.824900e+04	1.824900e+04	1.824900e+04
mean	24.232232	1.405978	8.606440e+05	2.930086e+05	2.951546e+05	2.983974e+04	2.396392e+05	1.821947e+05
std	15.481045	0.440877	3.453545e+06	1.264089e+06	1.204120e+06	1.074561e+05	9.892424e+05	7.461785e+05
min	0.000000	0.440877	8.456000e+01	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
25%	10.000000	1.100000	1.068385e+04	8.540700e+03	3.008780e+03	0.000000e+00	5.088840e+03	2.849420e+03
50%	24.000000	1.370000	1.073768e+05	8.645300e+03	2.951102e+04	1.849990e+02	3.974383e+04	2.636282e+04
75%	38.000000	1.860000	4.299293e+05	1.110202e+05	1.502099e+05	6.243420e+03	1.107834e+05	8.333767e+04
max	52.000000	3.250000	6.250955e+07	2.274932e+07	1.047057e+07	2.546439e+06	1.937315e+07	1.338459e+07

- We can see all columns having count **18249**. Looks like it doesn't contain missing values.

```
In [14]: df.isnull().sum() # Will show you null count for each column, but will not count Series() as null
```

```
Out[14]:
```

Unnamed: 0	Date	AveragePrice	Total Volume	4046	4225	4770	Total Bags	Small Bags	Large Bags	XLarge Bags	type	year	region
0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0	0	0	0	0
24	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0	0	0	0	0
27	0	0	0	0	0	0	0	0	0	0	0	0	0
28	0	0	0	0	0	0	0	0	0	0	0	0	0
29	0	0	0	0	0	0	0	0	0	0	0	0	0
30	0	0	0	0	0	0	0	0	0	0	0	0	0
31	0	0	0	0	0	0	0	0	0	0	0	0	0
32	0	0	0	0	0	0	0	0	0	0	0	0	0
33	0	0	0	0	0	0	0	0	0	0	0	0	0
34	0	0	0	0	0	0	0	0	0	0	0	0	0
35	0	0	0	0	0	0	0	0	0	0	0	0	0
36	0	0	0	0	0	0	0	0	0	0	0	0	0
37	0	0	0	0	0	0	0	0	0	0	0	0	0
38	0	0	0	0	0	0	0	0	0	0	0	0	0
39	0	0	0	0	0	0	0	0	0	0	0	0	0
40	0	0	0	0	0	0	0	0	0	0	0	0	0
41	0	0	0	0	0	0	0	0	0	0	0	0	0
42	0	0	0	0	0	0	0	0	0	0	0	0	0
43	0	0	0	0	0	0	0	0	0	0	0	0	0
44	0	0	0	0	0	0	0	0	0	0	0	0	0
45	0	0	0	0	0	0	0	0	0	0	0	0	0
46	0	0	0	0	0	0	0	0	0	0	0	0	0
47	0	0	0	0	0	0	0	0	0	0	0	0	0
48	0	0	0	0	0	0	0	0	0	0	0	0	0
49	0	0	0	0	0	0	0	0	0	0	0	0	0
50	0	0	0	0	0	0	0	0	0	0	0	0	0
51	0	0	0	0	0	0	0	0	0	0	0	0	0
52	0	0	0	0	0	0	0	0	0	0	0	0	0
53	0	0	0	0	0	0	0	0	0	0	0	0	0
54	0	0	0	0	0	0	0	0	0	0	0	0	0
55	0	0	0	0	0	0	0	0	0	0	0	0	0
56	0	0	0	0	0	0	0	0	0	0	0	0	0
57	0	0	0	0	0	0	0	0	0	0	0	0	0
58	0	0	0	0	0	0	0	0	0	0	0	0	0
59	0	0	0	0	0	0	0	0	0	0	0	0	0
60	0	0	0	0	0	0	0	0	0	0	0	0	0
61	0	0	0	0	0	0	0	0	0	0	0	0	0
62	0	0	0	0	0	0	0	0	0	0	0	0	0
63	0	0	0	0	0	0	0	0	0	0	0	0	0
64	0	0	0	0	0	0	0	0	0	0	0	0	0
65	0	0	0	0	0	0	0	0	0	0	0	0	0
66	0	0	0	0	0	0	0	0	0	0	0	0	0
67	0	0	0	0	0	0	0	0	0	0	0	0	0
68	0	0	0	0	0	0	0	0	0	0	0	0	0
69	0	0	0	0	0	0	0	0	0	0	0	0	0
70	0	0	0	0	0	0	0	0	0	0	0	0	0
71	0	0	0	0	0	0	0	0	0	0	0	0	0
72	0	0	0	0	0	0	0	0	0	0	0	0	0
73	0	0	0	0	0	0	0	0	0	0	0	0	0
74	0	0	0	0	0	0	0	0	0	0	0	0	0
75	0	0	0	0	0	0	0	0	0	0	0	0	0
76	0	0	0	0	0	0	0	0	0	0	0	0	0
77	0	0	0	0	0	0	0	0	0	0	0	0	0
78	0	0	0	0	0	0	0	0	0	0	0	0	0
79	0	0	0	0	0	0	0	0	0	0	0	0	0
80	0	0	0	0	0	0	0	0	0	0	0	0	0
81	0	0	0	0	0	0	0	0	0	0	0	0	0
82	0	0	0	0	0	0	0	0	0	0	0	0	0
83	0	0	0	0	0</								