



Time-Series Analysis

Prediction of Avocado-prices in the US

Ökonometrie 2
SS 2019



Group-members

Dornigg, Thomas (01551056)

Kafenda, Lukas (01607734)

Sunk, David (01605328)



Content

1. Project description
2. Exploratory Data Analysis of the Dataset
3. Testing for Stationarity
4. Correcting for Stationarity
5. Fitting appropriate model
 - a. Traditional approach (ARMA, ARIMA, SARIMA)
 - b. Facebook-Prophet Library
 - c. Comparison
6. Conclusion



1. Project and Dataset Description

- Kaggle dataset from 2018^[1]
- represents weekly retail scan data for national retail volume (units) and price
- based on actual retail sales of Hass avocados
- time-range: 04-01-2015 : 25-03-2018
- Research Questions:
 - Was the Avocadopocalypse of 2017 real? (taken from Kaggle)
 - What is the potential price which customers would have to pay for Conventional-Avocados in early 2018 respectively 2019?



2. Exploratory Data Analysis

- dataset consists of following features:
 - Date
 - Average Price
 - Type (organic/conventional)
 - Year
 - Total Volume
- avocado-type comparison (descriptive statistics, pie chart)
- evolution of avocado prices (line graph)
- comparison of avocado prices (density plot, boxplots, line graph)
- volume of sold avocados grouped by regions (frequency distribution bar chart)

Conventional

	Average Price	Total Volume	Region	Type
Types	numeric	numeric	constant	constant
Counts	169	169	169	169
Uniques	61	169	1	1
Missing	0	0	0	0
Missing %	0%	0%	0%	0%
Mean	1.09	33735038.97	NaN	NaN
StD	0.17	6118091.85	NaN	NaN
Minimum	0.76	21009730.21	NaN	NaN
Maximum	1.65	62505646.52	NaN	NaN
Mean abs. Deviation	0.13	4532699.23	NaN	NaN
25%	0.97	29761638.48	NaN	NaN
50%	1.04	32994014.16	NaN	NaN
75%	1.19	37026085.75	NaN	NaN
IQR	0.22	7264447.27	NaN	NaN
Skewness	1.13	1.23	NaN	NaN
Kurtosis	1.33	4.31	NaN	NaN

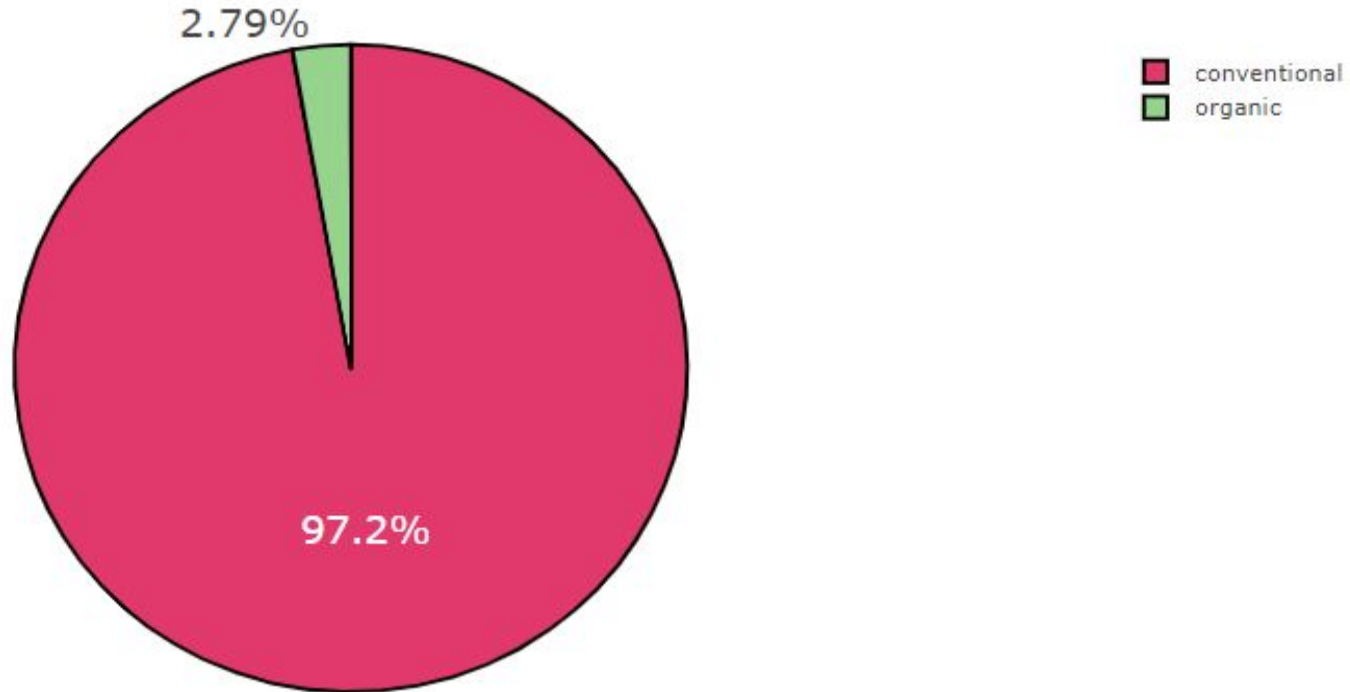
Organic

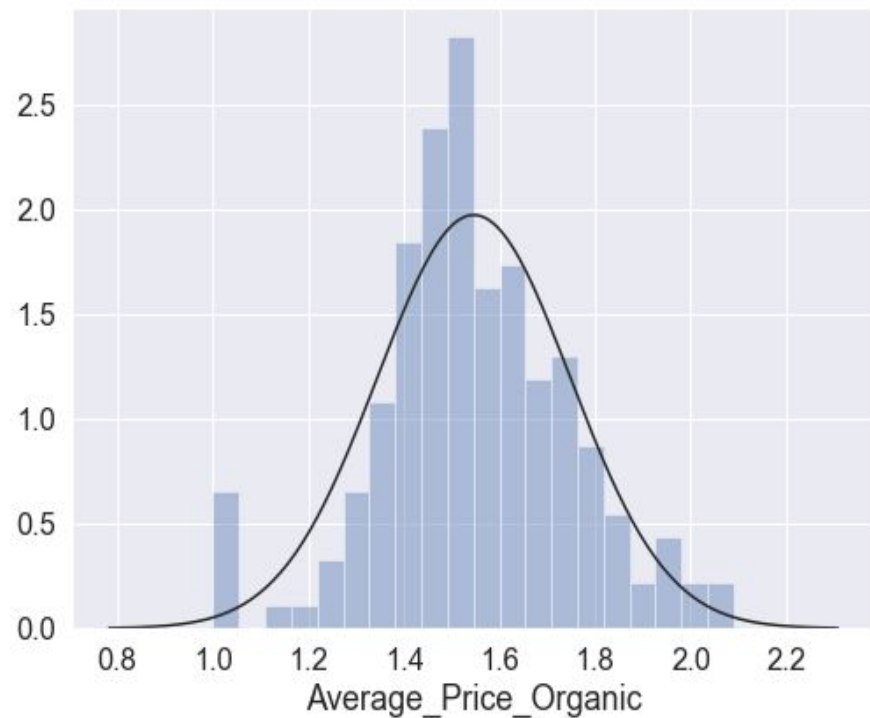
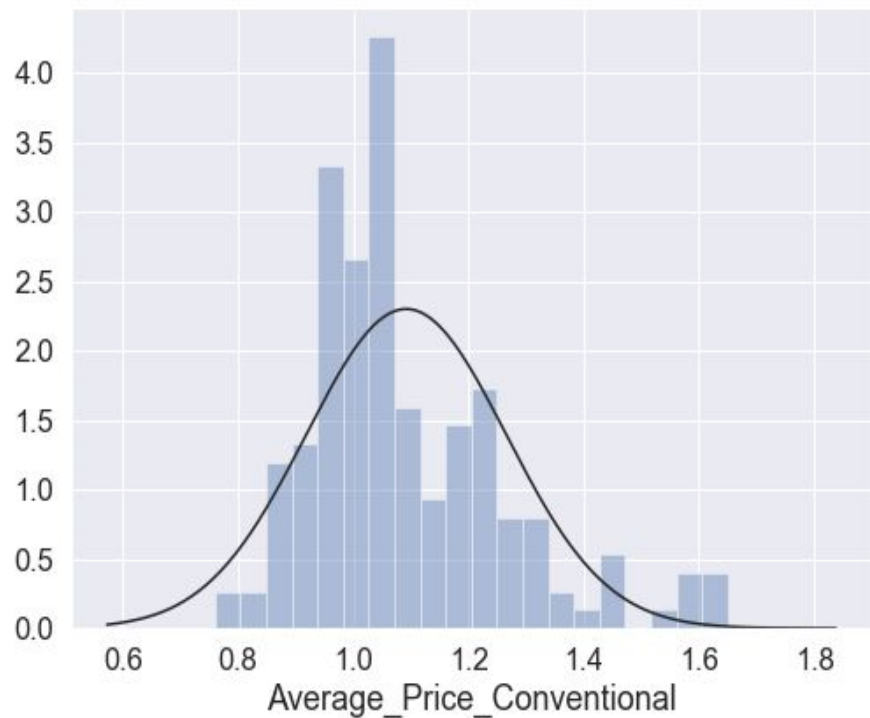
	Average Price	Total Volume	Region	Type
Types	numeric	numeric	constant	constant
Counts	169	169	169	169
Uniques	69	169	1	1
Missing	0	0	0	0
Missing %	0%	0%	0%	0%
Mean	1.55	967565.65	NaN	NaN
StD	0.20	302482.33	NaN	NaN
Minimum	1.00	501814.87	NaN	NaN
Maximum	2.09	1814929.97	NaN	NaN
Mean abs. Deviation	0.15	249361.63	NaN	NaN
25%	1.43	699763.35	NaN	NaN
50%	1.53	967886.13	NaN	NaN
75%	1.67	1148617.16	NaN	NaN
IQR	0.24	448853.81	NaN	NaN
Skewness	-0.13	0.47	NaN	NaN
Kurtosis	0.93	-0.44	NaN	NaN

