

Impact of Drought on Maple Tree Production

BU-NOAA Capstone Project Mid-Term Summary Report



Team 4B: Chris Chang, Jacinto Lemarroy, Mengxin Li, Shiyu Ye, Ying Li

Introduction and Business Problem

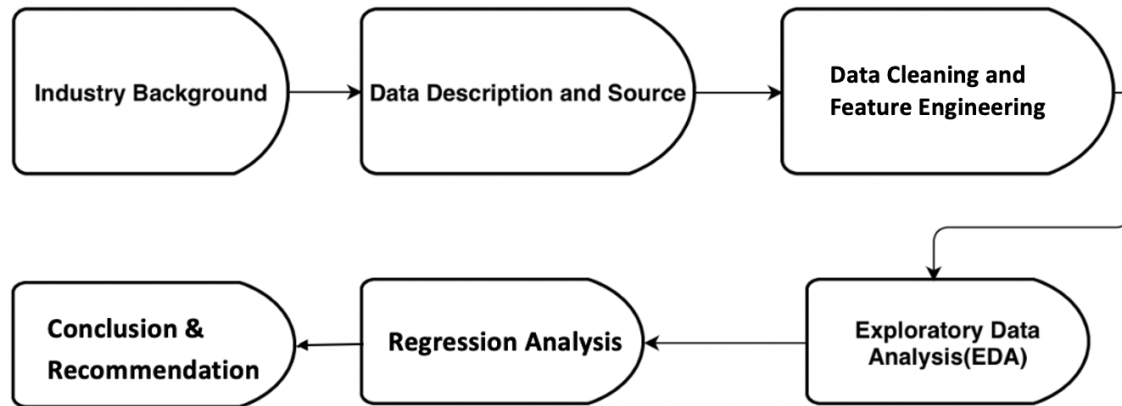
Maple and the maple industry are synonymous with the New England region's sugar houses and mountainsides with colorful leaves in the fall. The economic value of maple syrup hit \$130 million back in 2015(1.). However, the maple syrup industry is still rapidly growing. According to ReportLinker, the global maple syrup market is expected to grow by \$402.61 million during 2021 to 2025, progressing at a CAGR of 6%(2.). Beyond producing maple products, the maple industry contributes to the maple timber and fall foliage tourism industry, making it the fourth most valued agricultural commodity and the second most valued crop, closely behind greenhouses and nurseries (3.). However, the growth of maple trees and maple syrup production is sensitive to climate change. One main concern for maple tree growers and maple syrup producers is the impact of warming spring temperature during the critical sugaring period and increase in summer drought frequencies which affect the maple syrup yield and sugar content, and the flavor of maple products. Understanding how climate changes and other factors may impact this resource is essential to continue the management of the maple industry into the future.

This project aims to determine how maple syrup production in gallons is affected by drought and precipitation factors. Our primary focus will be on the Northeast region, precisely seven states, including Massachusetts, Maine, Vermont, New Hampshire, New York, Connecticut, and Pennsylvania. The New England region is well known for its maple syrup production, and hence we will delve deeper into that area. We will analyze variables such as the annual crop production of maple syrup and other parameters related to climate change. Furthermore, this report will further explain the impact of drought on maple tree production.

Methodology

Following the structure below, our report will first give industry background, followed by data description and source. Then we will focus on data cleaning and manipulation, exploratory data

analysis, and finally, we will expand on the next stage involving regression analysis. We considered machine learning methods, but we opted not to do so due to the lack of data.



Industry Background

The main objective of this project is to understand the factors affecting maple syrup production. Organizations such as NOAA and third-party websites already provide publicly available data regarding drought parameters for all states. In addition, there are vast information sources regarding drought and the different industries affected. We considered many sectors such as agriculture, manufacturing, energy, and tourism. Nonetheless, given that our team is based in Boston, we decided to focus on a sector that would be most relevant to our location.

After researching the aforementioned sectors, we arrived at three final project ideas. The first idea was how drought affects business in golf courses in New England. Next, we thought about how drought conditions affect hay prices for livestock in New England. The third idea was drought-affected maple tree production in New England. Given the data scarcity regarding the other sectors and the relevancy of the project idea, we concluded that agriculture data regarding maple tree production would be the most suitable option.

Initially, we brainstormed ways to merge the drought data with the maple tree production data. Then, we understood that the correlation between drought variables (E.g., DSCI, D0, D1, D2, D3, D4) and maple tree variables (E.g., gallons, value) could indicate a pattern that may paint an informative picture for us. After learning more about a potential correlation between drought and maple tree production, we could then be able to confidently create prediction models of maple production based on specific drought levels.

Data Description and Source

We collected valuable data from different sources to examine how the drought factors affect maple syrup production. For example, we collected data from official websites such as the National Oceanic and Atmospheric Administration (NOAA), the United States Department of Agriculture (USDA), and the National Drought Mitigation Center. We also reached out to organizations such as the New England Forestry Foundation to retrieve more information on maple syrup data, such as sap to syrup ratio and the number of tree taps. However, due to the limited and inconsistent data samples provided, we decided not to use the maple dataset for more accurate analysis.

Our dataset has 210 rows and 22 variables of different data types such as string and integer. It is based primarily on seven states (Massachusetts, Maine, Vermont, New Hampshire, New York, Connecticut, and Pennsylvania) from 1992 to 2021.

The attributes (columns) to focus on in our data analysis are:

1. **State:** State Name (Massachusetts, Maine, Vermont, New Hampshire, New York, Connecticut, and Pennsylvania)
2. **Years:** 1992 to 2021
3. **State ANSI:** American National Standards Institute (ANSI) Codes for States
4. **Gallons:** Maple Syrup Production in Gallons
5. **Value:** Maple Syrup Production in Dollars
6. **Average Temperature:** Average Temperature (Fahrenheit)
7. **Average Temperature (6-month average):** Average Temperature from September to February next year based on the growing period of maple trees
8. **Precipitation:** Average Precipitation (millimeters)
9. **Precipitation (6-month average):** Average Precipitation from September to February next year based on the growing period of maple trees
10. **None to D4:** Percent Area in U.S. Drought Monitor Categories
 - a. None: Absence of drought
 - b. D0: Abnormally dry
 - c. D1: Moderate Drought
 - d. D2: Severe Drought
 - e. D3: Extreme Drought
 - f. D4: Exceptional Drought
11. **DSCI (Drought Severity and Coverage Index):** An experimental method for converting drought levels from the U.S. Drought Monitor map to a single value for an area.
Formula: $1(D0) + 2(D1) + 3(D2) + 4(D3) + 5(D4) = DSCI$

12. **SPI (The Standard Precipitation Index)**: A relatively new drought index based only on the probability of precipitation
 - a. SPI drought: the degree of dry condition
 - b. SPI wet: the degree of wet condition
13. **EMNT (extreme minimum Temperature in air)**: Extreme minimum Temperature in the air throughout the year
14. **EMXT (extreme maximum Temperature in air)**: Extreme maximum Temperature in the air throughout the year
15. **Number of Taps**: The number of taphole for producing saps
16. **Yield per Taps**: The average sap yield (per gallon) for a taphole

Data Sources:

- Temperature & Precipitation:
<https://www.ncdc.noaa.gov/cag/statewide/time-series>
- Gallons & Value:
https://quickstats.nass.usda.gov/?long_desc__LIKE=maple+syrup&x=0&y=0
- DSCI & D0-D4:
<https://droughtmonitor.unl.edu/DmData/DataDownload/ComprehensiveStatistics.aspx>
- SPI:
<https://www.drought.gov/historical-information?dataset=1&selectedDateUSDM=20101221&selectedDateSpi=19580601>
- EMNT & EMXT:
<https://www.ncdc.noaa.gov/cdo-web/review>
- Number of Taps & Yield per Taps:
https://www.nass.usda.gov/Statistics_by_State/New_England/index.php

Data Preprocessing and Feature Engineering

As mentioned above, our dataset consists of 210 rows and 22 variables. However, there are some missing values in 64 rows of the dataset. The drought-related variables (none, D0, D1, D2, D3, D4, and DSCI) are unavailable for all states from 1992 to 1999. Considering the goal of this project is to examine the effect of drought on maple syrup production, we decided to combine existing data with another dataset and use machine learning to fill in missing values for most of the drought-related columns (none, D0, D1, D2, D3, D4). In the end, we decided to use the SPI dataset from USDA to make predictions. The SPI dataset has information about the standard deviation of drought-related information, such as standard deviation of drought level 0 (D0) and standard deviation of wet level 0 (W0) for northwestern states. After finalizing the input data, we used linear regression and random forest to make predictions. Linear regression outperforms the

random forest model since the dataset only has 120 rows of training data with a train/test split of (0.8:0.2). The accuracy score for the model is shown below:

