# Machine Learning 1 Project

## Avocado Pricing Dataset

Presented by Clifer Fernandez – April 2020 cohort

# Agenda

- Introduction
- Executive Summary
- Data Profiling and Analysis
- Machine Learning Model
- Conclusions
- Recommendations

# Introduction

- Dataset  -  Average Pricing for single Avocado
  - 10 numerical – price, sales volume, etc.
  - 3 categorical – type, region, year

## Our Task

- Select and Train a ML model to help predict the average price of a single avocado

# Executive Summary

- Original data set – 18249 records, no missing / null
- Analysis
  - Covers years 2015 – Q1 2018
  - Extremely high correlation between majority of variables
  - ***Small / medium size of conventionally grown, packed in small bags are preferred***
  - Increasing sales volume over the years
- Machine Learning Model
  - Regression models preferred
  - Random Forest Regressor model gives best result
- Recommendations
  - Other feature engineering methods
  - More powerful algorithms
    - Support Vector Machines (SVM)
    - Artificial Neural Networks (ANN)

# Data Profiling

- Data set
  - Import data
  - View records
  - Basic stats
  - Profiling

# Data Profiling - Pre-processing

- Pre-processing
  - Check data distribution
  - Review and convert to relevant data type
  - Review perceived outliers

```
 #   Column        Non-Null Count  Dtype
---  ------        --------------  -----
 0   date          18249 non-null  datetime64[ns]
 1   averageprice  18249 non-null  float64
 2   total volume  18249 non-null  int64
 3   4046          18249 non-null  int64
 4   4225          18249 non-null  int64
 5   4770          18249 non-null  int64
 6   total bags    18249 non-null  int64
 7   small bags    18249 non-null  int64
 8   large bags    18249 non-null  int64
 9   xlarge bags   18249 non-null  int64
 10  type          18249 non-null  category
 11  year          18249 non-null  category
 12  region        18249 non-null  category
 13  month         18249 non-null  category
 14  day           18249 non-null  category
dtypes: category(5), datetime64[ns](1), float64(1), int64(8)
```

# Data Profiling – Post Profiling

- **Post Profiling**
  - Rerun stats
    - profiling

**Dataset statistics**

| | |
|---|---|
| Number of variables | 16 |
| Number of observations | 18249 |
| Missing cells | 0 |
| Missing cells (%) | 0.0% |
| Duplicate rows | 0 |
| Duplicate rows (%) | 0.0% |
| Total size in memory | 1.6 MiB |
| Average record size in memory | 93.3 B |

**Variable types**

| | |
|---|---|
| NUM | 10 |
| CAT | 5 |
| DATE | 1 |

# Data Analysis

## Data spread

- 3+ years worth of data

- Average Price
  - Normal distribution
  - Target variable for prediction
- Total Volume
  - Right skew in data
  - Includes data points for 'Total US'

# Data Analysis

**Volume Sale by Avocado Grade**

- Small and medium size preferred over large

- X large is a niche product

- Spike in sale around Feb every year

# Data Analysis

**Bag size preference**

- Small bag size preferred

- XL size sold is very small numbers

# Data Analysis

**Annual variations in Average Price**

- range is between 1.2 – 1.5

- peaks during Sep – Oct every year

- 2017 prices were generally higher

  – attributed to a poor harvest leading to shortage in supply



| | averageprice | | | |
|---|---|---|---|---|
| **year** | **2015** | **2016** | **2017** | **2018** |
| **month** | | | | |
| **1** | 1.280 | 1.215 | 1.230 | 1.380 |
| **2** | 1.290 | 1.190 | 1.170 | 1.340 |
| **3** | 1.290 | 1.170 | 1.355 | 1.325 |
| **4** | 1.360 | 1.135 | 1.460 | NaN |
| **5** | 1.290 | 1.140 | 1.505 | NaN |
| **6** | 1.335 | 1.250 | 1.500 | NaN |
| **7** | 1.280 | 1.330 | 1.495 | NaN |
| **8** | 1.370 | 1.360 | 1.625 | NaN |
| **9** | 1.395 | 1.380 | 1.760 | NaN |
| **10** | 1.290 | 1.490 | 1.700 | NaN |
| **11** | 1.250 | 1.490 | 1.485 | NaN |
| **12** | 1.230 | 1.260 | 1.370 | NaN |

# Data Analysis

**Annual variations in Average Price by Type**

- Organic avocados are higher priced
  - Can be attributed to higher cost of production



| type | averageprice | |
| --- | --- | --- |
| month | conventional | organic |
| 1 | 1.060 | 1.530 |
| 2 | 1.020 | 1.540 |
| 3 | 1.110 | 1.540 |
| 4 | 1.120 | 1.580 |
| 5 | 1.060 | 1.560 |
| 6 | 1.135 | 1.640 |
| 7 | 1.190 | 1.690 |
| 8 | 1.190 | 1.755 |
| 9 | 1.230 | 1.830 |
| 10 | 1.310 | 1.760 |
| 11 | 1.180 | 1.680 |
| 12 | 1.060 | 1.570 |