Average Avocado Price from 2015 to 2018



Conventional vs. Organic Avocados grouped by Total Volume sold

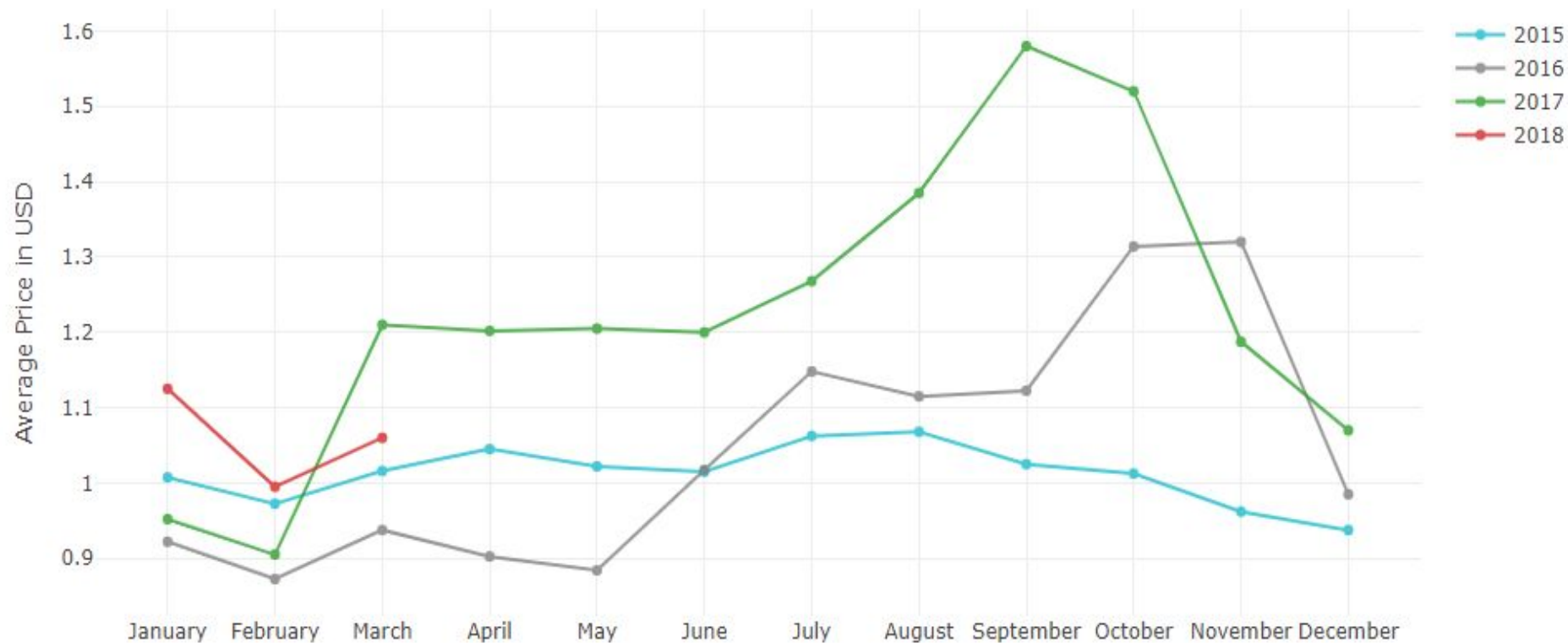




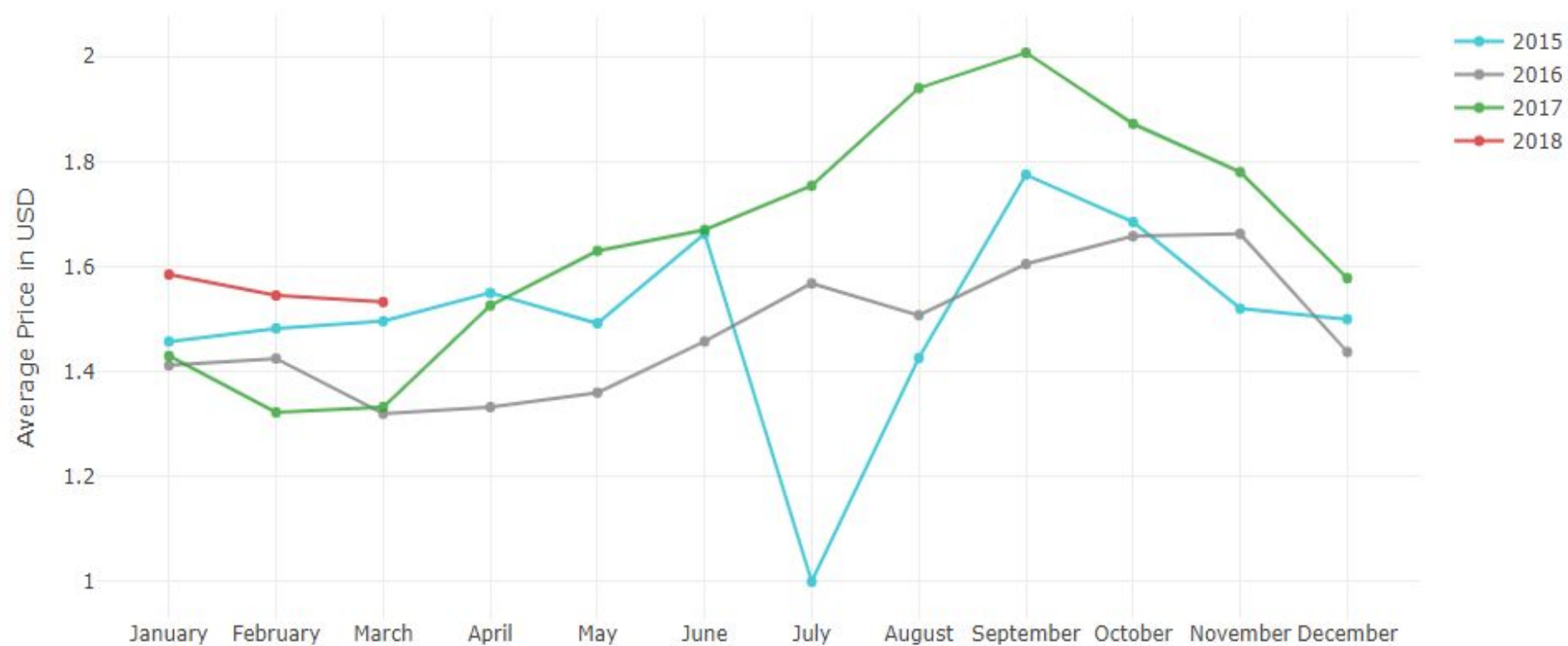
Price Distribution: Conventional vs. Organic



