## **Unraveling the Mystery of Avocado Prices: A Time Series Detective Story**

by Chloe Van
Posted on Dec 5, 2023

Imagine you're a detective in a digital world where every tick of the clock leaves a breadcrumb of data. Your mission is to solve the mystery hidden within these clues, predicting future events based on past patterns. This is the essence of time series, the focus of my recent machine learning assignment.

In this digital detective story, this case revolves around avocados, a seasonal fruit known for their versatility, health benefits, and tie to nature, affecting their availability and price. The first step in my investigation was examining the nature of the avocado dataset to understand its structure, analogous to a detective investigating various leads.

|                            |  |                              | Date").head()                  |                         |                              |                |                    |                        |                |             |                         |                      |                  |            |
|----------------------------|--|------------------------------|--------------------------------|-------------------------|------------------------------|----------------|--------------------|------------------------|----------------|-------------|-------------------------|----------------------|------------------|------------|
|                            | Date   | AveragePrice                 | Total Volume                   | 4046                    | 422                          | 5 4770         | ) Total Ba         | gs Small Bag           | s Large Bag    | s XLarge Ba | gs                      | type                 | year             | regi       |
| 0                          | 2015-01-04   | 1.64                         | 6182.81                        | 1561.30                 | 2958.1                       | 7 0.00         | 1663.3             | 34 1663.3              | 4 0.0          | 0 (         | 0.0 orç                 | ganic :              | 2015 West        | TexNewMexi |
| 51                         | 2015-01-04   | 1.00                         | 435021.49                      | 364302.39               | 23821.16                     | 82.1           | 5 46815.           | 79 16707.1             | 5 30108.6      | 4 (         | 0.0 convent             | ional :              | 2015             | Atlar      |
| 51                         | 2015-01-04   | 1.93                         | 17328.24                       | 2357.18                 | 12692.2                      | 1 9.4          | 7 2269.3           | 38 2269.3              | 8 0.0          | 0 (         | ).0 orç                 | ganic :              | 2015             | NewYo      |
| 51                         | 2015-01-04   | 1.02                         | 160130.15                      | 4007.4                  | 1 118435.79                  | 1201.50        | 36485.4            | 15 20325.4             | 16160.0        | 4 (         | 0.0 convent             | ional :              | 2015             | Indianapo  |
| 51                         | 2015-01-04   | 1.69                         | 34190.37                       | 3874.3                  | 1 14945.34                   | 4 32.4         | 5 15338.:          | 27 13793.2             | 2 1545.0       | 5 (         | 0.0 org                 | organic 2015         |                  | Pla        |
| O                          | lf_train.sor   | t_values(by="[               | Date", ascend:                 | ing=False)              | .head()                      |                |                    |                        |                |             |                         |                      |                  |            |
|                            | Date   | AveragePrice                 | Total Volume                   | 4046                    | 4225                         | 4770           | Total Bag          | s Small Bags           | Large Bage     | XLarge Bag  | s type                  | year                 |                  | region     |
| 4                          | 2017-09-24   | 1.42                         | 3627.18                        | 56.82                   | 95.31                        | 0.00           | 3475.0             | 5 3475.08              | 0.00           | 0.0         | 0 organic               | 2017                 |                  | Albany     |
| 4                          | 2017-09-24   | 2.05                         | 138458.82                      | 34599.22                | 63107.77                     | 1.57           | 40750.2            | 6 40700.4 <sup>-</sup> | 49.8           | 0.0         | 0 organic               | 2017                 |                  | California |
| 4                          | 2017-09-24   | 2.56                         | 18912.30                       | 507.34                  | 5068.56                      | 0.00           | 13336.4            | 0 4252.99              | 9083.4         | 0.0         | 0 organic               | 2017                 |                  | Atlanta    |
| 4                          | 2017-09-24   | 1.53                         | 62464.93                       | 2666.66                 | 6252.03                      | 94.49          | 53451.7            | 5 53451.78             | 0.00           | 0.0         | 0 organic               | 2017                 | BaltimoreW       | ashington  |
| 4                          | 2017-09-24   | 1.94                         | 1125443.38                     | 116338.61               | 237876.77                    | 1182.85        | 769536.9           | 0 639021.62            | 130401.43      | 113.8       | 5 organic               | 2017                 |                  | TotalUS    |
| C                          | lf_train.sor   | t_values(by=['               | 'region", "Dat                 | te"]).head              | (10)                         |                |                    |                        |                |             |                         |                      |                  |            |
|                            | Date   | AveragePrice                 | Total Volume                   | 4046                    | 4225                         | 4770 T         | otal Bags          | Small Bags I           | .arge Bags     | (Large Bags | type                    |                      |                  |            |
|                            | 2015-01-04   | 1.22                         | 40873.28                       | 2819.50                 | 28287.42                     | 49.90          | 9716.46            | 9186.93                | 529.53         | 0.0         | conventional            | 2015                 | Albany           |            |
| 51                         |  | 1.79                         | 1373.95                        | 57.42                   | 153.88                       | 0.00           | 1162.65            | 1162.65                | 0.00           | 0.0         | organic                 |                      |                  |            |
|                            | 2015-01-04   |                              |                                |                         | 31640.34                     | 127.12         | 8424.77            | 8036.04                | 388.73         | 0.0         | conventional            | 2015                 | Albany           |            |
| 51                         | 2015-01-11   | 1.24                         | 41195.08                       | 1002.85                 |                              |                |                    |                        | 0.00           | 0.0         | organic                 | 2015                 | Albany           |            |
| 0                          |  |                              | 41195.08<br>1182.56            | 1002.85<br>39.00        | 305.12                       | 0.00           | 838.44             | 838.44                 | 0.00           | 0.0         | organic                 |                      |                  |            |
| 0                          | 2015-01-11   | 1.24                         |                                |                         | 305.12                       | 0.00<br>135.77 | 838.44<br>11921.05 | 838.44<br>11651.09     | 269.96         |             | conventional            |                      | Albany           |            |
| 51<br>60<br>60             | 2015-01-11<br>2015-01-11                             | 1.24<br>1.77                 | 1182.56                        | 39.00                   | 305.12                       |                |                    |                        |                |             |                         | 2015                 |                  |            |
| 51<br>0<br>0<br>9          | 2015-01-11<br>2015-01-11<br>2015-01-18               | 1.24<br>1.77<br>1.17         | 1182.56<br>44511.28            | 39.00<br>914.14         | 305.12<br>31540.32<br>178.78 | 135.77         | 11921.05           | 11651.09               | 269.96         | 0.0         | conventional            | 2015<br>2015         | Albany           |            |
| 51<br>50<br>60<br>69<br>88 | 2015-01-11<br>2015-01-11<br>2015-01-18<br>2015-01-18 | 1.24<br>1.77<br>1.17<br>1.93 | 1182.56<br>44511.28<br>1118.47 | 39.00<br>914.14<br>8.02 | 305.12<br>31540.32<br>178.78 | 0.00           | 11921.05<br>931.67 | 11651.09<br>931.67     | 269.96<br>0.00 | 0.0         | conventional<br>organic | 2015<br>2015<br>2015 | Albany<br>Albany |            |