# Data Analysis

**Annual variations in Total Volume by Type**

- Conventionally farmed avocados are by far the most in demand

- Organic avocados are a very 'niche' sale item

# Data Analysis

**Top 5 cities with highest average price**

- Hartford Springfield

- New York

- San Francisco

- Philadelphia

- Chicago

- Note: The 'regional' areas were not selected for the listing although they are plotted in the graph

# Data Analysis

**Top 5 cities with highest total volumes**

- Los Angeles

- New York

- Dallas

- Houston

- Phoenix

- Note: The 'regional' areas were not selected for the listing although they are plotted in the graph. The 'high' tower is the Total US volume data point

# Data Analysis

## Parameter correlation

- Very high correlation among variables
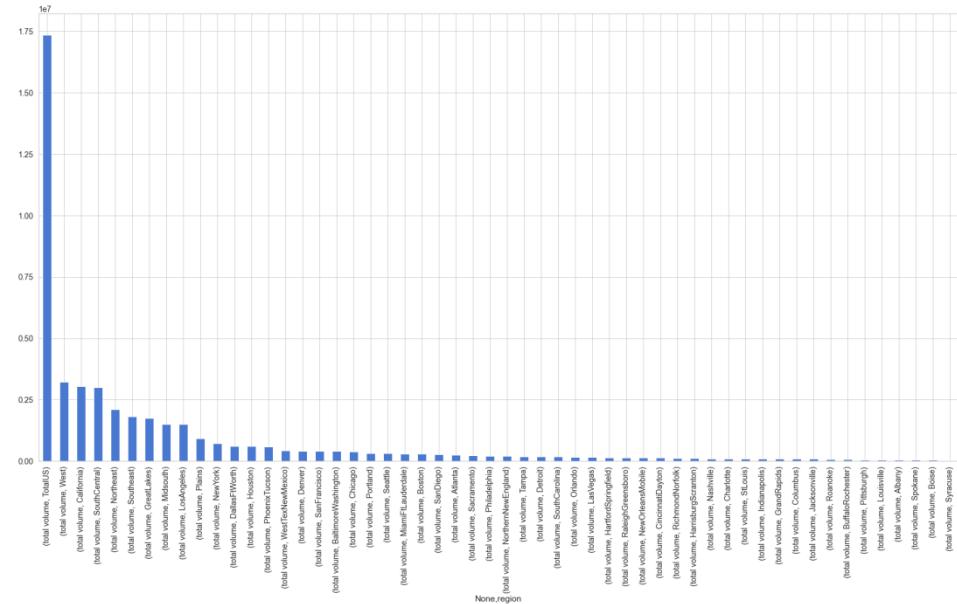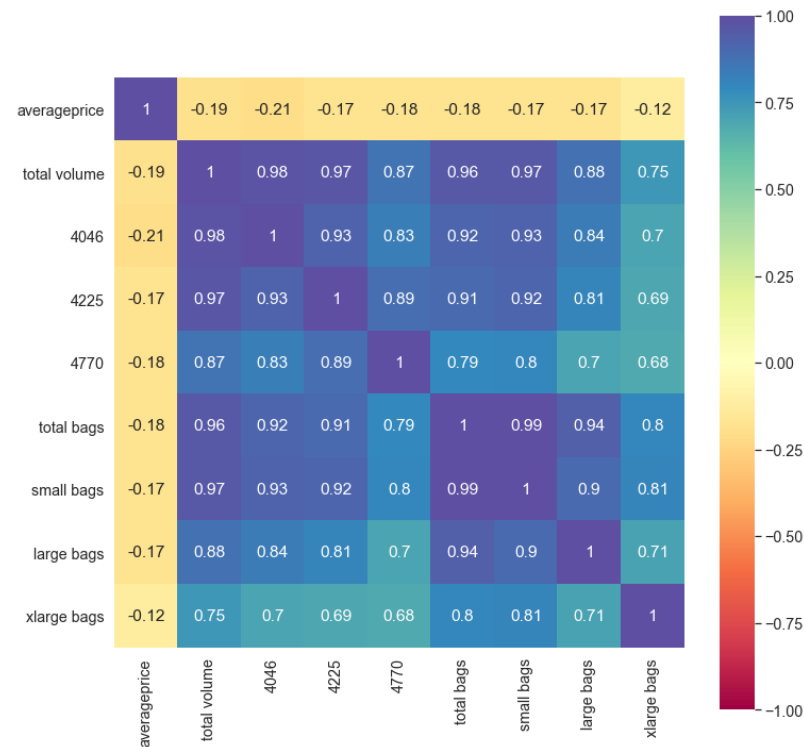  - Redundant information will be excluded for model building