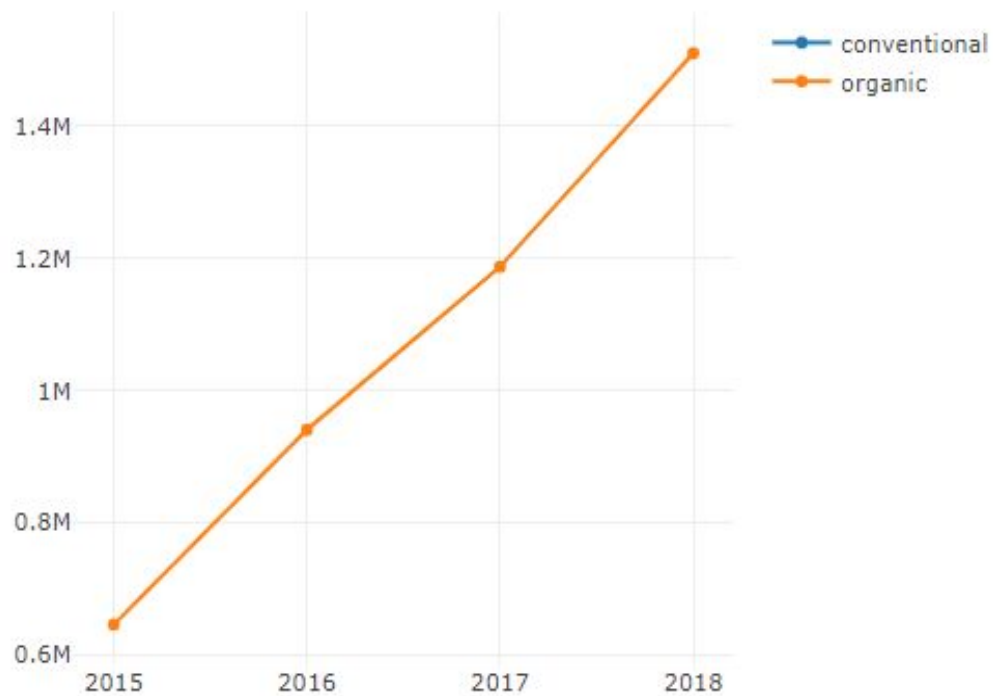
Price developement (grouped by months) for Conventional Avocados



Price developement (grouped by months) for Organic Avocados



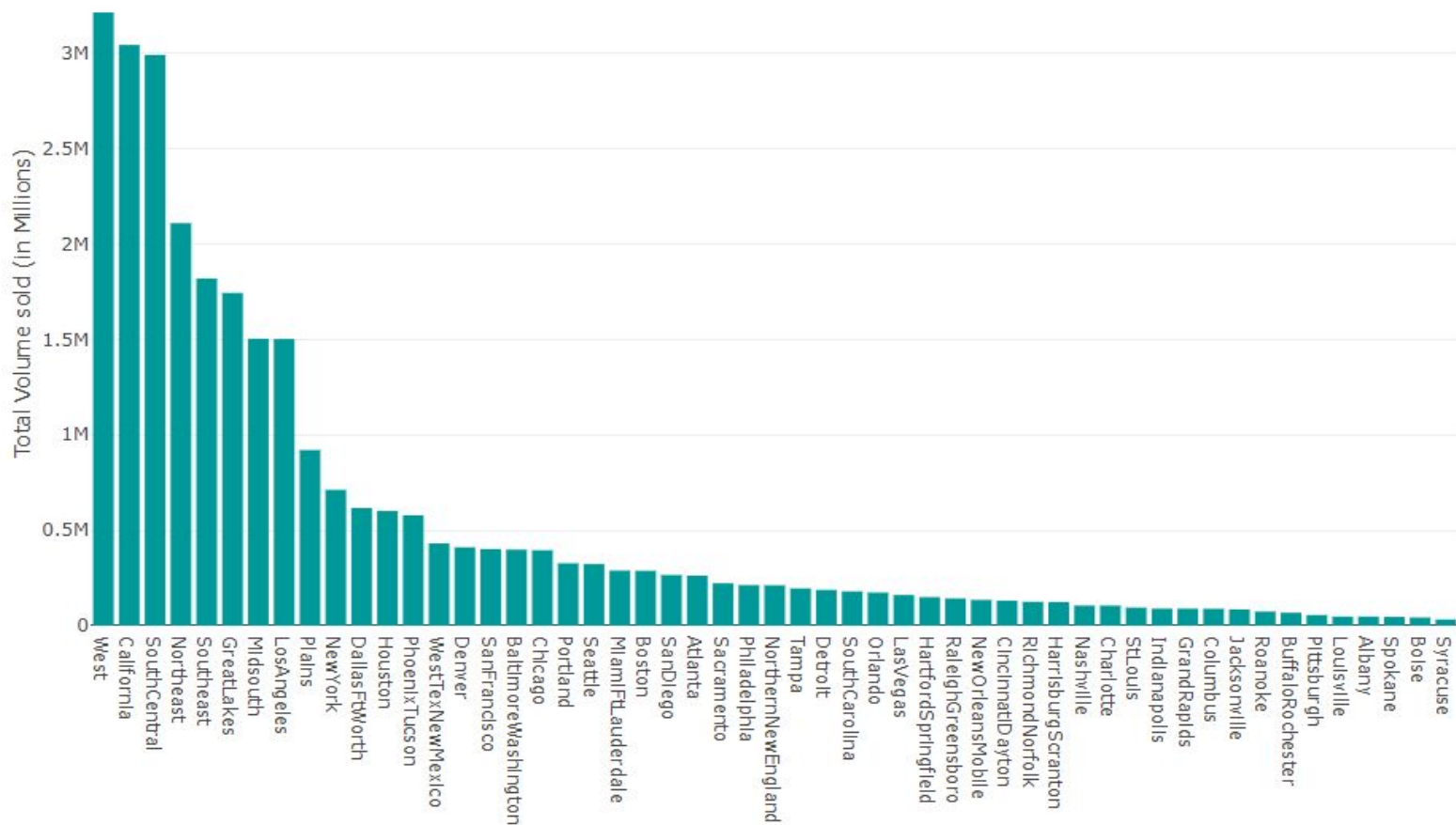
Average Volume Sold from 2015 to 2018



Average Price vs. Total Volume sold from 2015 to 2018



Total Volume of sold Avocados grouped by Regions (2015-2018)





3. Test for Stationarity

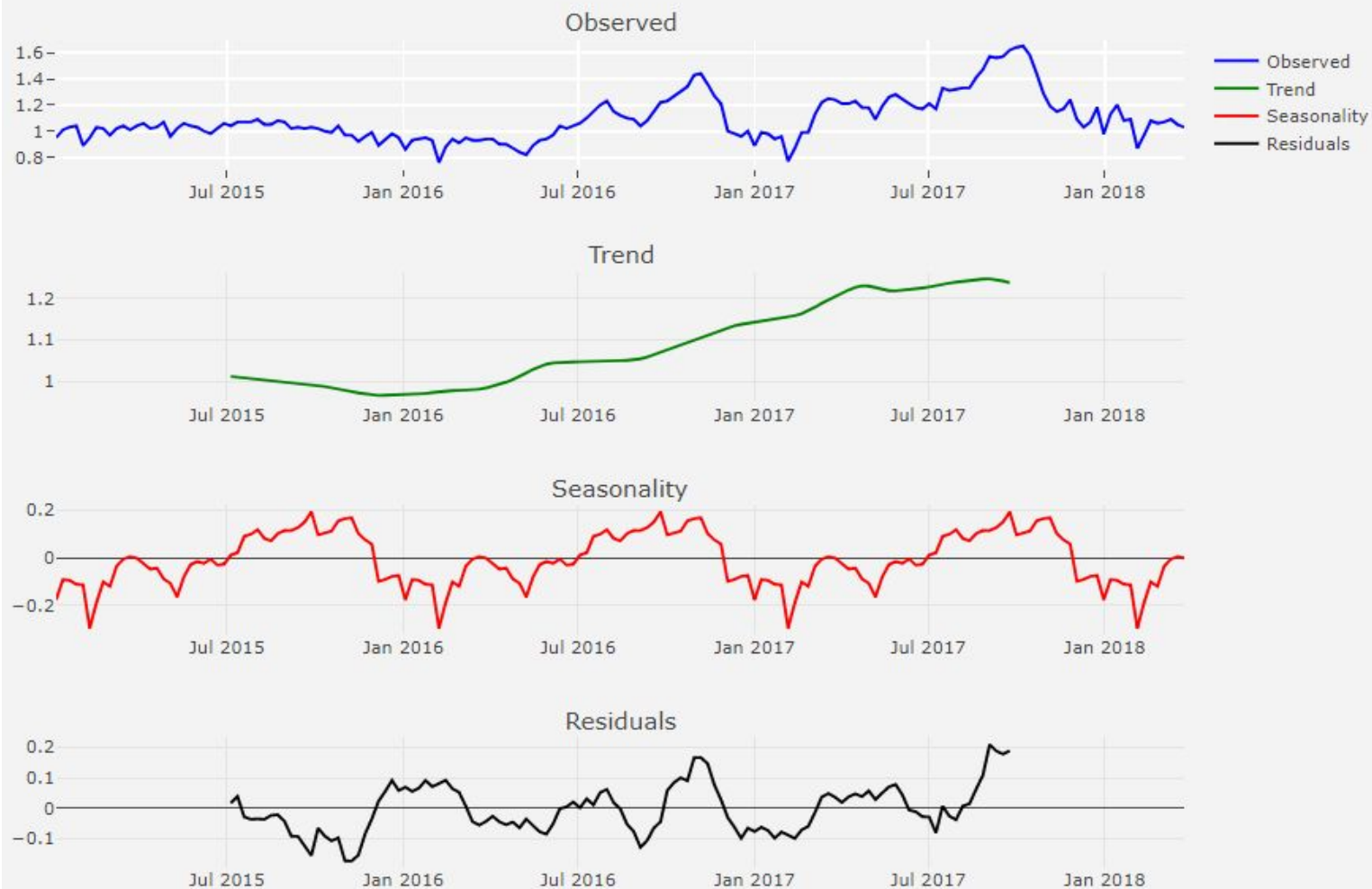
Stationarity is a statistical assumption that a time-series has:

- Constant mean
- Constant variance
- Autocovariance does not depend on time

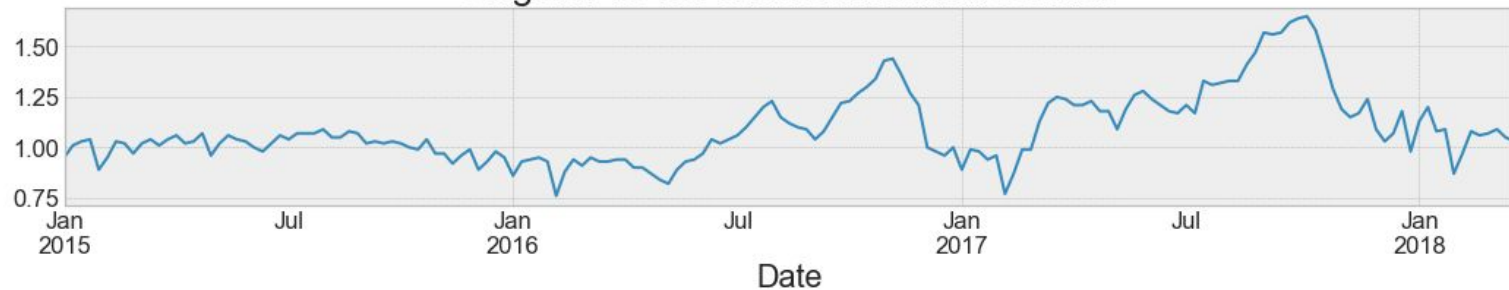


Testing both graphically and with statistical-tests

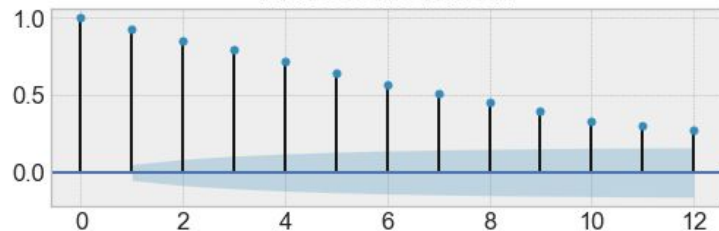
Decomposition of the Time Series for Conventional Avocados



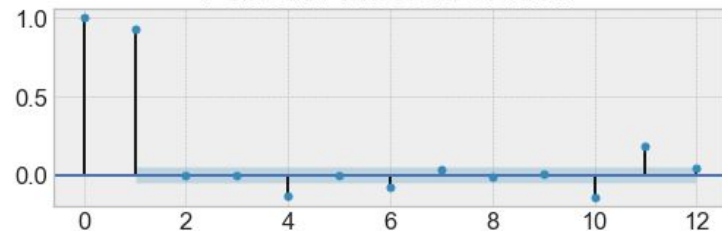
Original TS for Conventional Avocados



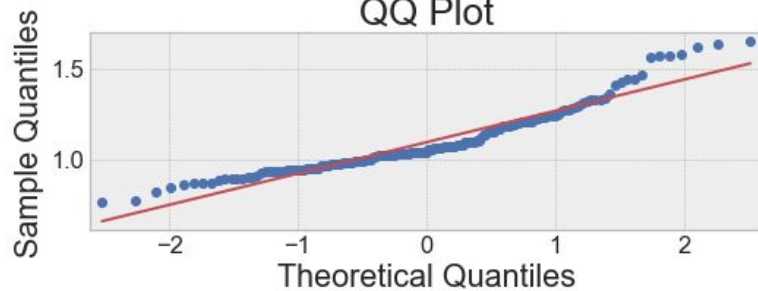
Autocorrelation



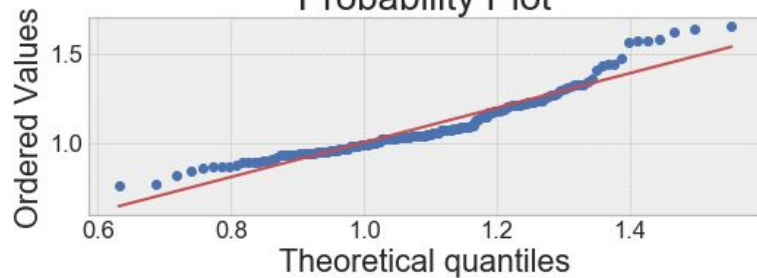
Partial Autocorrelation



QQ Plot



Probability Plot



Rolling mean and standard deviation for the Timeseries



Tests for Stationarity

1) Augmented Dickey-Fuller Test

H_0 : The Time-Series is not stationary

H_A : The Time-Series is stationary

2) KPSS-Test

H_0 : The Time-Series is stationary

H_A : The Time-Series is not stationary

Results for the conducted tests

Test-Statistic	p-value
Augmented-Dickey Fuller Test	0.098633
KPSS	0.034961



Assumption: Time-Series is not-stationary



4. Correcting for Stationarity

As seen from the graph and test-statistics before, the TS seems to be not stationary:

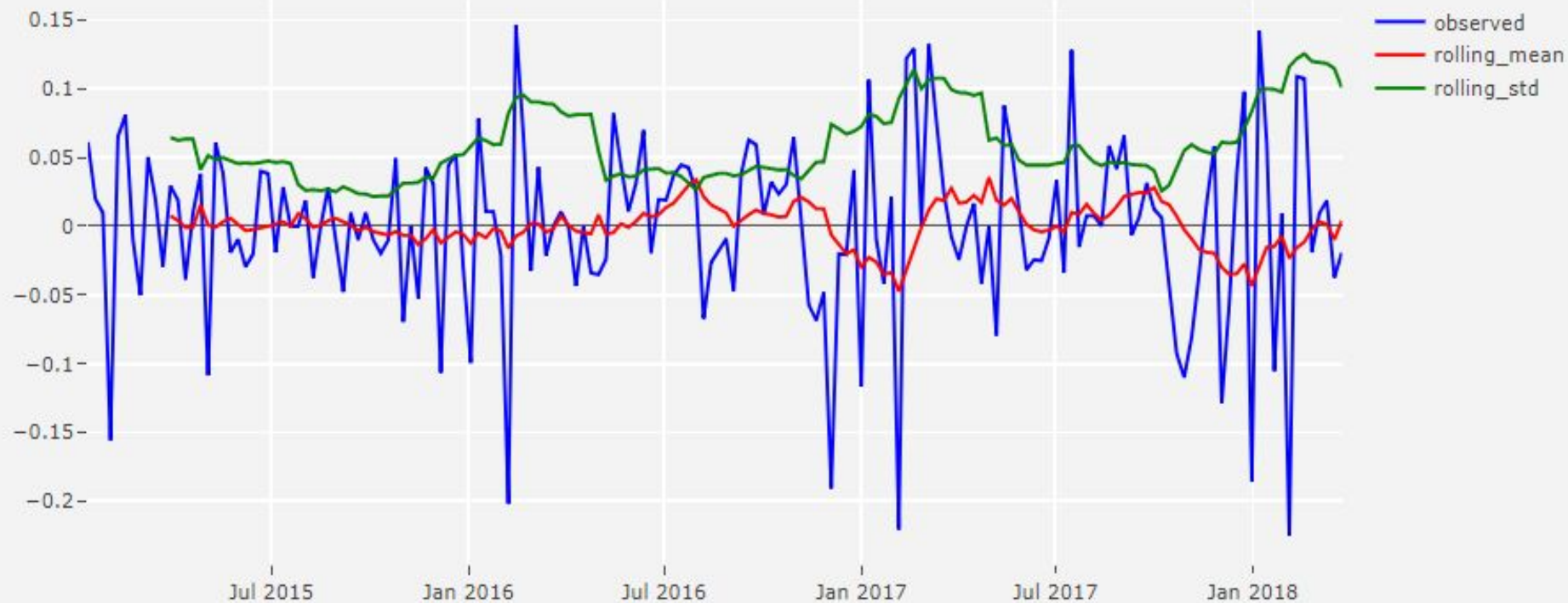
- Trend: mean is not constant over time
- Seasonality: variance is not constant over time