	None	D0	D1	D2
Train	0.864	0.732	0.719	0.643
Test	0.811	0.622	0.683	0.324

As the drought level increases from D0 to D2, the linear regression model becomes insufficient for making predictions. This is because severe drought is relatively rare, and the standard deviation of the drought variable is prone to changes with the tiny fluctuation in a single entry of data. We do not have enough data to address the overfitting issue currently. Since drought level 4 (D4) is 0 for all entries, we assumed D4 was also between 1992 and 1999. After making the assumption, we calculated the D3 value by subtracting all drought-related variables by 1 ($D3 = 1 - \text{None} - D0 - D1 - D2 - D4$). We put the negative predicted value to 0 before finalizing our dataset. Once we calculated the values for D0 to D4, we computed a weighted sum of the percentages in D0 through D4 to calculate DSCI. We used the following formula to compute DSCI: $1(D0) + 2(D1) + 3(D2) + 4(D3) + 5(D4) = \text{DSCI}$

We also checked the outliers for each independent variable. From the boxplot (Appendix 1), we can see our dataset has some outliers in a statistical sense, but they are actual data values. The variables such as precipitation and drought level are unpredictable and vary year by year, so we did not deal with those "outliers" for the further regression analysis.

Exploratory Data Analysis

Overall production of each state

We first look at the overall production of the seven states in the Northeast. From the production time series chart, Vermont, New York, and Maine are the states that produced the most maple syrup in gallons. We missed the value of gallons for Massachusetts and Connecticut from 2019 to 2021 because the USDA no longer collects and reports the data. (Appendix 2)

Also, maple syrup production decreased dramatically in 2010, 2012, and 2021. According to New England Agriculture Statistics, "the 2012 maple syrup season in New England was considered too warm. In March, many heatwaves ended the season for many and resulted in a significant drop in maple syrup production. An exception was Maine, where temperatures were cool enough in top-producing". Similarly, low production in 2021 was due to poor weather, which included warm weather in March that ended the season early, and low sugar content meant producers had to boil more sap to make a gallon of syrup.

Overall correlation heatmap

We created a correlation heatmap based on all the variables. The chart shows that gallons, value, and number of taps are highly correlated and have a positive relationship. At the same time, there is a moderate negative correlation between average temperature and both gallons and value. We need to consider the multicollinearity issue for further machine learning since all the drought-related variables such as precipitation, SPI-drought, SPI-wet, DSCI, and D0 to D3 were highly correlated. (Appendix 3)

Average Temp on Gallons

We found a negative correlation between average temperature and gallons. This supports a study from UW-Extension that warmer weather stops enzymes in tree roots to stop producing sugar in tree sap and could also lead to tree root damage, leading to a reduction in sap flow which leads to a decrease in gallon production. These two features have a slope of -3.84 and an R-Squared value of 18%, which means that the annual average temperature can explain 18% of the change in gallons. (Appendix 4)

Average Temp on Value

There is a negative correlation between average temperature and value, which makes sense since a reduction in gallons would reduce the total monetary value for syrup manufacturers. These two features have a slope of -1.10 and an R-Squared value of 15%, slightly lower than that between average temperature and gallons. This makes sense as gallons would be more directly correlated with temperature than monetary value temperature. (Appendix 5)

Extreme Maximum Temperature in Air on Gallons

Our analysis shows no clear relationship between extreme maximum temperature and gallon production of maple syrup. The majority of the values of gallons production are nearing zero production, which is expected given the extreme temperature in the air. (Appendix 6)

Extreme Minimum Temperature in Air on Gallons

Our analysis shows a very slight negative relationship between extreme minimum temperature and gallon production of maple syrup. The majority of the values of gallons production are nearing zero production, which is expected given the extreme temperature in the air. (Appendix 7)

Average Precipitation on Gallons

We made seven line charts comparing production in gallons and precipitation for the relationship between temperature and gallons. (Appendix 8)

For some states like Maine and New Hampshire, precipitation reduction does not necessarily mean an intense reduction in production. For example, during 2016, the precipitation in both

Maine and New Hampshire decreased dramatically while the output of those states during the year increased adversely. Other factors could affect maple syrup production, such as warmer temperatures and the trees. According to the USDA's 2016 maple syrup production report, 'Producers were encouraged to tap earlier this season by the warmer than normal temperatures.'

After considering maple trees' growing season, we found that the coefficient for 6-months-precipitation overproduction in gallons regression is almost six times that of months. Also, it is statistically significant (under the confidence level of 95%). Thus, we believe that precipitation is crucial to maple syrup production, especially during the growth of maple syrup.

SPI drought (including D0-D4)

Detailed into the level of drought, we first discovered D0 (abnormally drought area). The drought level versus production rate of each state (Maine as an example) is listed in Appendix 9.

During the early 2000s, there was an inverse relationship between D0 (Abnormally dry) levels of drought and gallon production. Over time, that relationship became less inverse to the point where even with high levels of D0, the drought did not impede high levels of maple syrup production. One should not generalize the relationship between D0 levels and maple syrup production in states in the Northeast, as not all states show the same behavior. Except for Massachusetts and Connecticut, D0 and maple syrup production seem to rise for most Northeast states.

Like D0 and production, the D1 (Moderate Drought) level also started inverse to production in the early 2000s. However, there was no clear relationship between the mid-2000s and late 2010s as both variables increased and decreased simultaneously. (Appendix 9 Maine as an example.)

During the early 2000s, D2 (Severe Drought) was high in most states. During the middle and late 2000s, D2 levels were essentially 0. During the mid-2010s, there was a spike in D2 levels again. In the late 2010s, there seemed to be an upward trend in D2 levels for most states, except Massachusetts and Connecticut. (Appendix 9)

It is difficult to generalize the relationship between D3 (Extreme Drought) levels and maple syrup production in the Northeast. Many states show even more different behavior in terms of units of D3 levels. For some states, D3 levels were essentially 0 for several years. One could argue that the reason for D3 levels to be minimal for many states is because extreme drought is not as common as lesser levels of drought in the Northeast. (Appendix 10)

Regression Analysis

We mainly ran regressions through each state and of all states together. And we first did linearity checks to try to eyeball the dataset. For some variables such as 'Yield per Taps' and 'Number of Taps,' the relationship with our target variable 'Gallons' is linear. This is because maple syrup production is calculated by multiplying yield per tap by the number of taps. (Appendix 11). However, the linearity relationship is not obvious for other variables such as 'SPI Wet' and 'Avg

Temperature.’ Further regression analysis should be made to investigate the real relationship between them. (Appendix 12)

1. Regressions for each state

We first tried running regression among each state. Since our dataset includes the temperature and precipitation of both the whole year and the half year(growing seasons for maple trees from September to January next year), we ran two regressions to see if considering maple trees' growing season does improve the regression results.

The R square of each state is listed in the appendix. (Appendix 13) Surprisingly, in most states, using precipitation during maple trees' growing season increases the R-squared value. Thus, in conclusion, we will use this finding to continue regressions on all states.

State	R-square	R-square(growing season)	Difference
ME	52.80%	52.60%	-0.20%
VT	76.40%	82.70%	6.30%
MH	80.10%	80.30%	0.20%
NY	57.50%	58.00%	0.50%
PA	37.60%	41.70%	4.10%
CT	49.90%	42.90%	-7.00%
MA	56.60%	56.60%	0.00%

After the above effort, we eliminated several insignificant variables that returned very high P-values in the previous regression. After checking the significance of these features, we decided to exclude D1, D2, D3, and D4 because of these variables’ collinearity with the D0. In addition, the VIF factor table listed below showed that any factor larger than 10 should be eliminated. However, we kept D0 and DSCI after several attempts for better regression results.