Figure 1: A snippet of the avocado dataset sorted by: Date ascending (top), Date descending (middle), Date & Region ascending (bottom)

My exploratory data analysis (Figure 1) revealed two distinct time series representing unique sequences of avocado prices over time. This was identified by observing measurements across regions and types on the same day. It was crucial to confirm that these time series had equally spaced measurements, similar to verifying the reliability of witnesses in an investigation. This consistency was vital to my analysis, as irregular data intervals could lead to misleading conclusions.

|           | <pre>df_hastarget = create_lag_feature(df_sort, "AveragePrice", +1, ["region", "type"], "AveragePriceNextWeek", clip=True) df_hastarget</pre> |              |              |         |          |        |            |            |            |             |              |      |                  |                      |
|-----------|---|--------------|--------------|---------|----------|--------|------------|------------|------------|-------------|--------------|------|------------------|----------------------|
|           |   |              |              |         |          |        |            |            |            |             |              |      |                  |                      |
|           | Date  | AveragePrice | Total Volume | 4046    | 4225     | 4770   | Total Bags | Small Bags | Large Bags | XLarge Bags | type         | year | region           | AveragePriceNextWeek |
|           | 2015-01-04  | 1.22         | 40873.28     | 2819.50 | 28287.42 | 49.90  | 9716.46    | 9186.93    | 529.53     | 0.0         | conventional | 2015 | Albany           | 1.24                 |
|           | 2015-01-11  | 1.24         | 41195.08     | 1002.85 | 31640.34 | 127.12 | 8424.77    | 8036.04    | 388.73     | 0.0         | conventional | 2015 | Albany           | 1.17                 |
|           | 2015-01-18  | 1.17         | 44511.28     | 914.14  | 31540.32 | 135.77 | 11921.05   | 11651.09   | 269.96     | 0.0         | conventional | 2015 | Albany           | 1.06                 |
|           | 2015-01-25  | 1.06         | 45147.50     | 941.38  | 33196.16 | 164.14 | 10845.82   | 10103.35   | 742.47     | 0.0         | conventional | 2015 | Albany           | 0.99                 |
|           | 2015-02-01  | 0.99         | 70873.60     | 1353.90 | 60017.20 | 179.32 | 9323.18    | 9170.82    | 152.36     | 0.0         | conventional | 2015 | Albany           | 0.99                 |
|           |   |              |              |         |          |        |            |            |            |             |              |      |                  |                      |
| 18243     | 2018-02-18  | 1.56         | 17597.12     | 1892.05 | 1928.36  | 0.00   | 13776.71   | 13553.53   | 223.18     | 0.0         | organic      | 2018 | WestTexNewMexico | 1.57                 |
| 18244     | 2018-02-25  |              | 18421.24     | 1974.26 | 2482.65  | 0.00   | 13964.33   | 13698.27   | 266.06     | 0.0         | organic      | 2018 | WestTexNewMexico | 1.54                 |
| 18245     | 2018-03-04  | 1.54         | 17393.30     | 1832.24 | 1905.57  | 0.00   | 13655.49   | 13401.93   | 253.56     | 0.0         | organic      | 2018 | WestTexNewMexico | 1.56                 |
| 18246     | 2018-03-11  | 1.56         | 22128.42     | 2162.67 | 3194.25  | 8.93   | 16762.57   | 16510.32   | 252.25     | 0.0         | organic      | 2018 | WestTexNewMexico | 1.56                 |
| 18247     | 2018-03-18  | 1.56         | 15896.38     | 2055.35 | 1499.55  | 0.00   | 12341.48   | 12114.81   | 226.67     | 0.0         | organic      | 2018 | WestTexNewMexico | 1.62                 |
| 18141 row | s × 14 columns  |              |              |         |          |        |            |            |            |             |              |      |                  |                      |