➡ Correct for Trend and Seasonality in order to make TS stationary

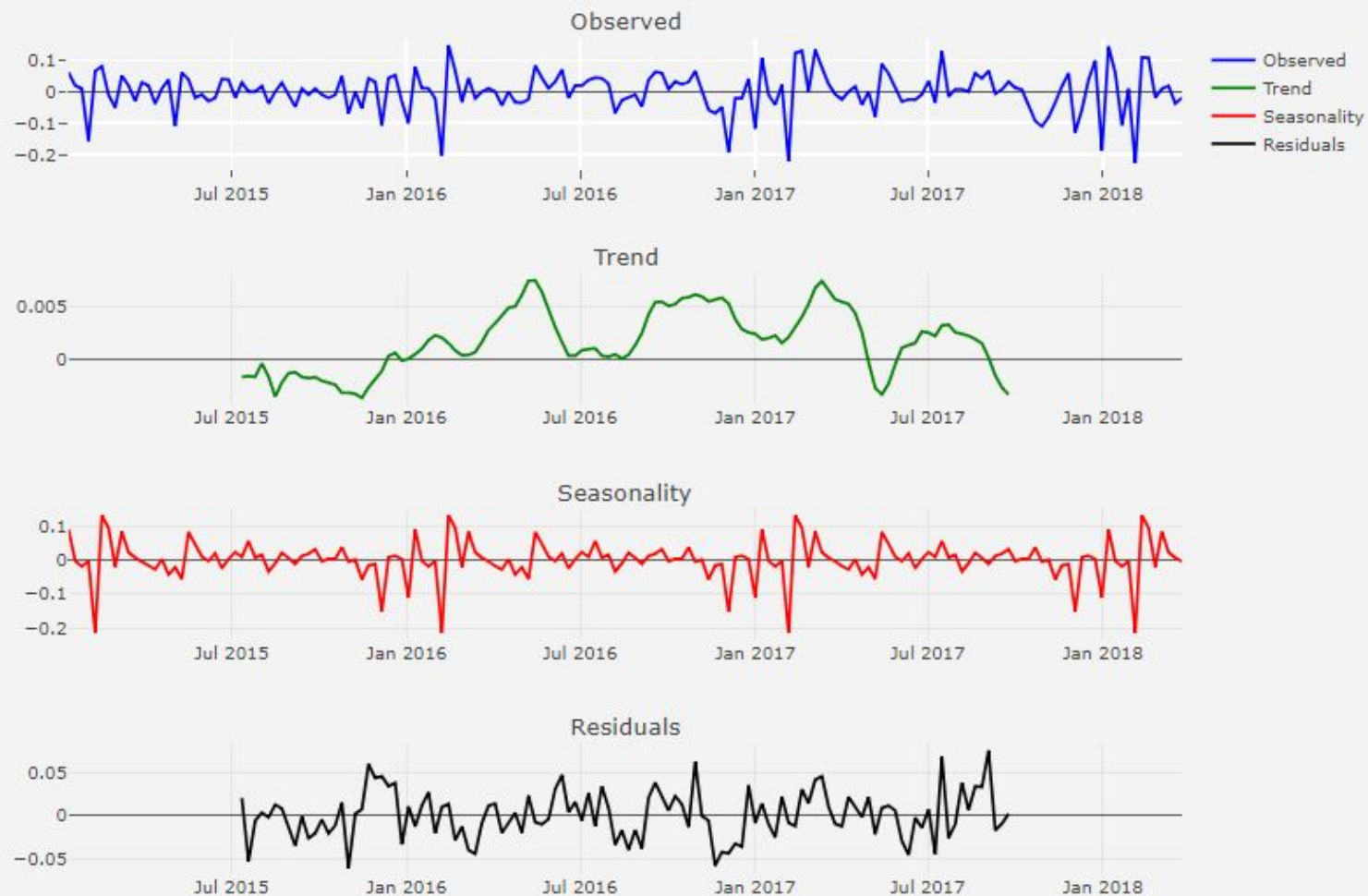
Possible solutions:

- Taking the Log
- Smoothing (rolling averages, e.g. weekly, monthly, ...)
- Differencing
- etc.

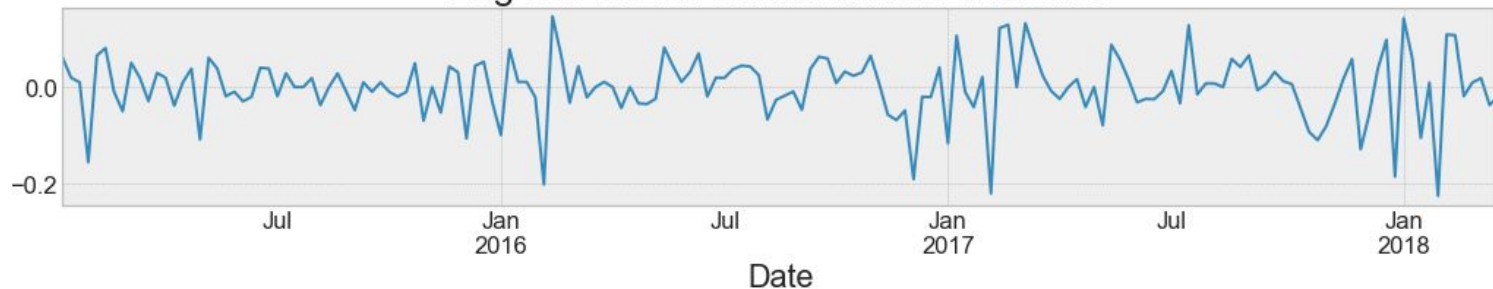
Rolling Mean and Standard Deviation for Log-Diff Time-Series