VIF Factor	features
1.810565e+11	const
3.800000e+00	Avg.Temperature_one_year
4.600000e+00	Avg.Temperature_six_month
1.140000e+01	Precipitation_one_year
8.800000e+00	Precipitation_six_month
8.530393e+09	none
inf	D0
inf	D1
inf	D2
inf	D3
NaN	D4
inf	DSCI
1.420000e+01	SPI-Drought
1.760000e+01	SPI-Wet
4.100000e+00	EMNT(in air)
3.700000e+00	EMXT(in air)
4.700000e+00	Number of Taps
6.300000e+00	Yield per Tap

After eliminating distracting variable factors, the R squared decreased slightly, but some previously insignificant variables (D0) became statistically significant. We see this R-squared decrease in states such as Maine. The SPI drought stayed substantial as well. Thus we concluded that on the state level, SPI drought, which measured the deviation of the drought effect and D0(abnormally dry) could affect the maple syrup production. However, the coefficient of D0 is negative, which indicates that in some way, appropriate drought is indeed a good sign for the production of maple syrup. The coefficient of the variable 'none,' meaning the percentage of the land that did not suffer from drought, is negative. However, the large P-value showed that this could be due to our limited number of data points, and we should not be concerned about that. (Appendix 14)

2. Regression on all states (Multivariate)

a. Multivariate Regression

The combined regression (Appendix 14) indicates that our independent variables explain around 44% of the variability of the maple syrup production in gallons. Moreover, our P-values show that only average temperature, precipitation, SPI values, and extreme max temperature are significant.

3. Regression with Lag Effect

Since the growth and harvest of maple syrup is a continuous process throughout the years, we want to analyze the effect of previous temperatures and precipitation on maple syrup production. Therefore, we added the last year's or last six months' temperature or precipitation as the new feature. In the end, we found that the previous year's precipitation impacts next year's production.

	Lag effect column	R_squared
One_Year	None	0.459
	Temperature	0.370
	Precipitation	0.467
	Both	0.388
Six_month	None	0.317
	Temperature	0.301
	Precipitation	0.303
	Both	0.272

Conclusion

This project aimed to understand and analyze the drought and precipitation determinants that affect maple syrup production in the Northeast states, specifically New England. This issue is essential since global warming has made it difficult for crop producers to obtain optimal production yields due to the extreme weather conditions they face each year. Even though there is more global warming awareness, fixing this issue will not be overnight. Therefore, governments must contemplate short and long-term strategies to alleviate the consequences of extreme weather conditions. Understanding drought and precipitation can help farmers, growers, and producers mitigate the risk of severe weather and precipitation conditions to prepare for the most efficient and effective harvest season and maximize yield production.

After running several different regression analyses, we came to a few conclusions:

- Average temperature has more correlation to gallons compared to precipitation to gallons. Specifically, the combination of average temperature (1-year) and precipitation (growing season) correlates most with production in gallons.
- Adjusted R-Squared almost always decreases when running the regressions with all independent variables.
- When removing variables that are not significant, adjusted R2 increased, but R2 decreased since we took out more variables.

- The lag effect indicates that incorporating the previous year's precipitation explains more variance in maple syrup production, but this is not the case for average temperature.

Limitations and Challenges

We ran into several obstacles along the way. Firstly, the overall regression model performance was insufficient since we have limited data available. In terms of dimension, our dataset did not include information about maple tree species, soil type, or estimation of planted maple trees. In terms of data volume, we do not have enough entries of data to run and tune the prediction model and more sophisticated machine learning models. We faced other data granularity issues when we attempted to add more depth and width to our data. The data sources had variables on different scales, levels, or years than our original dataset. We reached out to maple syrup connoisseurs, maple syrup companies, government agencies, and even farms, but we received little to no response in most cases.

Given these aforementioned circumstances, it was difficult for us to draw causal conclusions about how drought affects maple syrup production or make predictions about maple syrup production.

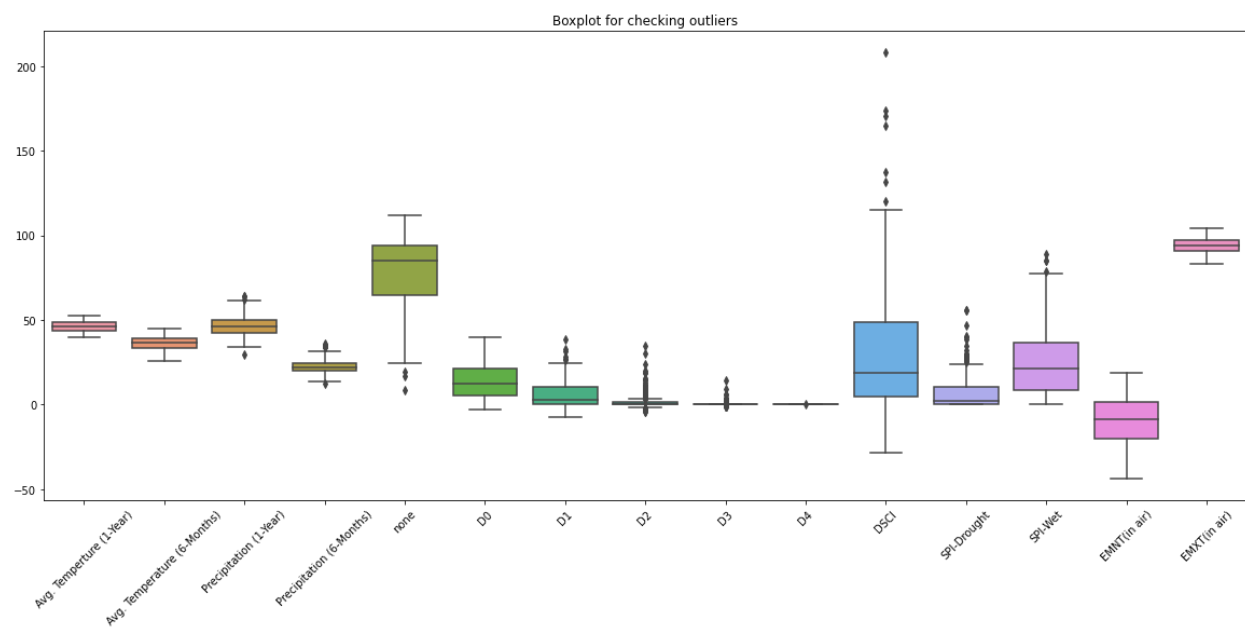
Potential Improvements

Moving forward, we thought of a few modifications that could have improved our analysis and our project overall:

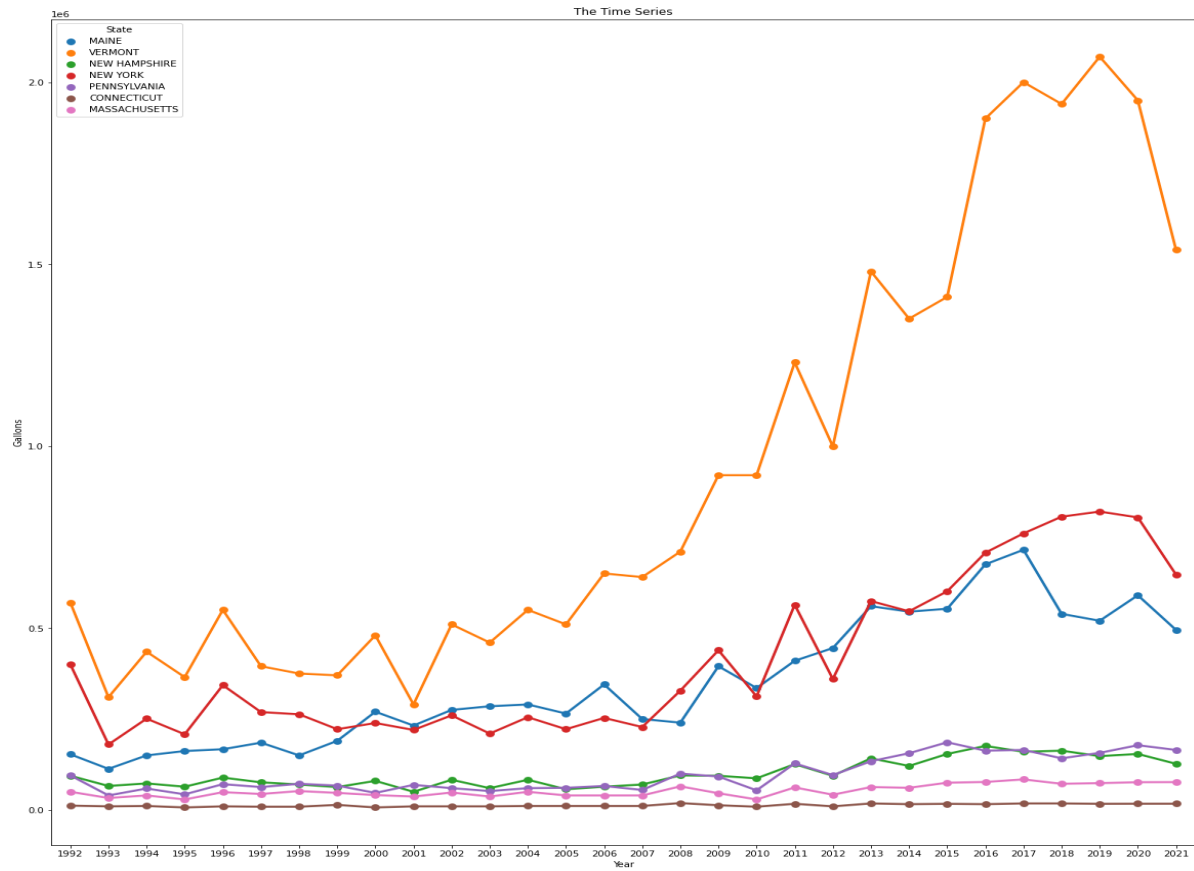
- Add more features (Soil, Geological variables, etc.) to improve the model.
- Find granular data (monthly, weekly, daily) to perform a more accurate analysis.
- Design dashboard (Tableau, PowerBI) for data visualization.
- Create a database to store and update data since our data is continuous.
- Expand drought impact analysis on other crops - e.g., apples, corn, soybeans.