Figure 2: A snippet of creating a lag feature named AveragePriceNextWeek

I created an 'AveragePriceNextWeek' feature, a lagged version of the 'AveragePrice', to uncover patterns for future insights (Figure 2). By providing a retrospective lens, this lagged feature helped me identify how past prices influenced avocado prices the following week. To ensure the strength of my investigation, I excluded rows where the 'target' data was missing, thereby preserving the integrity and accuracy of my findings akin to a detective ensuring that only complete and reliable evidence is considered in a case.

I then developed models to forecast next week's average avocado prices. As a foundational step, I established a baseline model, assuming that 'if today's price is X, then tomorrow's price will likely be the same.' This model's R² scores serve as a reference point to compare more complex models. As a benchmark, if a complex model couldn't outperform this baseline model, it must not be a very good model.

An R² score measures how well my "detective's hypothesis" (in this case, my baseline and linear regression models) matches the actual "facts of the case" (the real-world target of avocado prices). An R² score ranges from 0 to 1, where scores closer to 1 indicate that this model does a great job of predicting avocado prices. On the other hand, a score closer to 0 would mean that this model doesn't do a good job capturing the data trends and makes poor predictions.

Next, I experimented with a few approaches with Linear Regression, which helps me draw a straight line through a scatter plot of data points - each representing an avocado's price on a specific day. This line attempts to capture the general trend of data, showing me the direction of avocado prices moving over time. I chose Linear Regression for its simplicity and effectiveness,

providing a clear and straightforward way to understand trends with continuous values as targets.

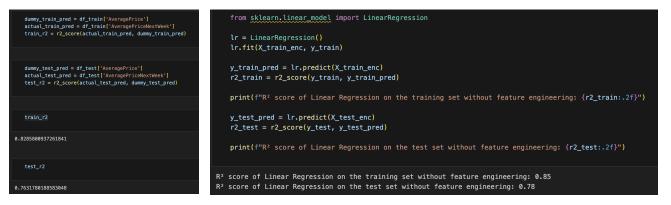


Figure 3: Left: snippet of baseline model R2 scores, Right: snippet of linear regression model R2 scores

As shown in Figure 3, the R² scores from Linear Regression only slightly improved compared to the baseline, suggesting that this basic approach (without feature engineering) was too straightforward. This model couldn't grasp the periodic nature of the data, which lay in how I encoded Date as integers. This approach limited the model's effectiveness, which could not flex around the more complex patterns of time.