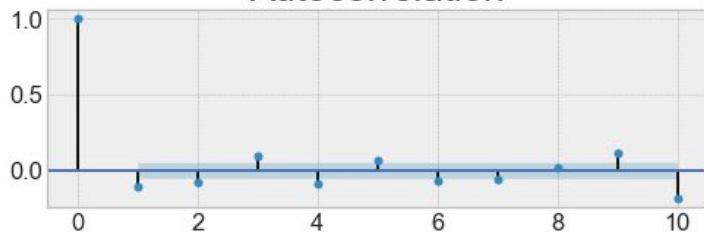
Decomposition of Log-Diff Time-Series



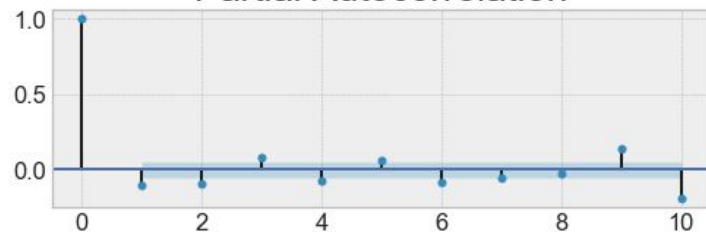
Log-Diff TS for Conventional Avocados



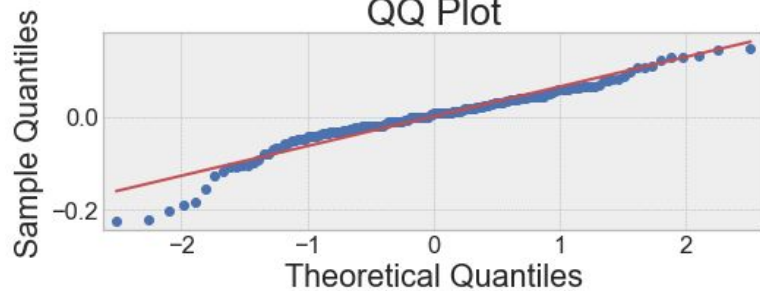
Autocorrelation



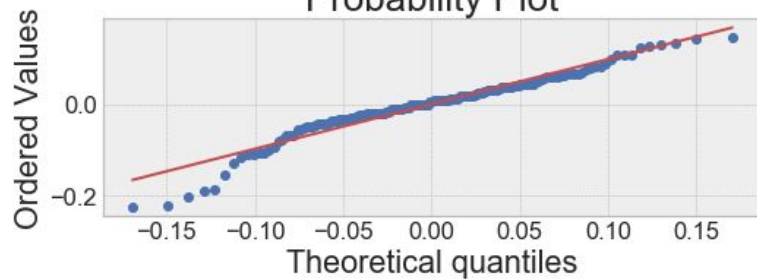
Partial Autocorrelation



QQ Plot



Probability Plot

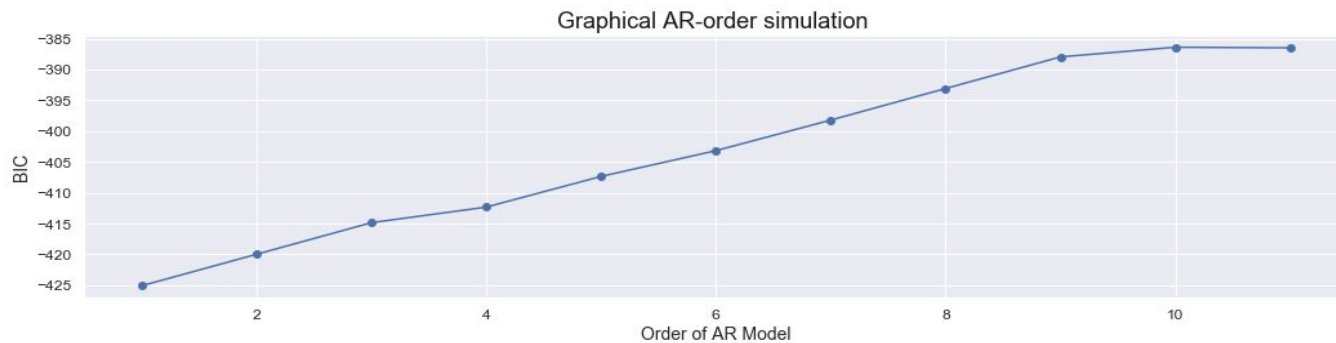
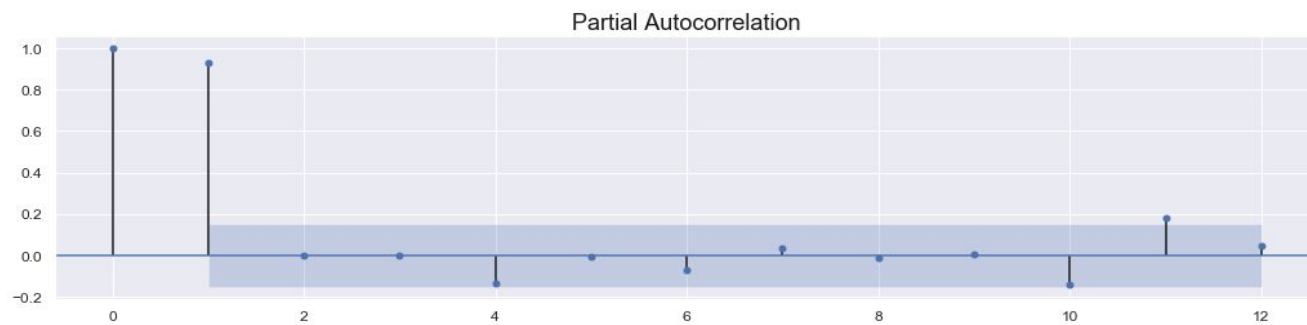
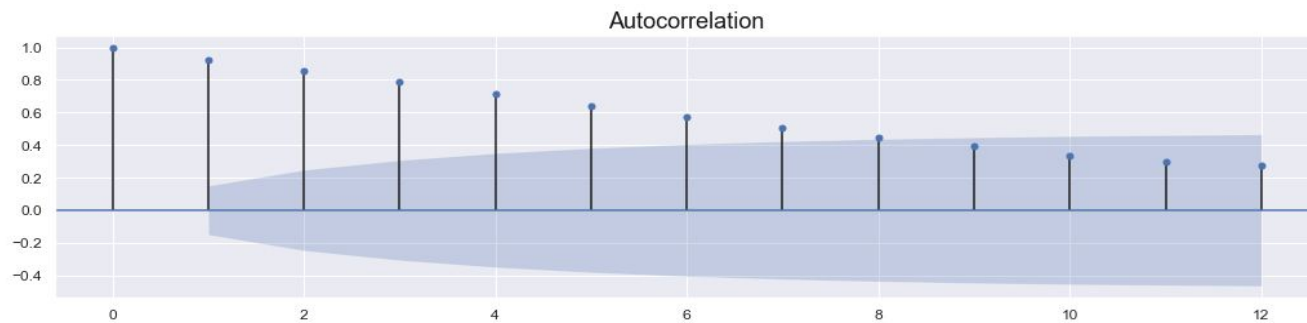


Results for conducted tests on Log-Diff-Time-Series

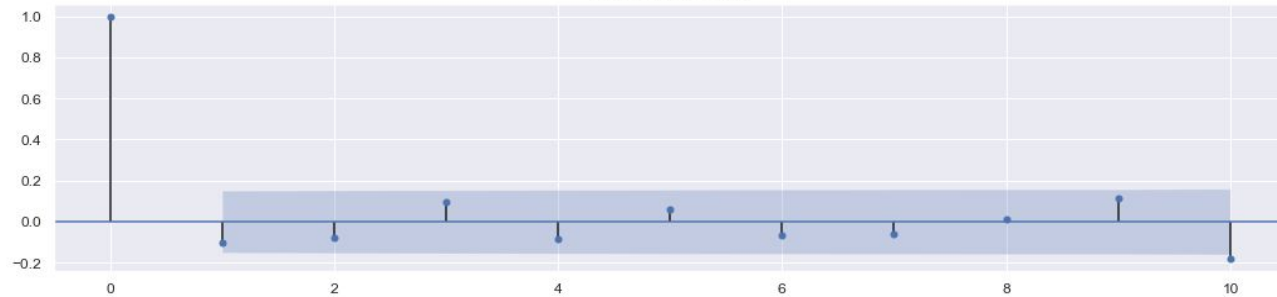
Test-Statistic	p-value
Augmented-Dickey Fuller Test	0.000144
KPSS	0.10