Appendix

Appendix 1: Boxplot for checking outliers

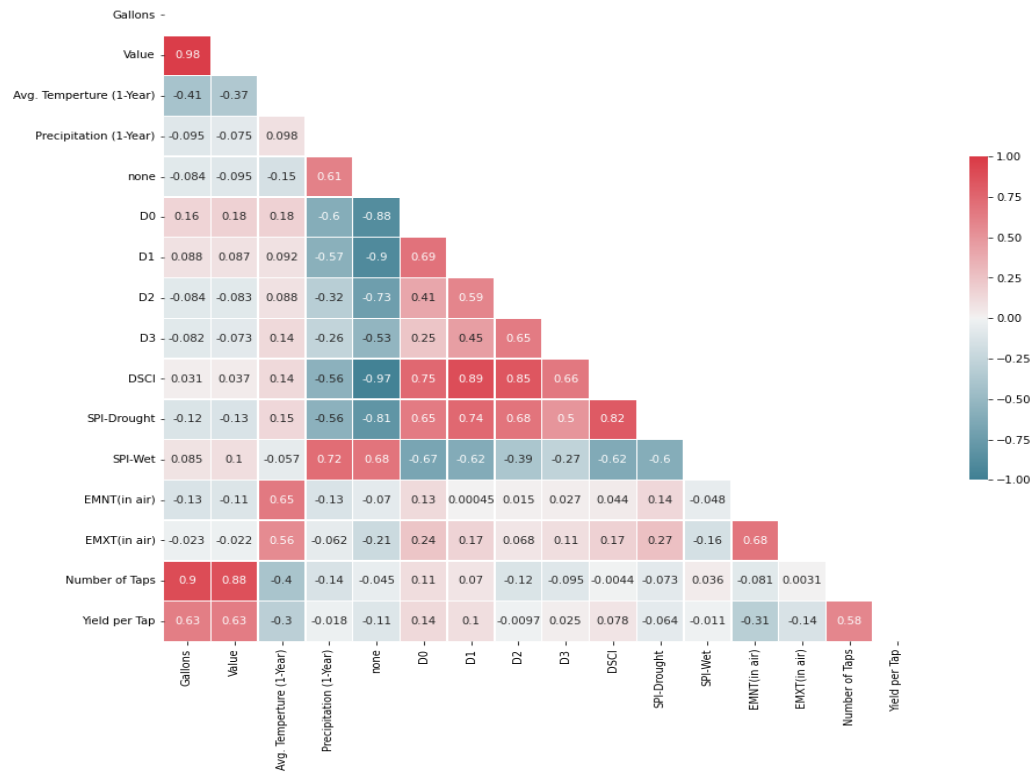


Appendix 2: Overall production of each states through 1992 to now

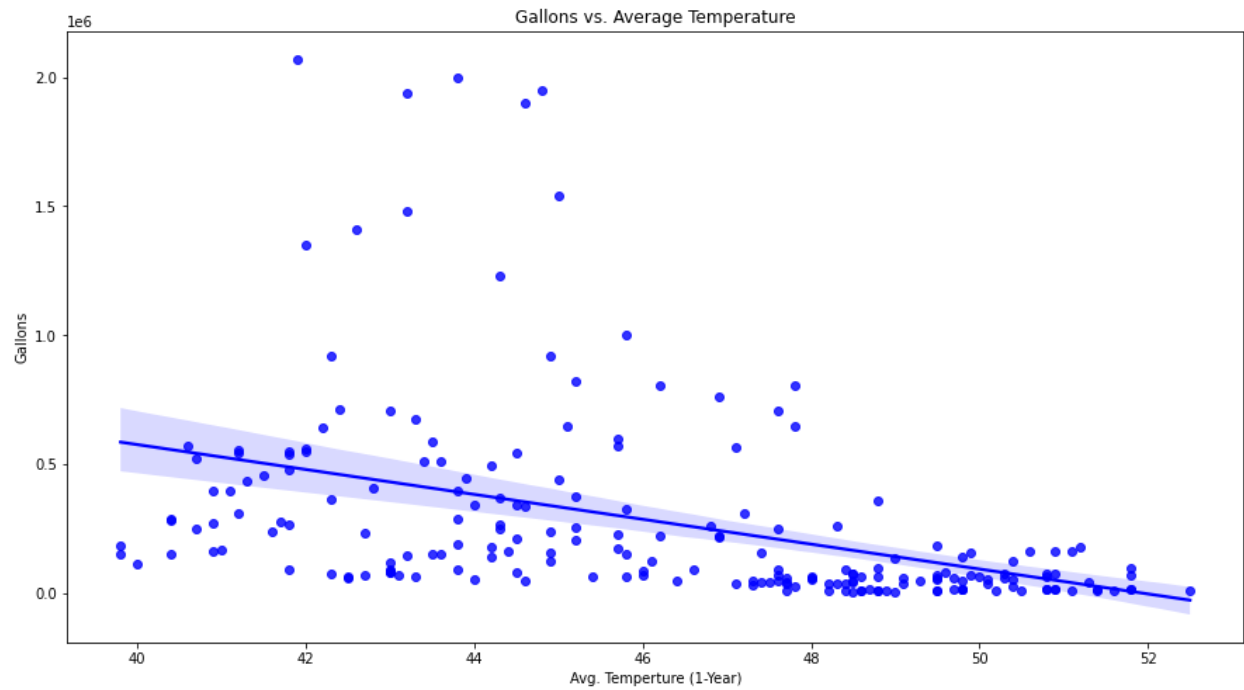


Appendix 3: Variable Heatmap

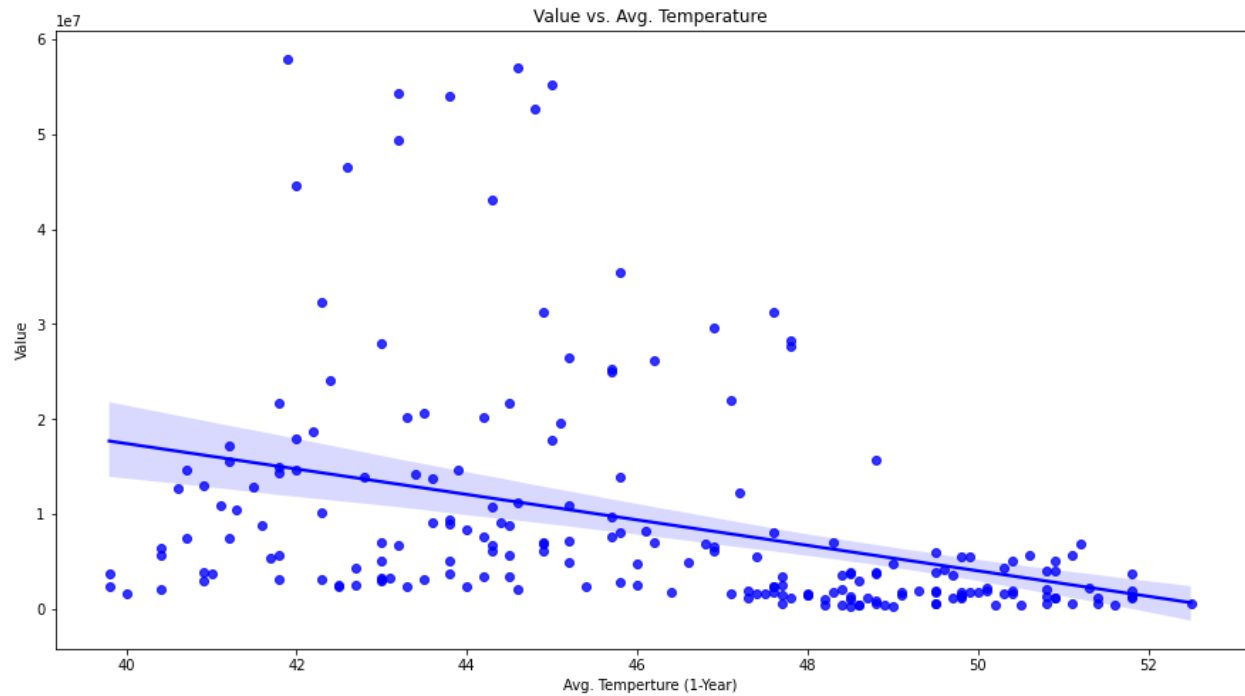
Heatmap-correlation matrix



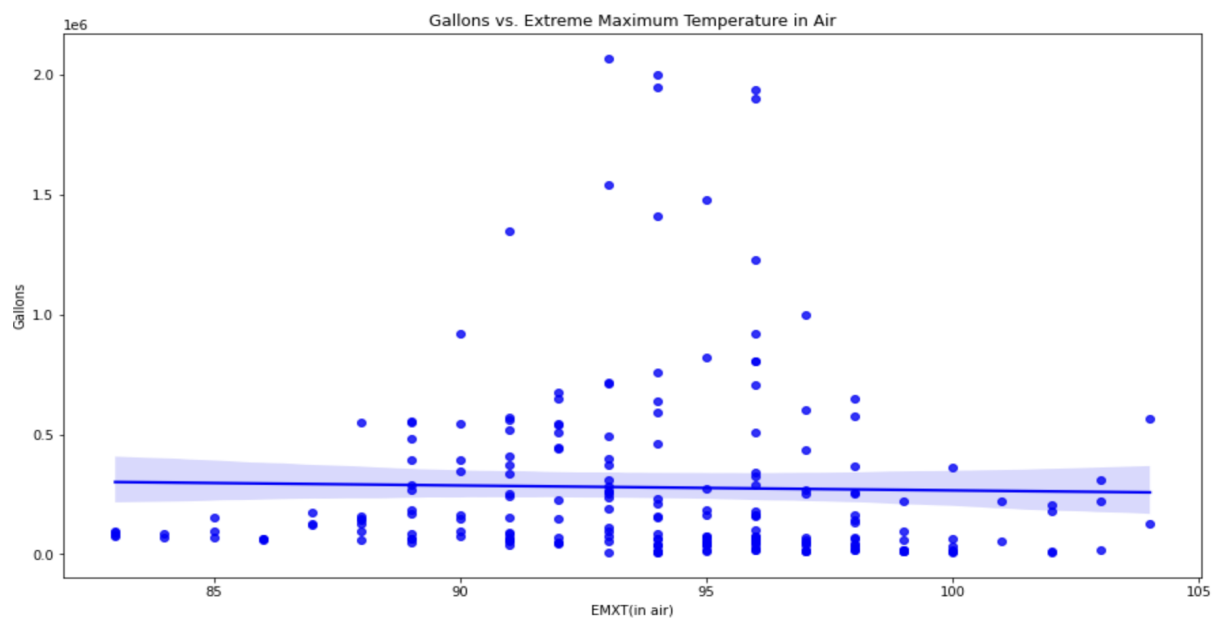