|  | averageprice | total volume | 4046 | 4225 | 4770 | total bags | small bags | large bags | xlarge bags |
|---|---|---|---|---|---|---|---|---|---|
| **averageprice** | 1 | -0.19 | -0.21 | -0.17 | -0.18 | -0.18 | -0.17 | -0.17 | -0.12 |
| **total volume** | -0.19 | 1 | 0.98 | 0.97 | 0.87 | 0.96 | 0.97 | 0.88 | 0.75 |
| **4046** | -0.21 | 0.98 | 1 | 0.93 | 0.83 | 0.92 | 0.93 | 0.84 | 0.7 |
| **4225** | -0.17 | 0.97 | 0.93 | 1 | 0.89 | 0.91 | 0.92 | 0.81 | 0.69 |
| **4770** | -0.18 | 0.87 | 0.83 | 0.89 | 1 | 0.79 | 0.8 | 0.7 | 0.68 |
| **total bags** | -0.18 | 0.96 | 0.92 | 0.91 | 0.79 | 1 | 0.99 | 0.94 | 0.8 |
| **small bags** | -0.17 | 0.97 | 0.93 | 0.92 | 0.8 | 0.99 | 1 | 0.9 | 0.81 |
| **large bags** | -0.17 | 0.88 | 0.84 | 0.81 | 0.7 | 0.94 | 0.9 | 1 | 0.71 |
| **xlarge bags** | -0.12 | 0.75 | 0.7 | 0.69 | 0.68 | 0.8 | 0.81 | 0.71 | 1 |

# Machine Learning Model

**Regression Models Selected**

- Linear Regression

- Decision Tree Regressor

- Random Forest Regressor

# Machine Learning Model

**Model building**

- Train – Test Split dataset

- Scaling of the data

- Running the regression models

- Model evaluation

# Machine Learning Model

| Criteria | Linear Regression | | Decision Tree Regressor | | Random Forest Regressor | | Random Forest Regressor - GridsearchCV | |
|---|---|---|---|---|---|---|---|---|
| **Model Evaluation Matrix** | | | | | | | | |
| | Train | Test | Train | Test | Train | Test | Train | Test |
| **Mean Absolute Error** | 0.180 | 0.179 | 0.000 | 0.127 | 0.036 | 0.100 | 0.036 | 0.100 |
| **Mean Square Error** | 0.057 | 0.057 | 0.000 | 0.038 | 0.002 | 0.020 | 0.002 | 0.020 |
| **Root Mean Square Error** | 0.238 | 0.238 | 0.000 | 0.196 | 0.054 | 0.144 | 0.053 | 0.143 |
| **$R^2$** | 0.649 | 0.640 | 1.000 | 0.756 | 0.982 | 0.868 | 0.982 | 0.869 |
| **Adjusted $R^2$** | 0.647 | 0.630 | 1.000 | 0.749 | 0.981 | 0.865 | 0.982 | 0.865 |
| **Comment** | The Random Forest Regressor models have given the best Adjusted R2 values and are therefore the preferred models. The model obtained after Hyperparameter tuning (GridSearchCV) is very close to our original model which implies that our original model was good enough.<br><br>The best parameters returned by GridSearchCV are as follows:<br><br>• 'max_depth': None<br>• 'max_features': 'auto',<br>• 'n_estimators': 300 | | | | | | | |

# Conclusions

- Dataset for 3 full years and Q1 of the fourth
- Average Price is normally distributed
- Small / medium sized avocados in small sized bags are preferred
- Annual variation in price, peaking in Sept – Oct
- Conventionally farmed option preferred
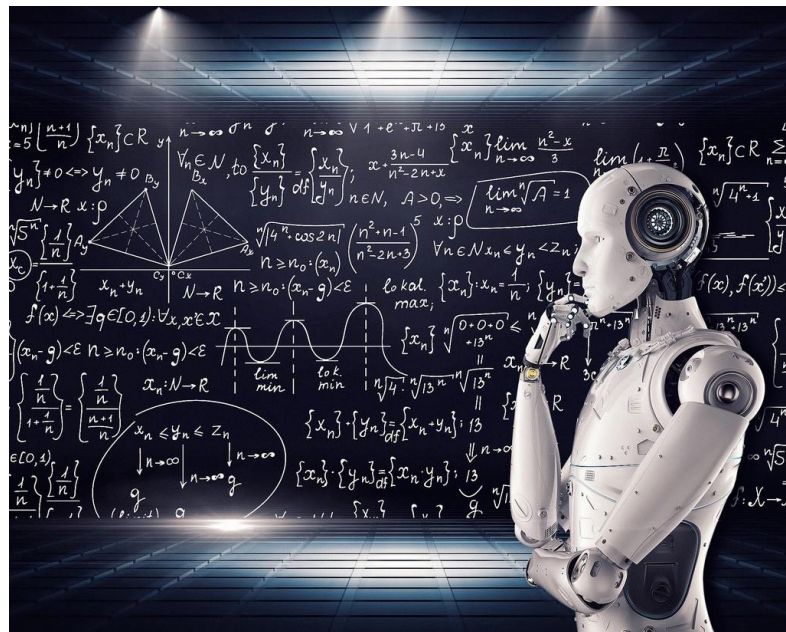  - Perhaps due to lower costs

# Conclusions

- Pricing can be predicted using Machine Learning Models

- Regression Models selected
  - Random Forest Regressor provided the best model

# Recommendations

- Model Improvement
  - More powerful Algorithms
    - Support Vector Machines (SVM)
    - Artificial Neural Networks (ANN)

# Thank You!

# Bon Appetite!