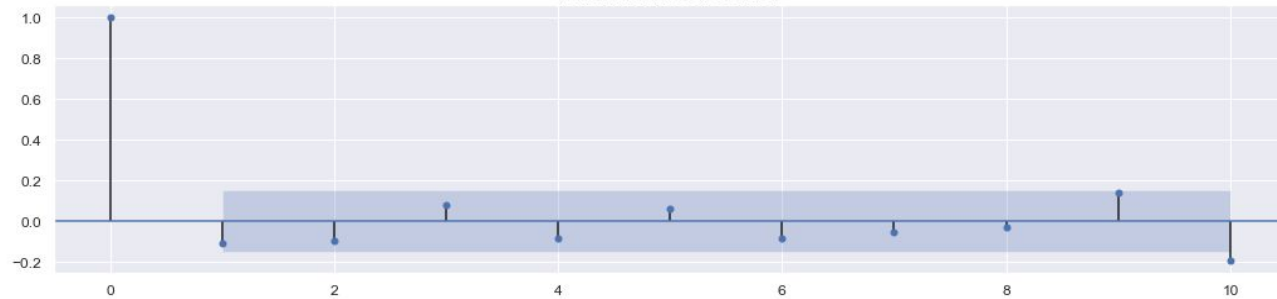
Assumption: Time-Series is stationary



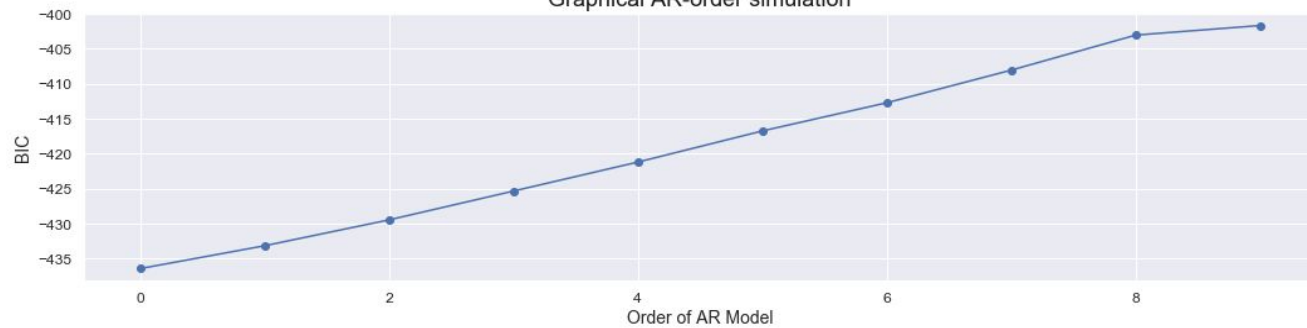
Autocorrelation



Partial Autocorrelation



Graphical AR-order simulation

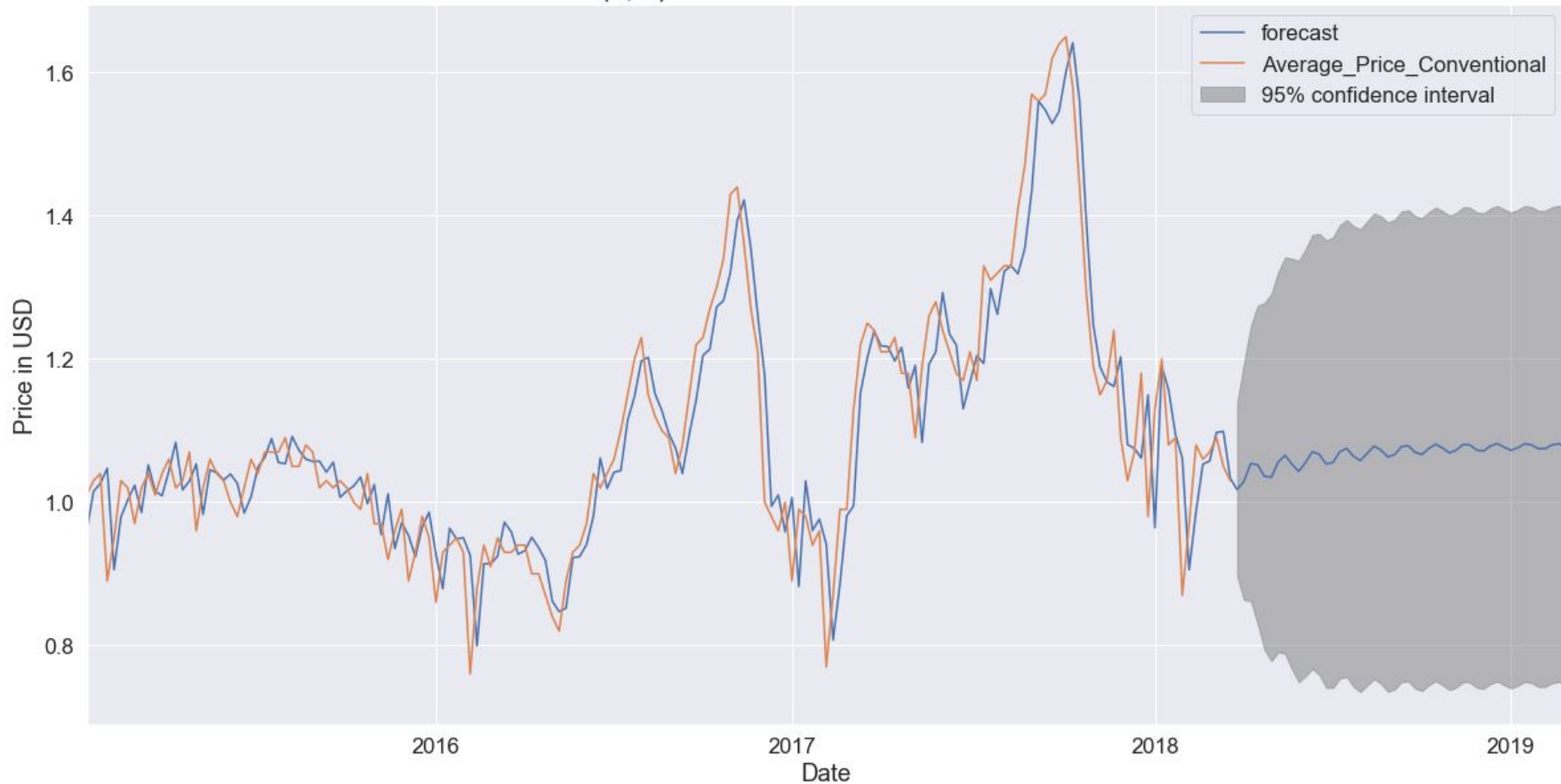




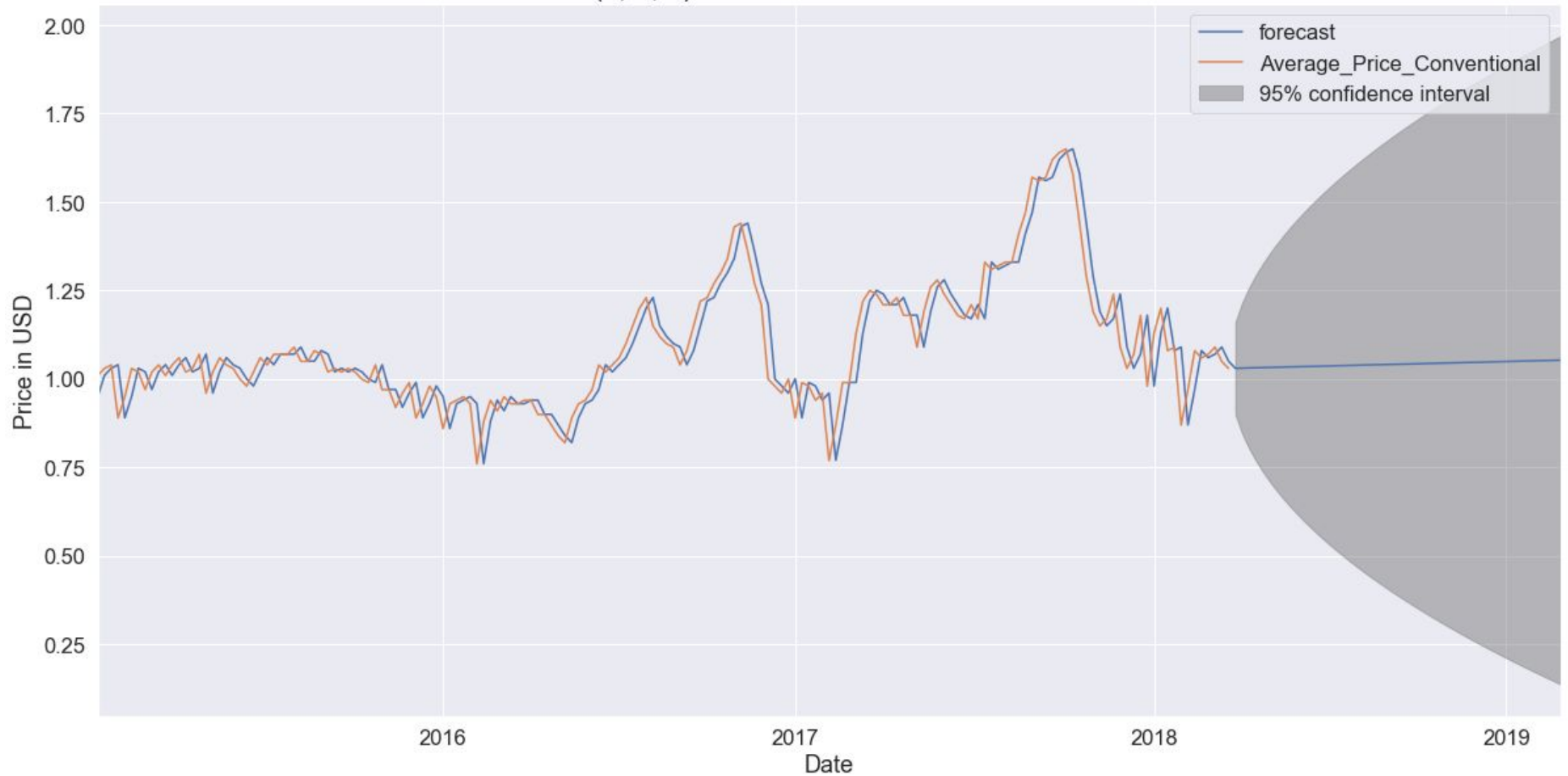
5.1 Model-Fitting: Traditional approach

- since model-fits for log-diff TS yielded in bad results, we tried to fit a model for the original TS
- following fits were conducted:
 - ARMA (4, 3)
 - ARIMA (0, 1, 0)
 - SARIMA (1, 1, 0) x (1, 1, 0, 52)
 -
- model evaluation with statistical metrics:
 - RMSE (= Root Mean Square Error)
 - MAPE (= Mean Absolute Percentage Error)
 - MAE (= Mean Absolute Error)

ARMA (4, 3) Forecast for Conventional Avocados



ARIMA (0, 1, 0) Forecast for Conventional Avocados





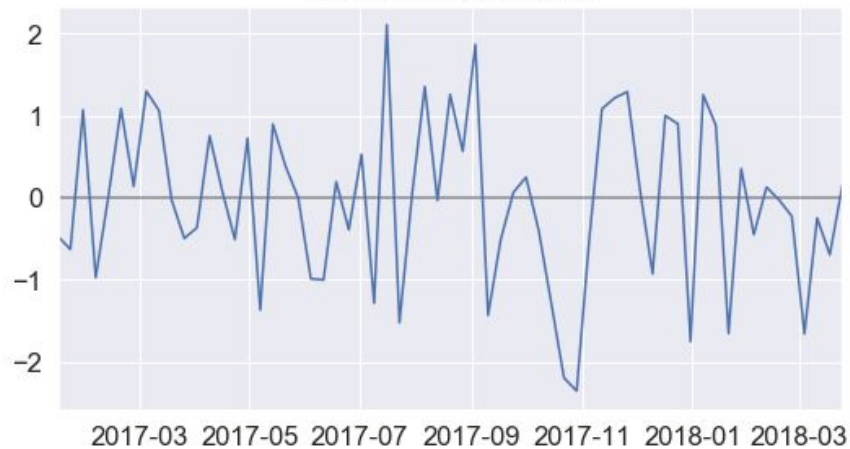
Short Detour: SARIMA

- ARIMA can model data with a trend, but not with a seasonal component
- extension of ARIMA: **Seasonal Autoregressive Integrated Moving Average** (= SARIMA)
- configuration of SARIMA:
 - Trend Elements:
 - p: Trend autoregressive order
 - d: Trend difference order
 - q: Trend moving average order
 - Seasonal Elements:
 - P: Seasonal autoregressive order
 - D: Seasonal difference order
 - Q: Seasonal moving average order
 - m: Number of time steps (e.g.: 12 for monthly data, 52 for weekly data etc.)