Appendix 4: Regression line of average temperature on production in gallons



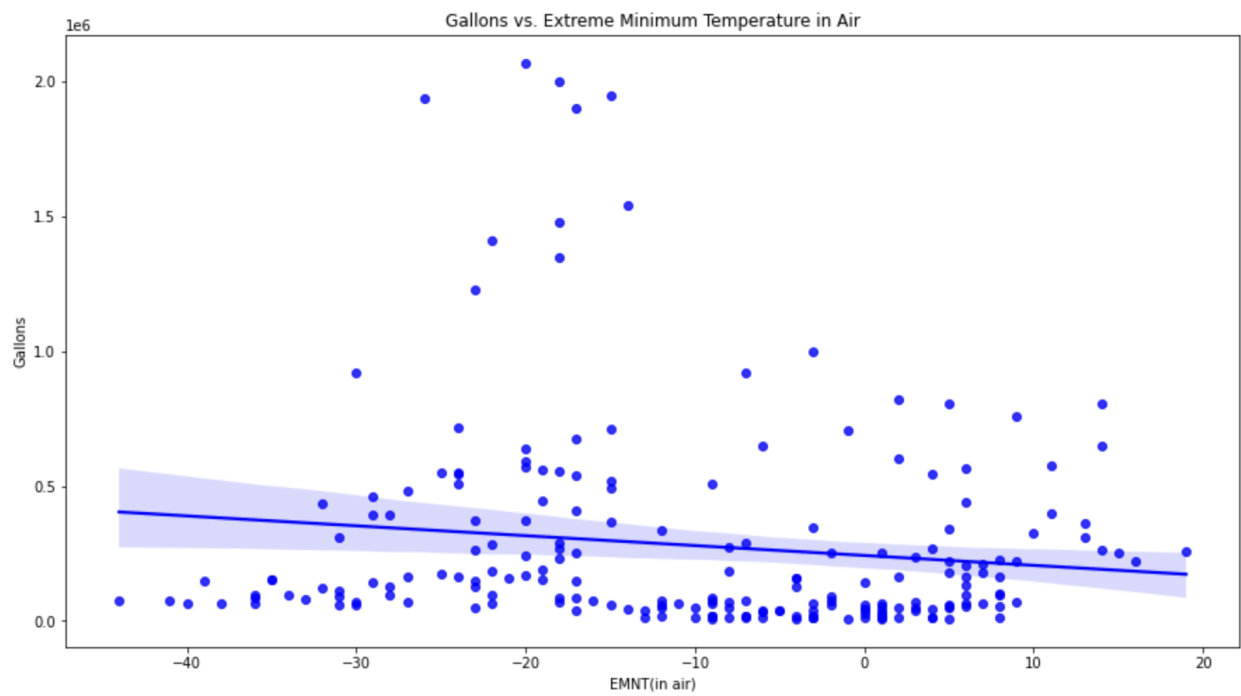
Appendix 5: Regression line of average temperature on production in values



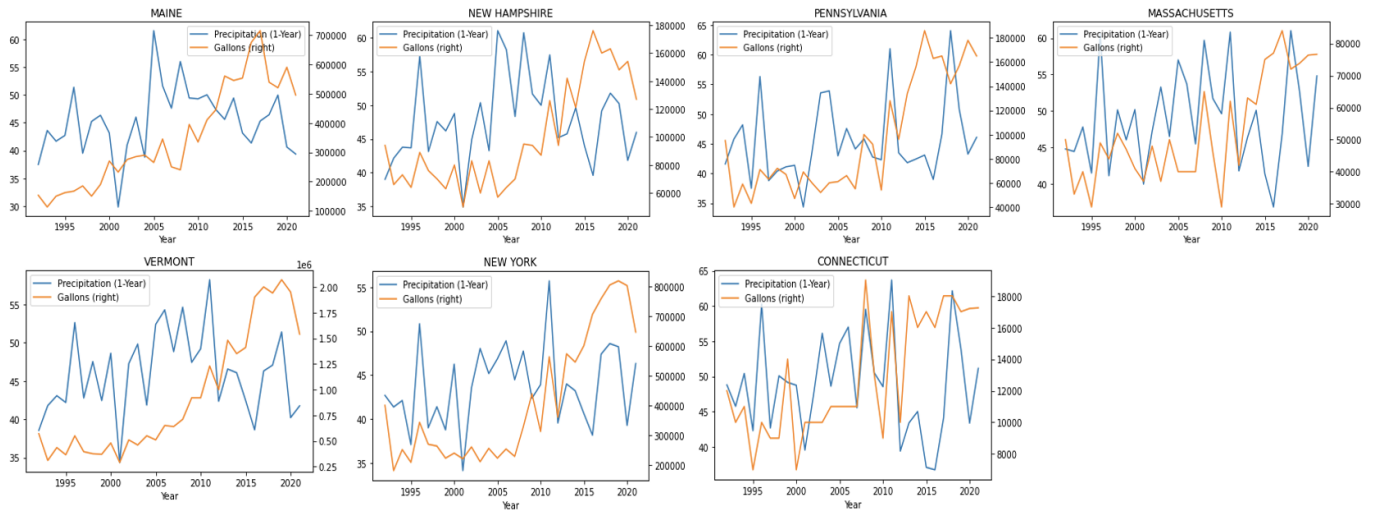
Appendix 6: Regression line of extreme maximum temperature on production in gallons