To dive deeper, I experimented with various methods of encoding time. The first method involved extracting the "month" from the Date to capture the fluctuations in avocado prices throughout the year. Given that avocados are seasonal fruits, the relevance of the months seemed significant in understanding the relationship between price and availability, especially under the influence of supply and demand laws. The improvement in my results confirmed that this approach was more insightful, as seen in the improved R<sup>2</sup> scores below (Figure 4).

```
# Linear Regression with feature engineering extracting Month from Date

X_train_enc, y_train, X_test_enc, y_test, preprocessor = preprocess_features(
    df_month_train,
    df_month_test,
    numeric_features,
    categorical_features + ["Month"],
    drop_features,
    target
}

lr_month = LinearRegression()
lr_month.fit(X_train_enc, y_train)

y_train_pred = lr_month.predict(X_train_enc)
    r2_train = r2_score(y_train, y_train_pred)

print(f"R2 score of Linear Regression on the training set with Month: {r2_train:.2f}")

y_test_pred = lr_month.predict(X_test_enc)
    r2_test = r2_score(y_test, y_test_pred)

print(f"R2 score of Linear Regression on the test set with Month: {r2_test:.2f}")

Number of OHE features: 68
Total columns in new_columns: 78
Transformed data shape: (15441, 78)
R2 score of Linear Regression on the training set with Month: 0.85
R2 score of Linear Regression on the test set with Month: 0.85
R3 score of Linear Regression on the test set with Month: 0.85
R3 score of Linear Regression on the test set with Month: 0.85
R3 score of Linear Regression on the test set with Month: 0.80
```

```
# Linear Regression with feature engineering extracting Season from Month and Date

X_train_enc, y_train, X_test_enc, y_test, preprocessor = preprocess_features(
    df_season_train,
    df_season_train,
    df_season_train,
    df_season_test,
    numeric_features,
    cateoprical_features + ["Season"],
    drop_features,
    target
)

Ir_season = LinearRegression()
Ir_season.fit(X_train_enc, y_train)

y_train_pred = Ir_season.predict(X_train_enc)
    r2_train = r2_score(y_train, y_train_pred)

print(f"R2 score of Linear Regression on the training set with Season: {r2_train:.2f}")

y_test_pred = Ir_season.predict(X_test_enc)
    r2_test = r2_score(y_test, y_test_pred)

print(f"R2 score of Linear Regression on the test set with Season: {r2_test:.2f}")

Number of OHE features: 60
Total columns in new_columns: 70
Transformed data shape: (15441, 70)
R2 score of Linear Regression on the training set with Season: 0.79
R3 score of Linear Regression on the test set with Season: 0.79
```

Figure 4: Left: snippet of linear regression extracting month R<sup>2</sup> scores, Right: snippet of linear regression extracting seasons R<sup>2</sup> scores

The second method involved extracting "seasons" to categorize the data more broadly. However, this was like a detective making general assumptions. Seasons can have significant variability among themselves, where conditions in early summer versus late summer can differ significantly. The overgeneralization slightly reduced the model's effectiveness, as shown in the test set R² scores (Figure 4).

```
# Linear Regression with feature engineering extracting days_since
   X_train_enc, y_train, X_test_enc, y_test, preprocessor = preprocess_features(
       df_days_since_train,
       df_days_since_test,
       numeric_features + ["Days_since"],
       categorical features.
       drop_features,
       target
   lr_days_since = LinearRegression()
   lr_days_since.fit(X_train_enc, y_train)
   y_train_pred = lr_days_since.predict(X_train_enc)
   r2_train = r2_score(y_train, y_train_pred)
   print(f"R2 score of Linear Regression on the training set with days since: {r2 train:.2f}")
   y_test_pred = lr_days_since.predict(X_test_enc)
   r2_test = r2_score(y_test, y_test_pred)
   print(f"R2 score of Linear Regression on the test set with days_since: {r2_test:.2f}")
Number of OHE features: 56
Total columns in new columns: 67
Transformed data shape: (15441, 67)
R<sup>2</sup> score of Linear Regression on the training set with days_since: 0.85
R^{z} score of Linear Regression on the test set with days_since: 0.84
```

Figure 5: snippet of linear regression extracting days\_since R2 scores

The most effective method was extracting the "days since" feature, as shown by the best R<sup>2</sup> scores (Figure 5). Using the first Date of the data set as a reference date, this approach was like a detective piecing together an entire history, recognizing that patterns and behaviors evolve over long periods. This could potentially capture influences such as changing agricultural practices, climate change and increased market demand. I reduced seasonal bias by preventing the model from pigeonholing into expectations based on the time of the year only.

While the results are promising, there are several reasons why my model may not be as robust as they appear. One primary concern is the risk of overfitting, which means the model has become too tailored to the nuances of the patterns in this data set. Suppose the model is too finely tuned to this dataset with features like 'days since' and 'month.' In that case, it may not perform well with new data or under different avocado market conditions.

Figure 6: all unique regions in the data set

During my exploratory data analysis, I uncovered location overlaps in 'region' such as TotalUS and individual US states (Figure 6). This overlap could skew my results as the same data might be counted more than once. This could lead to overemphasizing certain information and result in misleading conclusions about avocado prices.

This analysis and features available in the data set (Figure 1) do not explicitly include factors such as economic trends and consumer preferences. By not incorporating these economic indicators, this model may miss key drivers of price changes, additionally causing misleading predictions.

In summary, experimenting with lag features and encoding Date in my time series data was crucial. It enabled me to uncover the intricate patterns of how avocado prices change over time, analogous to a detective using various tools to discover the deeper layers of a complex case.