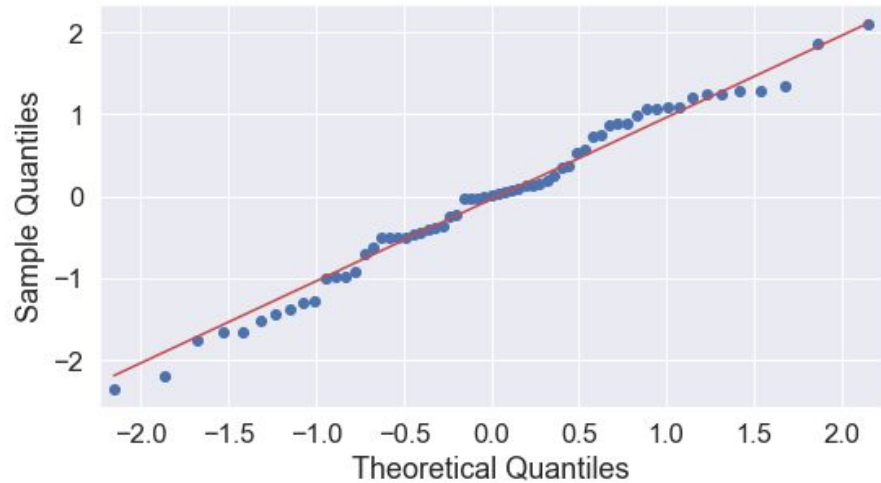
Average, Predicted and Forecasted Price of Conventional Avocados with SARIMAX (1, 1, 0)x(1, 1, 0, 52)



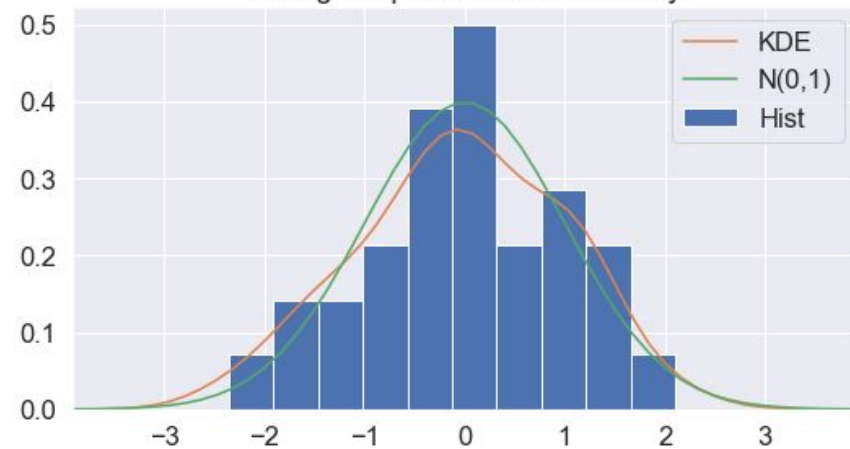
Standardized residual



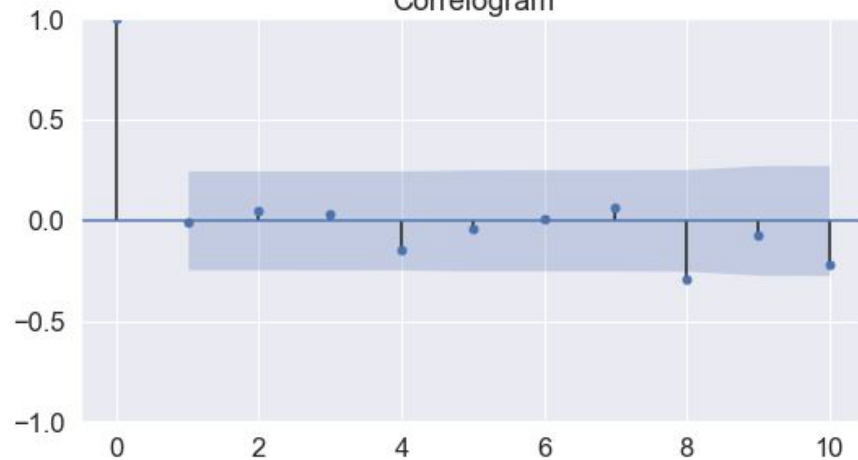
Normal Q-Q



Histogram plus estimated density



Correlogram

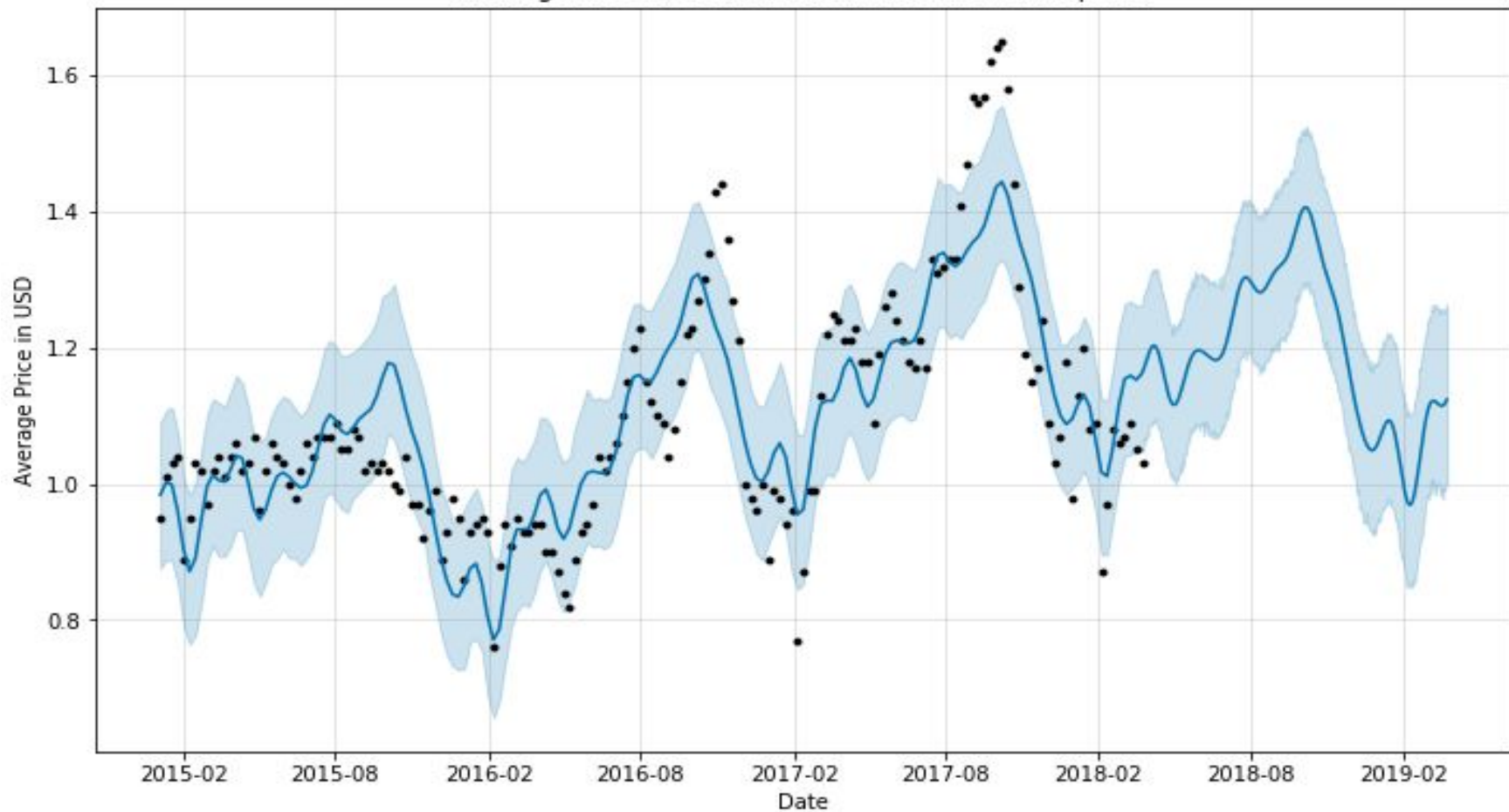




5.2 Forecast with FBProphet

- created by the Core Data Science Team (> 650 Data Scientists) at *Facebook*
- developed for Python and R
- library was open-sourced in 2017
- forecasts for TS-data based on additive models where non-linear trends are fit with yearly, weekly, and daily seasonality, plus holiday effects
- **does not require manual pre-processing steps** (e.g. differencing to make the data stationary)
- works best for TS which has strong seasonal effects and several seasons of historical data

Average Price of Conventional Avocados (FBProphet)





5.3 Comparison

	RMSE	MAE	MAPE
SARIMA	0.109	0.055	0.057
FBProphet	0.086	0.0675	0.061

Comparison betw. SARIMAX and FBProphet Forecast





6. Conclusion

- First Research Question:
 - Yes, as clearly seen from the graphs, 2017 was *the year* for avocados
 - Average price and volume of sold avocados rose tremendously
- Second Research Question:
 - Price range for predicted avocados between 2018 and early 2019: 0.80 - 1.10 USD



Sources

URL: <https://www.kaggle.com/neuromusic/avocado-prices>, Zugriff 27.6.2019

URL: <https://hassavocadoboard.com/>, Zugriff 27.6.2019