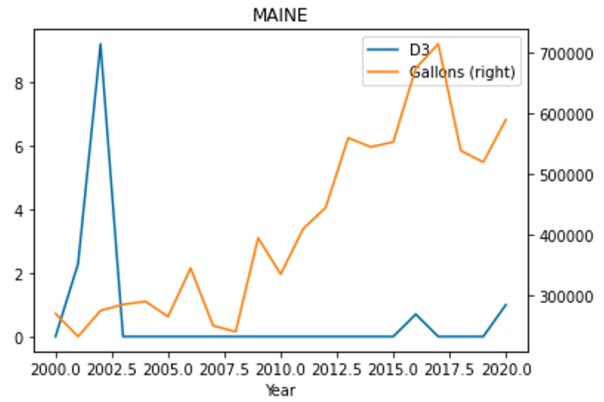
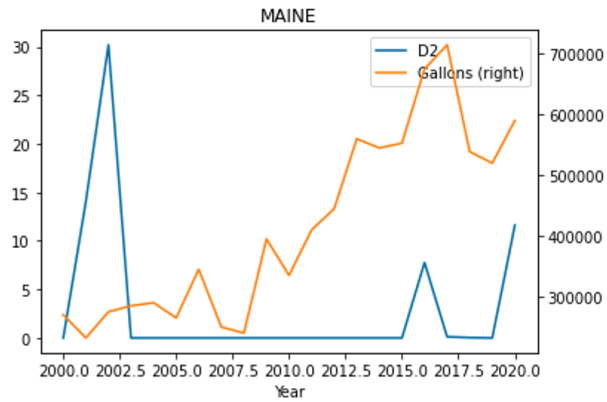
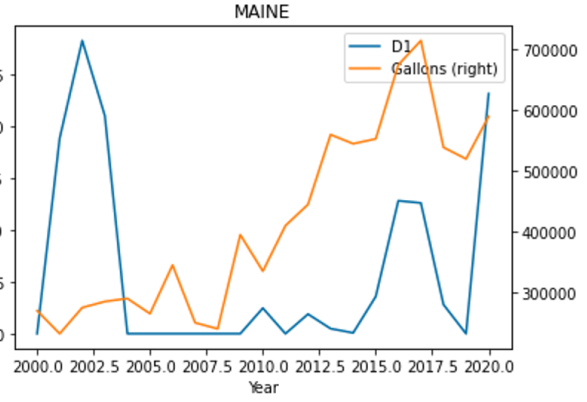
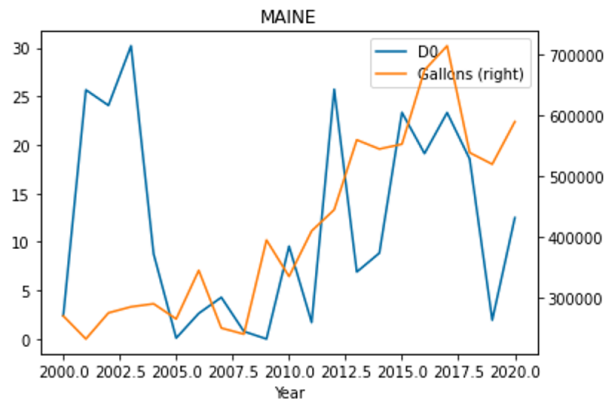
Appendix 7: Regression line of extreme maximum temperature on production in gallons



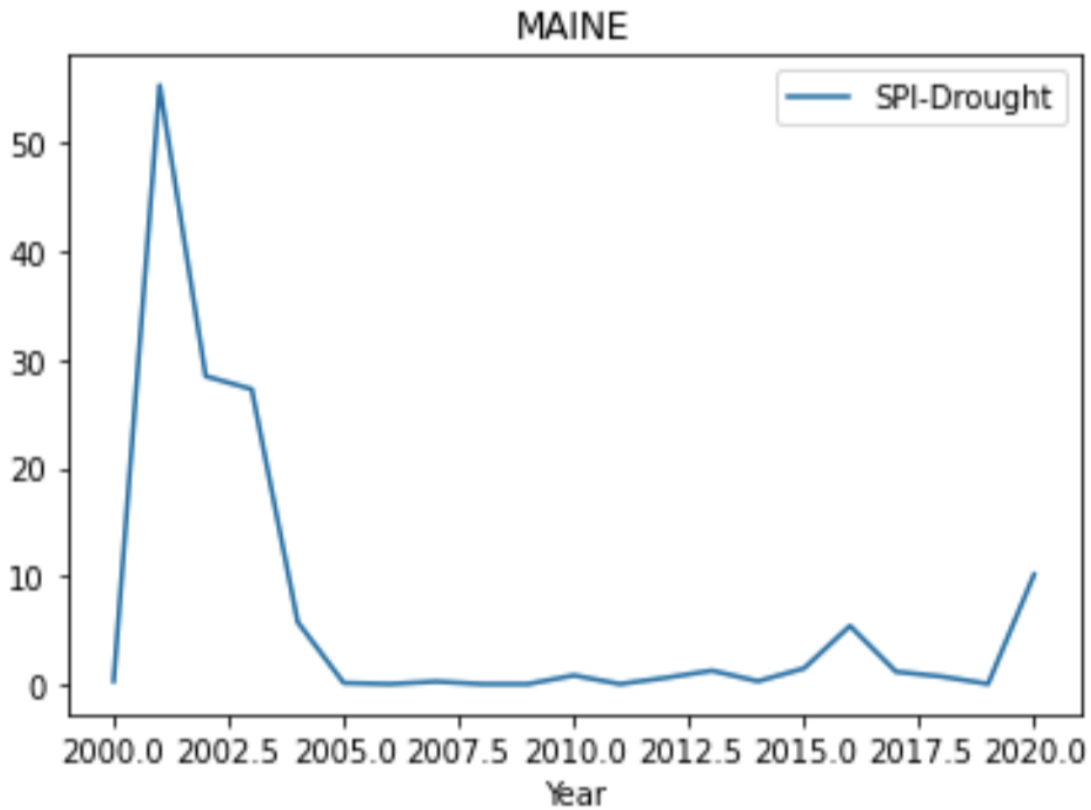
Appendix 8:



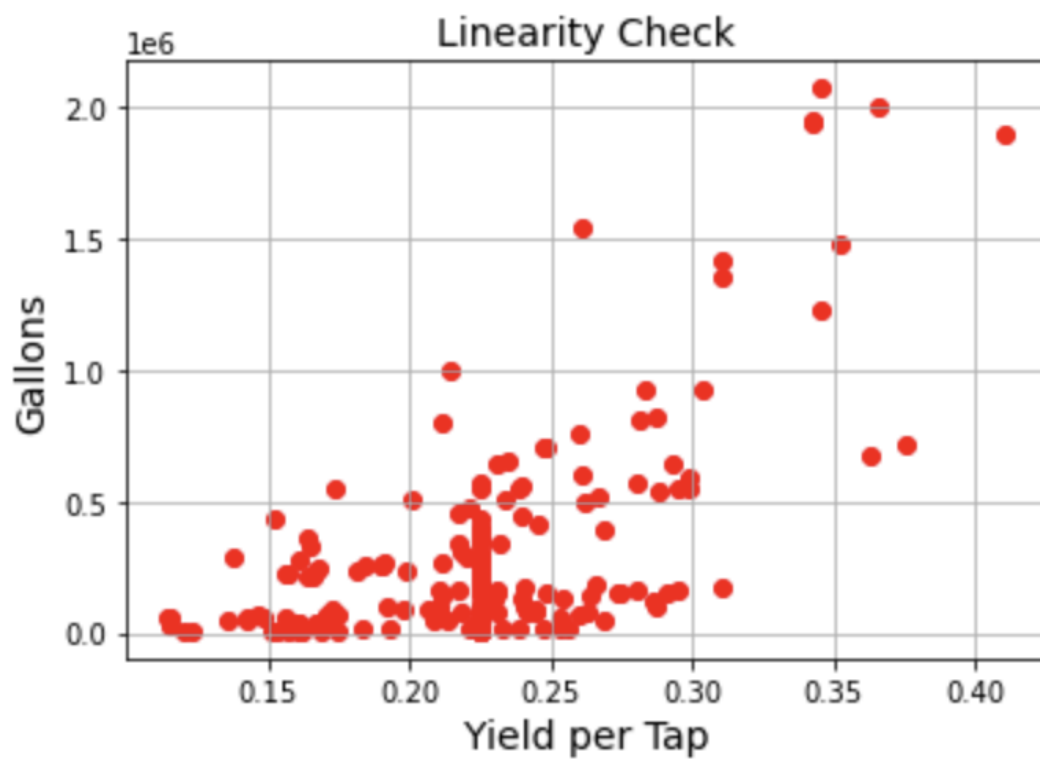
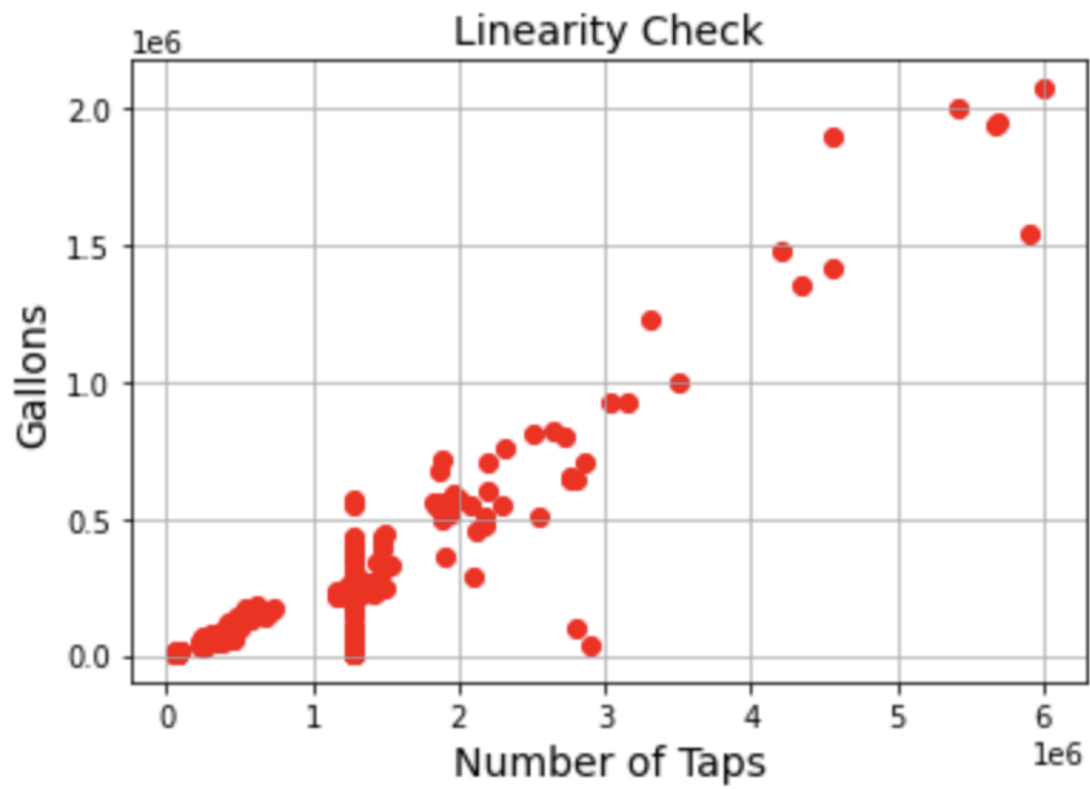
Appendix 9:



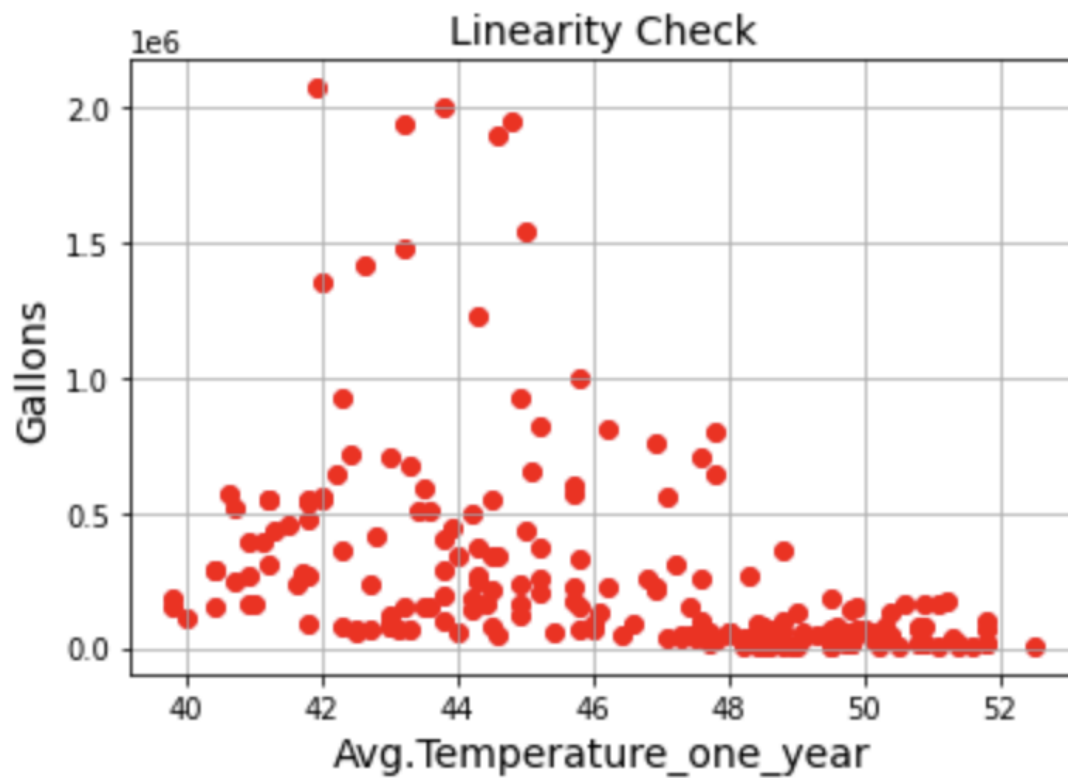
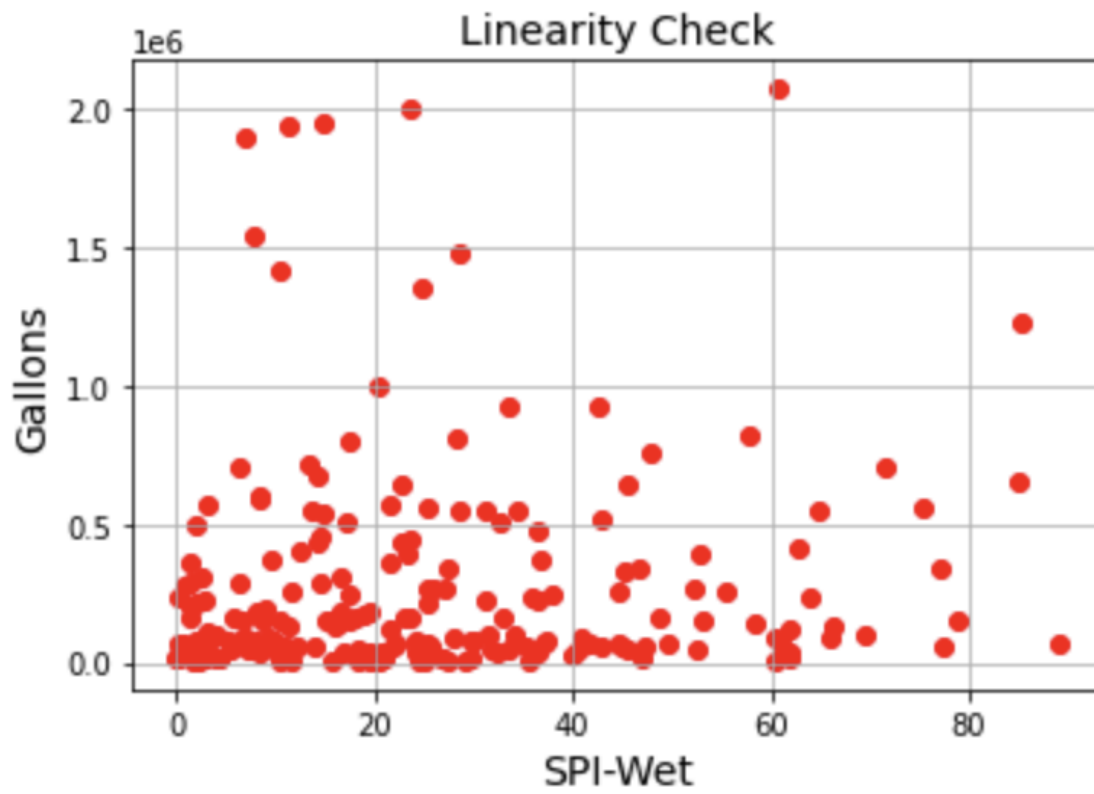
Appendix 10:



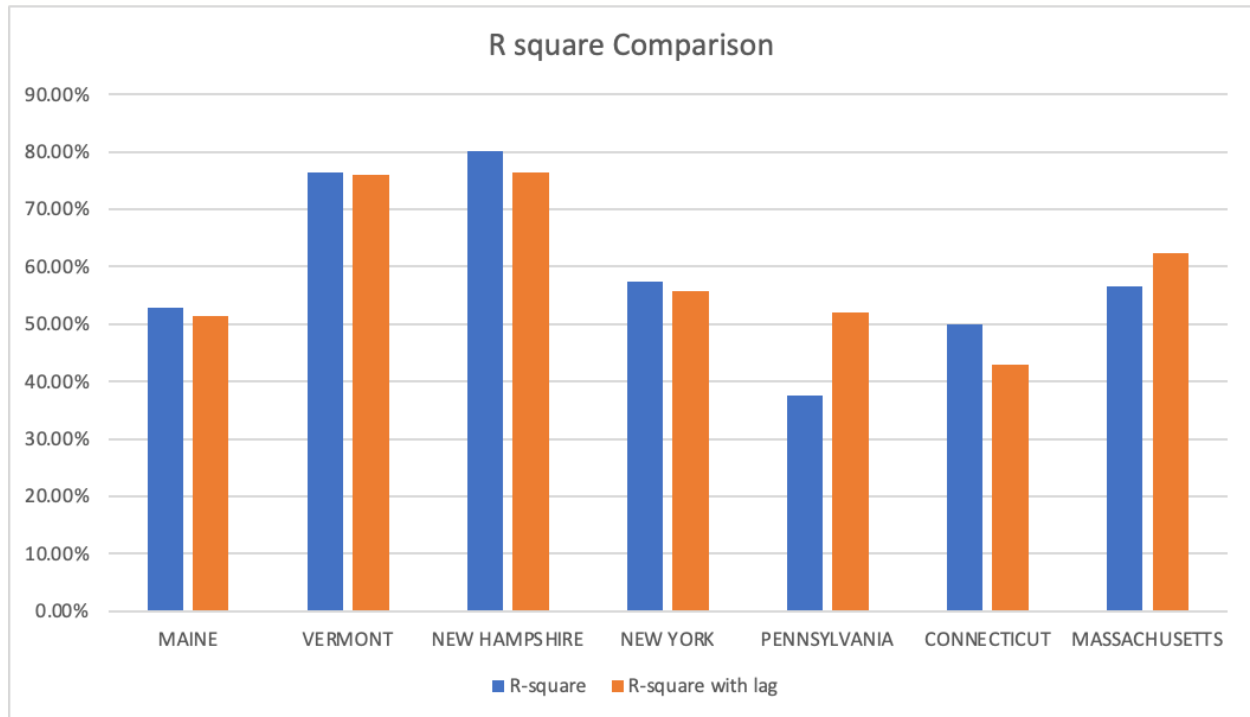
Appendix 11:



Appendix 12:



Appendix 13:



Appendix 14:

OLS Regression Results						
Dep. Variable:	y	R-squared:	0.528			
Model:	OLS	Adj. R-squared:	0.194			
Method:	Least Squares	F-statistic:	1.583			
Date:	Sat, 30 Apr 2022	Prob (F-statistic):	0.188			
Time:	19:29:55	Log-Likelihood:	-392.85			
No. Observations:	30	AIC:	811.7			
Df Residuals:	17	BIC:	829.9			
Df Model:	12					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	6.2e+08	1.27e+10	0.049	0.962	-2.62e+10	2.74e+10
Avg.Temperature_one_year	-3.086e+04	3.69e+04	-0.836	0.415	-1.09e+05	4.7e+04
Precipitation_one_year	-7331.6186	9886.484	-0.742	0.468	-2.82e+04	1.35e+04
none	-6.188e+06	1.27e+08	-0.049	0.962	-2.74e+08	2.62e+08
D0	1.127e+13	1.17e+13	0.964	0.349	-1.34e+13	3.6e+13
D1	2.255e+13	2.34e+13	0.964	0.349	-2.68e+13	7.19e+13
D2	3.382e+13	3.51e+13	0.964	0.349	-4.02e+13	1.08e+14
D3	4.51e+13	4.68e+13	0.964	0.349	-5.36e+13	1.44e+14
DSCI	-1.127e+13	1.17e+13	-0.964	0.349	-3.6e+13	1.34e+13
SPI-Drought	-1.928e+04	7168.590	-2.689	0.016	-3.44e+04	-4152.966
SPI-Wet	5438.2765	3591.113	1.514	0.148	-2138.309	1.3e+04
EMNT(in air)	3190.4185	6839.786	0.466	0.647	-1.12e+04	1.76e+04
EMXT(in air)	6206.0932	2.3e+04	0.270	0.790	-4.23e+04	5.47e+04
Omnibus:	0.431	Durbin-Watson:	0.962			
Prob(Omnibus):	0.806	Jarque-Bera (JB):	0.362			
Skew:	0.243	Prob(JB):	0.835			
Kurtosis:	2.771	Cond. No.	3.29e+11			

Appendix 15:

OLS Regression Results						
=====						
Dep. Variable:	Gallons	R-squared:	0.473			
Model:	OLS	Adj. R-squared:	0.438			
Method:	Least Squares	F-statistic:	13.51			
Date:	Sun, 01 May 2022	Prob (F-statistic):	3.82e-21			
Time:	19:40:58	Log-Likelihood:	-2933.1			
No. Observations:	210	AIC:	5894.			
Df Residuals:	196	BIC:	5941.			
Df Model:	13					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

const	-1.111e+09	5.51e+09	-0.202	0.840	-1.2e+10	9.75e+09
Avg.Temperature_one_year	-7.574e+04	8818.466	-8.588	0.000	-9.31e+04	-5.83e+04
Precipitation_six_month	-1.864e+04	7552.788	-2.468	0.014	-3.35e+04	-3743.322
none	1.113e+07	5.51e+07	0.202	0.840	-9.75e+07	1.2e+08
D0	-6.437e+11	2.82e+12	-0.228	0.820	-6.21e+12	4.92e+12
D1	-1.287e+12	5.65e+12	-0.228	0.820	-1.24e+13	9.85e+12
D2	-9.656e+11	4.24e+12	-0.228	0.820	-9.32e+12	7.39e+12
D2	-9.656e+11	4.24e+12	-0.228	0.820	-9.32e+12	7.39e+12
D3	-2.575e+12	1.13e+13	-0.228	0.820	-2.48e+13	1.97e+13
D4	-3.219e+12	1.41e+13	-0.228	0.820	-3.11e+13	2.46e+13
DSCI	6.437e+11	2.82e+12	0.228	0.820	-4.92e+12	6.21e+12
SPI-Drought	-1.909e+04	3613.510	-5.284	0.000	-2.62e+04	-1.2e+04
SPI-Wet	8488.4794	1562.595	5.432	0.000	5406.821	1.16e+04
EMNT(in air)	3422.8589	2292.318	1.493	0.137	-1097.916	7943.633
EMXT(in air)	2.243e+04	6767.305	3.315	0.001	9088.080	3.58e+04